

Package ‘aifedducation’

August 24, 2025

Type Package

Title Artificial Intelligence for Education

Version 1.1.1

Description In social and educational settings, the use of Artificial Intelligence (AI) is a challenging task. Relevant data is often only available in handwritten forms, or the use of data is restricted by privacy policies. This often leads to small data sets. Furthermore, in the educational and social sciences, data is often unbalanced in terms of frequencies. To support educators as well as educational and social researchers in using the potentials of AI for their work, this package provides a unified interface for neural nets in 'PyTorch' to deal with natural language problems. In addition, the package ships with a shiny app, providing a graphical user interface. This allows the usage of AI for people without skills in writing python/R scripts. The tools integrate existing mathematical and statistical methods for dealing with small data sets via pseudo-labeling (e.g. Cascante-Bonilla et al. (2020) <[doi:10.48550/arXiv.2001.06001](https://doi.org/10.48550/arXiv.2001.06001)>) and imbalanced data via the creation of synthetic cases (e.g. Islam et al. (2012) <[doi:10.1016/j.asoc.2021.108288](https://doi.org/10.1016/j.asoc.2021.108288)>). Performance evaluation of AI is connected to measures from content analysis which educational and social researchers are generally more familiar with (e.g. Berding & Pargmann (2022) <[doi:10.30819/5581](https://doi.org/10.30819/5581)>, Gwet (2014) <ISBN:978-0-9708062-8-4>, Krippendorff (2019) <[doi:10.4135/9781071878781](https://doi.org/10.4135/9781071878781)>). Estimation of energy consumption and CO₂ emissions during model training is done with the 'python' library 'codcarbon'. Finally, all objects created with this package allow to share trained AI models with other people.

License GPL-3

URL <https://fberding.github.io/aifedducation/>

BugReports <https://github.com/cran/aifedducation/issues>

Depends R (>= 3.5.0)

Imports doParallel, foreach, iotarelr(>= 0.1.5), methods, Rcpp (>= 1.0.10), reshape2, reticulate (>= 1.42.0), rlang, stringi, utils

Suggests bslib, DT, fs, future, ggplot2, knitr, pkgdown, promises, readtext, readxl, rmarkdown, shiny(>= 1.9.0), shinyFiles, shinyWidgets, shinyCSSloaders, sortable, testthat (>= 3.0.0)

LinkingTo Rcpp, RcppArmadillo

VignetteBuilder knitr

Config/testthat.edition 3

Encoding UTF-8

LazyData true

RoxygenNote 7.3.2

SystemRequirements PyTorch (see vignette ``Get started'')

Config/Needs/website rmarkdown

NeedsCompilation yes

Author Berding Florian [aut, cre] (ORCID: <<https://orcid.org/0000-0002-3593-1695>>), Tykhonova Yuliia [aut] (ORCID: <<https://orcid.org/0009-0006-9015-1006>>), Pargmann Julia [ctb] (ORCID: <<https://orcid.org/0000-0003-3616-0172>>), Leube Anna [ctb] (ORCID: <<https://orcid.org/0009-0001-6949-1608>>), Riebenbauer Elisabeth [ctb] (ORCID: <<https://orcid.org/0000-0002-8535-3694>>), Rebmann Karin [ctb], Slopinski Andreas [ctb]

Maintainer Berding Florian <florian.berding@uni-hamburg.de>

Repository CRAN

Date/Publication 2025-08-23 22:00:26 UTC

Contents

.AIFEBaseTransformer	5
.AIFEBertTransformer	14
.AIFEFunnelTransformer	19
.AIFELongformerTransformer	24
.AIFEModernBertTransformer	29
.AIFEMpnetTransformer	33
.AIFERobertaTransformer	38
add_missing_args	43
AIFE BaseModel	44
AIFETrType	49
aife_transformer.load_model_mlm	49
aife_transformer.load_tokenizer	50
aife_transformer.make	51
auto_n_cores	51
build_documentation_for_model	52
build_layer_stack_documentation_for_vignette	53
calc_standard_classification_measures	53

calc_tokenizer_statistics	54
cat_message	55
check_adjust_n_samples_on_CI	55
check_aif_py_modules	56
check_all_args	57
check_class_and_type	57
ClassifiersBasedOnTextEmbeddings	58
class_vector_to_py_dataset	63
clean_pytorch_log_transformers	63
cohens_kappa	64
create_dir	65
create_object	65
create_synthetic_units_from_matrix	66
data.frame_to_py_dataset	67
DataManagerClassifier	67
EmbeddedText	73
fleiss_kappa	78
generate_args_for_tests	79
generate_embeddings	80
generate_id	80
generate_tensors	81
get_alpha_3_codes	82
get_batches_index	82
get_called_args	83
get_coder_metrics	83
get_current_args_for_print	85
get_depr_obj_names	85
get_desc_for_core_model_architecture	86
get_file_extension	86
get_fixed_test_tensor	87
get_layer_documentation	87
get_magnitude_values	88
get_n_chunks	89
get_parameter_documentation	90
get_param_def	90
get_param_dict	91
get_param_doc_desc	92
get_py_package_version	93
get_py_package_versions	93
get_synthetic_cases_from_matrix	94
get_TEClassifiers_class_names	95
get_test_data_for_classifiers	95
gwet_ac	96
imdb_movie_reviews	97
install_aifedducation	97
install_aifedducation_studio	98
install_py_modules	99
kendalls_w	100

knnor	101
knnor_is_same_class	102
kripp_alpha	102
LargeDataSetBase	103
LargeDataSetForText	106
LargeDataSetForTextEmbeddings	110
load_all_py_scripts	116
load_from_disk	116
load_py_scripts	117
long_load_target_data	117
matrix_to_array_c	118
ModelsBasedOnTextEmbeddings	119
output_message	121
prepare_r_array_for_dataset	122
prepare_session	122
print_message	123
py_dataset_to_embeddings	123
random_bool_on_CI	124
read_log	125
read_loss_log	125
reduce_to_unique	126
reset_log	127
reset_loss_log	127
run_py_file	128
save_to_disk	128
set_transformers_logger	129
start_aifedducation_studio	130
summarize_args_for_long_task	130
TEClassifierParallel	131
TEClassifierParallelPrototype	137
TEClassifierProtoNet	144
TEClassifierRegular	148
TEClassifiersBasedOnProtoNet	150
TEClassifiersBasedOnRegular	156
TEClassifierSequential	159
TEClassifierSequentialPrototype	165
TEFeatureExtractor	171
tensor_list_to_numpy	175
tensor_to_matrix_c	176
tensor_to_numpy	176
TextEmbeddingModel	177
to_categorical_c	186
update_aifedducation	186
write_log	187

.AIFEBaseTransformer *Base R6 class for creation and definition of .AIFE*Transformer-like classes*

Description

This base class is used to create and define .AIFE*Transformer-like classes. It serves as a skeleton for a future concrete transformer and cannot be used to create an object of itself (an attempt to call new-method will produce an error).

See p.1 Base Transformer Class in [Transformers for Developers](#) for details.

Create

The `create`-method is a basic algorithm that is used to create a new transformer, but cannot be called directly.

Train

The `train`-method is a basic algorithm that is used to train and tune the transformer but cannot be called directly.

Concrete transformer implementation

There are already implemented concrete (child) transformers (e.g. BERT, DeBERTa-V2, etc.), to implement a new one see p.4 Implement A Custom Transformer in [Transformers for Developers](#)

Public fields

`params` A list containing transformer's parameters ('static', 'dynamic' and 'dependent' parameters)

`list()` containing all the transformer parameters. Can be set with `set_model_param()`.

'Static' parameters:

Regardless of the transformer, the following parameters are always included:

- `text_dataset`
- `sustain_track`
- `sustain_iso_code`
- `sustain_region`
- `sustain_interval`
- `trace`
- `pytorch_safetensors`
- `log_dir`
- `log_write_interval`

'Dynamic' parameters:

In the case of `create` it also contains (see `create`-method for details):

- `model_dir`
- `vocab_size`

- max_position_embeddings
- hidden_size
- hidden_act
- hidden_dropout_prob
- attention_probs_dropout_prob
- intermediate_size
- num_attention_heads

In the case of **train** it also contains (see train-method for details):

- output_dir
- model_dir_path
- p_mask
- whole_word
- val_size
- n_epoch
- batch_size
- chunk_size
- min_seq_len
- full_sequences_only
- learning_rate
- n_workers
- multi_process
- keras_trace
- pytorch_trace

'Dependent' parameters:

Depending on the transformer and the method used class may contain different parameters:

- vocab_do_lower_case
- num_hidden_layer
- add_prefix_space
- etc.

temp A list containing temporary transformer's parameters

list() containing all the temporary local variables that need to be accessed between the step functions. Can be set with **set_model_temp()**.

For example, it can be a variable **tok_new** that stores the tokenizer from **steps_for_creation\$create_tokenizer_dr**. To train the tokenizer, access the variable **tok_new** in **steps_for_creation\$calculate_vocab** through the **temp** list of this class.

Methods

Public methods:

- [.AIFEBaseTransformer\\$new\(\)](#)
- [.AIFEBaseTransformer\\$init_transformer\(\)](#)
- [.AIFEBaseTransformer\\$set_title\(\)](#)
- [.AIFEBaseTransformer\\$set_model_param\(\)](#)

- .AIFEBaseTransformer\$set_model_temp()
- .AIFEBaseTransformer\$set_SFC_check_max_pos_emb()
- .AIFEBaseTransformer\$set_SFC_create_tokenizer_draft()
- .AIFEBaseTransformer\$set_SFC_calculate_vocab()
- .AIFEBaseTransformer\$set_SFC_save_tokenizer_draft()
- .AIFEBaseTransformer\$set_SFC_create_final_tokenizer()
- .AIFEBaseTransformer\$set_SFC_create_transformer_model()
- .AIFEBaseTransformer\$set_required_SFC()
- .AIFEBaseTransformer\$set_SFT_load_existing_model()
- .AIFEBaseTransformer\$set_SFT_cuda_empty_cache()
- .AIFEBaseTransformer\$set_SFT_create_data_collator()
- .AIFEBaseTransformer\$create()
- .AIFEBaseTransformer\$train()
- .AIFEBaseTransformer\$clone()

Method new(): An object of this class cannot be created. Thus, method's call will produce an error.

Usage:

```
.AIFEBaseTransformer$new()
```

Returns: This method returns an error.

Method init_transformer(): Method to execute while initializing a new transformer.

Usage:

```
.AIFEBaseTransformer$init_transformer(title, init_trace)
```

Arguments:

title string A new title.

init_trace bool option to show prints. If TRUE (by default) - messages will be shown, otherwise (FALSE) - hidden.

Returns: This method returns nothing.

Method set_title(): Setter for the title. Sets a new value for the **title** private attribute.

Usage:

```
.AIFEBaseTransformer$set_title(title)
```

Arguments:

title string A new title.

Returns: This method returns nothing.

Method set_model_param(): Setter for the parameters. Adds a new parameter and its value to the **params** list.

Usage:

```
.AIFEBaseTransformer$set_model_param(param_name, param_value)
```

Arguments:

`param_name` string Parameter's name.

`param_value` any Parameter's value.

Returns: This method returns nothing.

Method `set_model_temp()`: Setter for the temporary model's parameters. Adds a new temporary parameter and its value to the `temp` list.

Usage:

```
.AIFEBaseTransformer$set_model_temp(temp_name, temp_value)
```

Arguments:

`temp_name` string Parameter's name.

`temp_value` any Parameter's value.

Returns: This method returns nothing.

Method `set_SFC_check_max_pos_emb()`: Setter for the `check_max_pos_emb` element of the private `steps_for_creation` list. Sets a new fun function as the `check_max_pos_emb` step.

Usage:

```
.AIFEBaseTransformer$set_SFC_check_max_pos_emb(fun)
```

Arguments:

`fun` function() A new function.

Returns: This method returns nothing.

Method `set_SFC_create_tokenizer_draft()`: Setter for the `create_tokenizer_draft` element of the private `steps_for_creation` list. Sets a new fun function as the `create_tokenizer_draft` step.

Usage:

```
.AIFEBaseTransformer$set_SFC_create_tokenizer_draft(fun)
```

Arguments:

`fun` function() A new function.

Returns: This method returns nothing.

Method `set_SFC_calculate_vocab()`: Setter for the `calculate_vocab` element of the private `steps_for_creation` list. Sets a new fun function as the `calculate_vocab` step.

Usage:

```
.AIFEBaseTransformer$set_SFC_calculate_vocab(fun)
```

Arguments:

`fun` function() A new function.

Returns: This method returns nothing.

Method `set_SFC_save_tokenizer_draft()`: Setter for the `save_tokenizer_draft` element of the private `steps_for_creation` list. Sets a new fun function as the `save_tokenizer_draft` step.

Usage:

```
.AIFEBaseTransformer$set_SFC_save_tokenizer_draft(fun)
```

Arguments:

fun function() A new function.

Returns: This method returns nothing.

Method set_SFC_create_final_tokenizer(): Setter for the create_final_tokenizer element of the private steps_for_creation list. Sets a new fun function as the create_final_tokenizer step.

Usage:

.AIFEBaseTransformer\$set_SFC_create_final_tokenizer(fun)

Arguments:

fun function() A new function.

Returns: This method returns nothing.

Method set_SFC_create_transformer_model(): Setter for the create_transformer_model element of the private steps_for_creation list. Sets a new fun function as the create_transformer_model step.

Usage:

.AIFEBaseTransformer\$set_SFC_create_transformer_model(fun)

Arguments:

fun function() A new function.

Returns: This method returns nothing.

Method set_required_SFC(): Setter for all required elements of the private steps_for_creation list. Executes setters for all required creation steps.

Usage:

.AIFEBaseTransformer\$set_required_SFC(required_SFC)

Arguments:

required_SFC list() A list of all new required steps.

Returns: This method returns nothing.

Method set_SFT_load_existing_model(): Setter for the load_existing_model element of the private steps_for_training list. Sets a new fun function as the load_existing_model step.

Usage:

.AIFEBaseTransformer\$set_SFT_load_existing_model(fun)

Arguments:

fun function() A new function.

Returns: This method returns nothing.

Method set_SFT_cuda_empty_cache(): Setter for the cuda_empty_cache element of the private steps_for_training list. Sets a new fun function as the cuda_empty_cache step.

Usage:

.AIFEBaseTransformer\$set_SFT_cuda_empty_cache(fun)

Arguments:

`fun function()` A new function.

Returns: This method returns nothing.

Method `set_SFT_create_data_collator()`: Setter for the `create_data_collator` element of the private `steps_for_training` list. Sets a new `fun` function as the `create_data_collator` step. Use this method to make a custom data collator for a transformer.

Usage:

```
.AIFEBaseTransformer$set_SFT_create_data_collator(fun)
```

Arguments:

`fun function()` A new function.

Returns: This method returns nothing.

Method `create()`: This method creates a transformer configuration based on the child-transformer architecture and a vocabulary using the python libraries `transformers` and `tokenizers`.

This method **adds** the following parameters to the `temp` list:

- `log_file`
- `raw_text_dataset`
- `pt_safe_save`
- `value_top`
- `total_top`
- `message_top`

This method **uses** the following parameters from the `temp` list:

- `log_file`
- `raw_text_dataset`
- `tokenizer`

Usage:

```
.AIFEBaseTransformer$create(
  model_dir,
  text_dataset,
  vocab_size,
  max_position_embeddings,
  hidden_size,
  num_attention_heads,
  intermediate_size,
  hidden_act,
  hidden_dropout_prob,
  attention_probs_dropout_prob,
  sustain_track,
  sustain_iso_code,
  sustain_region,
  sustain_interval,
  trace,
  pytorch_safetensors,
```

```

        log_dir,
        log_write_interval
    )
Arguments:
model_dir string Path to the directory where the model should be saved. Allowed values:
any
text_dataset LargeDataSetForText LargeDataSetForText Object storing textual data.
vocab_size int Size of the vocabulary. Allowed values:  $1000 \leq x \leq 5e+05$ 
max_position_embeddings int Number of maximum position embeddings. This parameter
also determines the maximum length of a sequence which can be processed with the model.
Allowed values:  $10 \leq x \leq 4048$ 
hidden_size int Number of neurons in each layer. This parameter determines the dimension-
ality of the resulting text embedding. Allowed values:  $1 \leq x \leq 2048$ 
num_attention_heads int determining the number of attention heads for a self-attention layer.
Only relevant if attention_type='multihead' Allowed values:  $0 \leq x$ 
intermediate_size int determining the size of the projection layer within a each transformer
encoder. Allowed values:  $1 \leq x$ 
hidden_act string Name of the activation function. Allowed values: 'gelu', 'relu', 'silu',
'gelu_new'
hidden_dropout_prob double Ratio of dropout. Allowed values:  $0 \leq x \leq 0.6$ 
attention_probs_dropout_prob double Ratio of dropout for attention probabilities. Al-
lowed values:  $0 \leq x \leq 0.6$ 
sustain_track bool If TRUE energy consumption is tracked during training via the python
library 'codecarbon'.
sustain_iso_code string ISO code (Alpha-3-Code) for the country. This variable must be
set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List\_of\_ISO\_3166\_country\_codes. Allowed values: any
sustain_region string Region within a country. Only available for USA and Canada See the
documentation of codecarbon for more information. https://mlco2.github.io/codecarbon/parameters.html Allowed values: any
sustain_interval int Interval in seconds for measuring power usage. Allowed values:  $1 \leq
x$ 
trace bool TRUE if information about the estimation phase should be printed to the console.
pytorch_safetensors bool


- TRUE: a 'pytorch' model is saved in safetensors format.
- FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format
(.bin).


log_dir string Path to the directory where the log files should be saved. If no logging is
desired set this argument to NULL. Allowed values: any
log_write_interval 'r get_param_doc_desclog_write_interval
Returns: This method does not return an object. Instead, it saves the configuration and vocab-
ulary of the new model to disk.

```

Method train(): This method can be used to train or fine-tune a transformer based on BERT architecture with the help of the python libraries `transformers`, `datasets`, and `tokenizers`.
This method **adds** the following parameters to the temp list:

- log_file
- loss_file
- from_pt
- from_tf
- load_safe
- raw_text_dataset
- pt_safe_save
- value_top
- total_top
- message_top

This method **uses** the following parameters from the `temp` list:

- log_file
- raw_text_dataset
- tokenized_dataset
- tokenizer

Usage:

```
.AIFEBaseTransformer$train(
    output_dir,
    model_dir_path,
    text_dataset,
    p_mask,
    whole_word,
    val_size,
    n_epoch,
    batch_size,
    chunk_size,
    full_sequences_only,
    min_seq_len,
    learning_rate,
    sustain_track,
    sustain_iso_code,
    sustain_region,
    sustain_interval,
    trace,
    pytorch_trace,
    pytorch_safetensors,
    log_dir,
    log_write_interval
)
```

Arguments:

`output_dir` string Path to the directory where the model should be saved. Allowed values:
any
`model_dir_path` string Path to the directory where the original model is stored. Allowed
values: any
`text_dataset` LargeDataSetForText [LargeDataSetForText](#) Object storing textual data.

`p_mask` double Ratio that determines the number of words/tokens used for masking. Allowed values: $0 < x < 1$

`whole_word` bool * TRUE: whole word masking should be applied.

- FALSE: token masking is used.

`val_size` double between 0 and 1, indicating the proportion of cases which should be used for the validation sample during the estimation of the model. The remaining cases are part of the training data. Allowed values: $0 < x < 1$

`n_epoch` int Number of training epochs. Allowed values: $1 \leq x$

`batch_size` int Size of the batches for training. Allowed values: $1 \leq x$

`chunk_size` int Maximum length of every sequence. Must be equal or less the global maximum size allowed by the model. Allowed values: $100 \leq x$

`full_sequences_only` bool TRUE for using only chunks with a sequence length equal to `chunk_size`.

`min_seq_len` int Only relevant if `full_sequences_only` = FALSE. Value determines the minimal sequence length included in training process. Allowed values: $10 \leq x$

`learning_rate` double Initial learning rate for the training. Allowed values: $0 < x \leq 1$

`sustain_track` bool If TRUE energy consumption is tracked during training via the python library 'codecarbon'.

`sustain_iso_code` string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes. Allowed values: any

`sustain_region` string Region within a country. Only available for USA and Canada See the documentation of codecarbon for more information. <https://mlco2.github.io/codecarbon/parameters.html> Allowed values: any

`sustain_interval` int Interval in seconds for measuring power usage. Allowed values: $1 \leq x$

`trace` bool TRUE if information about the estimation phase should be printed to the console.

`pytorch_trace` “r get_param_doc_desc("pytorch_trace")“

`pytorch_safetensors` bool

- TRUE: a 'pytorch' model is saved in safetensors format.
- FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format (.bin).

`log_dir` string Path to the directory where the log files should be saved. If no logging is desired set this argument to NULL. Allowed values: any

`log_write_interval` ‘r get_param_doc_desclog_write_interval

Returns: This method does not return an object. Instead, it saves the configuration and vocabulary of the new model to disk.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

`.AIFEBaseTransformer$clone(deep = FALSE)`

Arguments:

`deep` Whether to make a deep clone.

References

Hugging Face transformers documentation:

- [BERT](#)
- [DeBERTa](#)
- [Funnel](#)
- [Longformer](#)
- [RoBERTa](#)
- [MPNet](#)

See Also

Other R6 classes for transformers: [.AIFEBertTransformer](#), [.AIFEFunnelTransformer](#), [.AIFELongformerTransformer](#), [.AIFEMpnetTransformer](#), [.AIFERobertaTransformer](#)

[.AIFEBertTransformer](#) *Child R6 class for creation and training of BERT transformers*

Description

This class has the following methods:

- `create`: creates a new transformer based on BERT.
- `train`: trains and fine-tunes a BERT model.

Create

New models can be created using the `.AIFEBertTransformer$create` method.

Train

To train the model, pass the directory of the model to the method `.AIFEBertTransformer$train`.

Pre-Trained models that can be fine-tuned using this method are available at <https://huggingface.co/>.

The model is trained using dynamic masking, as opposed to the original paper, which used static masking.

Super class

`aifedducation::AIFEBaseTransformer -> .AIFEBertTransformer`

Methods

Public methods:

- .AIFEBertTransformer\$new()
- .AIFEBertTransformer\$create()
- .AIFEBertTransformer\$train()
- .AIFEBertTransformer\$clone()

Method new(): Creates a new transformer based on BERT and sets the title.

Usage:

```
.AIFEBertTransformer$new(init_trace = TRUE)
```

Arguments:

init_trace bool option to show prints. If TRUE (by default) - messages will be shown, otherwise (FALSE) - hidden.

Returns: This method returns nothing.

Method create(): This method creates a transformer configuration based on the BERT base architecture and a vocabulary based on WordPiece by using the python libraries transformers and tokenizers.

This method adds the following '*dependent*' parameters to the base class's inherited params list:

- vocab_do_lower_case
- num_hidden_layer

Usage:

```
.AIFEBertTransformer$create(
  model_dir,
  text_dataset,
  vocab_size = 30522,
  vocab_do_lower_case = FALSE,
  max_position_embeddings = 512,
  hidden_size = 768,
  num_hidden_layer = 12,
  num_attention_heads = 12,
  intermediate_size = 3072,
  hidden_act = "GELU",
  hidden_dropout_prob = 0.1,
  attention_probs_dropout_prob = 0.1,
  sustain_track = FALSE,
  sustain_iso_code = NULL,
  sustain_region = NULL,
  sustain_interval = 15,
  trace = TRUE,
  pytorch_safetensors = TRUE,
  log_dir = NULL,
  log_write_interval = 2
)
```

Arguments:

```

model_dir string Path to the directory where the model should be saved. Allowed values:
any

text_dataset LargeDataSetForText LargeDataSetForText Object storing textual data.

vocab_size int Size of the vocabulary. Allowed values: 1000 <= x <= 5e+05

vocab_do_lower_case bool TRUE if all words/tokens should be lower case.

max_position_embeddings int Number of maximum position embeddings. This parameter
also determines the maximum length of a sequence which can be processed with the model.
Allowed values: 10 <= x <= 4048

hidden_size int Number of neurons in each layer. This parameter determines the dimension-
ality of the resulting text embedding. Allowed values: 1 <= x <= 2048

num_hidden_layer int Number of hidden layers. Allowed values: 1 <= x

num_attention_heads int determining the number of attention heads for a self-attention layer.
Only relevant if attention_type='multihead' Allowed values: 0 <= x

intermediate_size int determining the size of the projection layer within a each transformer
encoder. Allowed values: 1 <= x

hidden_act string Name of the activation function. Allowed values: 'gelu', 'relu', 'silu',
'gelu_new'

hidden_dropout_prob double Ratio of dropout. Allowed values: 0 <= x <= 0.6

attention_probs_dropout_prob double Ratio of dropout for attention probabilities. Al-
lowed values: 0 <= x <= 0.6

sustain_track bool If TRUE energy consumption is tracked during training via the python
library 'codecarbon'.

sustain_iso_code string ISO code (Alpha-3-Code) for the country. This variable must be
set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List\_of\_ISO\_3166\_country\_codes. Allowed values: any

sustain_region string Region within a country. Only available for USA and Canada See the
documentation of codecarbon for more information. https://mlco2.github.io/codecarbon/parameters.html Allowed values: any

sustain_interval int Interval in seconds for measuring power usage. Allowed values: 1 <=
x

trace bool TRUE if information about the estimation phase should be printed to the console.

pytorch_safetensors bool
• TRUE: a 'pytorch' model is saved in safetensors format.
• FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format
(.bin).

log_dir string Path to the directory where the log files should be saved. If no logging is
desired set this argument to NULL. Allowed values: any

log_write_interval int Time in seconds determining the interval in which the logger should
try to update the log files. Only relevant if log_dir is not NULL. Allowed values: 1 <= x

```

Returns: This method does not return an object. Instead, it saves the configuration and vocabulary of the new model to disk.

Method train(): This method can be used to train or fine-tune a transformer based on BERT architecture with the help of the python libraries transformers, datasets, and tokenizers.

Usage:

```
.AIFEBertTransformer$train(
  output_dir,
  model_dir_path,
  text_dataset,
  p_mask = 0.15,
  whole_word = TRUE,
  val_size = 0.1,
  n_epoch = 1,
  batch_size = 12,
  chunk_size = 250,
  full_sequences_only = FALSE,
  min_seq_len = 50,
  learning_rate = 0.003,
  sustain_track = FALSE,
  sustain_iso_code = NULL,
  sustain_region = NULL,
  sustain_interval = 15,
  trace = TRUE,
  pytorch_trace = 1,
  pytorch_safetensors = TRUE,
  log_dir = NULL,
  log_write_interval = 2
)
```

Arguments:

`output_dir` string Path to the directory where the model should be saved. Allowed values:
any

`model_dir_path` string Path to the directory where the original model is stored. Allowed
values: any

`text_dataset` LargeDataSetForText [LargeDataSetForText](#) Object storing textual data.

`p_mask` double Ratio that determines the number of words/tokens used for masking. Allowed
values: $0 < x < 1$

`whole_word` bool * TRUE: whole word masking should be applied.

- FALSE: token masking is used.

`val_size` double between 0 and 1, indicating the proportion of cases which should be used for
the validation sample during the estimation of the model. The remaining cases are part of
the training data. Allowed values: $0 < x < 1$

`n_epoch` int Number of training epochs. Allowed values: $1 \leq x$

`batch_size` int Size of the batches for training. Allowed values: $1 \leq x$

`chunk_size` int Maximum length of every sequence. Must be equal or less the global maxi-
mum size allowed by the model. Allowed values: $100 \leq x$

`full_sequences_only` bool TRUE for using only chunks with a sequence length equal to
chunk_size.

`min_seq_len` int Only relevant if `full_sequences_only = FALSE`. Value determines the min-
imal sequence length included in training process. Allowed values: $10 \leq x$

`learning_rate` double Initial learning rate for the training. Allowed values: $0 < x \leq 1$

`sustain_track` bool If TRUE energy consumption is tracked during training via the python
library 'codecarbon'.

sustain_iso_code string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes. Allowed values: any

sustain_region string Region within a country. Only available for USA and Canada See the documentation of codecarbon for more information. <https://mlco2.github.io/codecarbon/parameters.html> Allowed values: any

sustain_interval int Interval in seconds for measuring power usage. Allowed values: $1 \leq x$

trace bool TRUE if information about the estimation phase should be printed to the console.

pytorch_trace int ml_trace=0 does not print any information about the training process from pytorch on the console. Allowed values: $0 \leq x \leq 1$

pytorch_safetensors bool

- TRUE: a 'pytorch' model is saved in safetensors format.
- FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format (.bin).

log_dir string Path to the directory where the log files should be saved. If no logging is desired set this argument to NULL. Allowed values: any

log_write_interval int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if log_dir is not NULL. Allowed values: $1 \leq x$

Returns: This method does not return an object. Instead the trained or fine-tuned model is saved to disk.

Method clone(): The objects of this class are cloneable with this method.

Usage:

`.AIFEBertTransformer$clone(deep = FALSE)`

Arguments:

deep Whether to make a deep clone.

Note

This model uses a WordPiece tokenizer like BERT and can be trained with whole word masking. The transformer library may display a warning, which can be ignored.

References

Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In J. Burstein, C. Doran, & T. Solorio (Eds.), Proceedings of the 2019 Conference of the North (pp. 4171–4186). Association for Computational Linguistics. doi:[10.18653/v1/N191423](https://doi.org/10.18653/v1/N191423)

Hugging Face documentation

- https://huggingface.co/docs/transformers/model_doc/bert
- https://huggingface.co/docs/transformers/model_doc/bert#transformers.BertForMaskedLM
- https://huggingface.co/docs/transformers/model_doc/bert#transformers.TFBertForMaskedLM

See Also

Other R6 classes for transformers: [.AIFEBaseTransformer](#), [.AIFEFunnelTransformer](#), [.AIFELongformerTransformer](#), [.AIFEMpnetTransformer](#), [.AIFERobertaTransformer](#)

.AIFEFunnelTransformer

Child R6 class for creation and training of Funnel transformers

Description

This class has the following methods:

- `create`: creates a new transformer based on Funnel.
- `train`: trains and fine-tunes a Funnel model.

Create

New models can be created using the `.AIFEFunnelTransformer$create` method.

Model is created with `separate_cls = TRUE`, `truncate_seq = TRUE`, and `pool_q_only = TRUE`.

Train

To train the model, pass the directory of the model to the method `.AIFEFunnelTransformer$train`.

Pre-Trained models which can be fine-tuned with this function are available at <https://huggingface.co/>.

Training of the model makes use of dynamic masking.

Super class

`aifedducation::AIFEBaseTransformer -> AIFEFunnelTransformer`

Methods

Public methods:

- `.AIFEFunnelTransformer$new()`
- `.AIFEFunnelTransformer$create()`
- `.AIFEFunnelTransformer$train()`
- `.AIFEFunnelTransformer$clone()`

Method `new()`: Creates a new transformer based on Funnel and sets the title.

Usage:

`.AIFEFunnelTransformer$new(init_trace = TRUE)`

Arguments:

`init_trace` bool option to show prints. If TRUE (by default) - messages will be shown, otherwise (FALSE) - hidden.

Returns: This method returns nothing.

Method create(): This method creates a transformer configuration based on the Funnel transformer base architecture and a vocabulary based on WordPiece using the python `transformers` and `tokenizers` libraries.

This method adds the following '*dependent*' parameters to the base class's inherited `params` list:

- `vocab_do_lower_case`
- `target_hidden_size`
- `block_sizes`
- `num_decoder_layers`
- `pooling_type`
- `activation_dropout`

Usage:

```
.AIFEFunnelTransformer$create(
  model_dir,
  text_dataset,
  vocab_size = 30522,
  vocab_do_lower_case = FALSE,
  max_position_embeddings = 512,
  hidden_size = 768,
  target_hidden_size = 64,
  block_sizes = c(4, 4, 4),
  num_attention_heads = 12,
  intermediate_size = 3072,
  num_decoder_layers = 2,
  pooling_type = "Mean",
  hidden_act = "GELU",
  hidden_dropout_prob = 0.1,
  attention_probs_dropout_prob = 0.1,
  activation_dropout = 0,
  sustain_track = TRUE,
  sustain_iso_code = NULL,
  sustain_region = NULL,
  sustain_interval = 15,
  trace = TRUE,
  pytorch_safetensors = TRUE,
  log_dir = NULL,
  log_write_interval = 2
)
```

Arguments:

`model_dir` string Path to the directory where the model should be saved. Allowed values:
any

`text_dataset` LargeDataSetForText [LargeDataSetForText](#) Object storing textual data.

`vocab_size` int Size of the vocabulary. Allowed values: $1000 \leq x \leq 5e+05$

`vocab_do_lower_case` bool TRUE if all words/tokens should be lower case.

`max_position_embeddings` int Number of maximum position embeddings. This parameter also determines the maximum length of a sequence which can be processed with the model. Allowed values: $10 \leq x \leq 4048$

`hidden_size` int Number of neurons in each layer. This parameter determines the dimensionality of the resulting text embedding. Allowed values: $1 \leq x \leq 2048$

`target_hidden_size` int Number of neurons in the final layer. This parameter determines the dimensionality of the resulting text embedding. Allowed values: $1 \leq x$

`block_sizes` vector<vector<int>> determining the number and sizes of each block.

`num_attention_heads` int determining the number of attention heads for a self-attention layer. Only relevant if `attention_type='multihead'`. Allowed values: $0 \leq x$

`intermediate_size` int determining the size of the projection layer within a each transformer encoder. Allowed values: $1 \leq x$

`num_decoder_layers` int Number of decoding layers. Allowed values: $1 \leq x$

`pooling_type` string Type of pooling.

- "mean" for pooling with mean.
- "max" for pooling with maximum values. Allowed values: 'Mean', 'Max'

`hidden_act` string Name of the activation function. Allowed values: 'gelu', 'relu', 'silu', 'gelu_new'

`hidden_dropout_prob` double Ratio of dropout. Allowed values: $0 \leq x \leq 0.6$

`attention_probs_dropout_prob` double Ratio of dropout for attention probabilities. Allowed values: $0 \leq x \leq 0.6$

`activation_dropout` double Dropout probability between the layers of the feed-forward blocks. Allowed values: $0 \leq x \leq 0.6$

`sustain_track` bool If TRUE energy consumption is tracked during training via the python library 'codecarbon'.

`sustain_iso_code` string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes. Allowed values: any

`sustain_region` string Region within a country. Only available for USA and Canada See the documentation of codecarbon for more information. <https://mlco2.github.io/codecarbon/parameters.html> Allowed values: any

`sustain_interval` int Interval in seconds for measuring power usage. Allowed values: $1 \leq x$

`trace` bool TRUE if information about the estimation phase should be printed to the console.

`pytorch_safetensors` bool

- TRUE: a 'pytorch' model is saved in safetensors format.
- FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format (.bin).

`log_dir` string Path to the directory where the log files should be saved. If no logging is desired set this argument to NULL. Allowed values: any

`log_write_interval` int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if `log_dir` is not NULL. Allowed values: $1 \leq x$

Returns: This method does not return an object. Instead, it saves the configuration and vocabulary of the new model to disk.

Method train(): This method can be used to train or fine-tune a transformer based on Funnel Transformer architecture with the help of the python libraries transformers, datasets, and tokenizers.

Usage:

```
.AIFEFunnelTransformer$train(
    output_dir,
    model_dir_path,
    text_dataset,
    p_mask = 0.15,
    whole_word = TRUE,
    val_size = 0.1,
    n_epoch = 1,
    batch_size = 12,
    chunk_size = 250,
    full_sequences_only = FALSE,
    min_seq_len = 50,
    learning_rate = 0.003,
    sustain_track = TRUE,
    sustain_iso_code = NULL,
    sustain_region = NULL,
    sustain_interval = 15,
    trace = TRUE,
    pytorch_trace = 1,
    pytorch_safetensors = TRUE,
    log_dir = NULL,
    log_write_interval = 2
)
```

Arguments:

output_dir string Path to the directory where the model should be saved. Allowed values:
any

model_dir_path string Path to the directory where the original model is stored. Allowed
values: any

text_dataset LargeDataSetForText [LargeDataSetForText](#) Object storing textual data.

p_mask double Ratio that determines the number of words/tokens used for masking. Allowed
values: $0 < x < 1$

whole_word bool * TRUE: whole word masking should be applied.

- FALSE: token masking is used.

val_size double between 0 and 1, indicating the proportion of cases which should be used for
the validation sample during the estimation of the model. The remaining cases are part of
the training data. Allowed values: $0 < x < 1$

n_epoch int Number of training epochs. Allowed values: $1 \leq x$

batch_size int Size of the batches for training. Allowed values: $1 \leq x$

chunk_size int Maximum length of every sequence. Must be equal or less the global maxi-
mum size allowed by the model. Allowed values: $100 \leq x$

full_sequences_only bool TRUE for using only chunks with a sequence length equal to
chunk_size.

```

min_seq_len int Only relevant if full_sequences_only = FALSE. Value determines the minimal sequence length included in training process. Allowed values: 10 <= x
learning_rate double Initial learning rate for the training. Allowed values: 0 < x <= 1
sustain_track bool If TRUE energy consumption is tracked during training via the python library 'codecarbon'.
sustain_iso_code string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List\_of\_ISO\_3166\_country\_codes. Allowed values: any
sustain_region string Region within a country. Only available for USA and Canada See the documentation of codecarbon for more information. https://mlco2.github.io/codecarbon/parameters.html Allowed values: any
sustain_interval int Interval in seconds for measuring power usage. Allowed values: 1 <= x
trace bool TRUE if information about the estimation phase should be printed to the console.
pytorch_trace int ml_trace=0 does not print any information about the training process from pytorch on the console. Allowed values: 0 <= x <= 1
pytorch_safetensors bool
    • TRUE: a 'pytorch' model is saved in safetensors format.
    • FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format (.bin).
log_dir string Path to the directory where the log files should be saved. If no logging is desired set this argument to NULL. Allowed values: any
log_write_interval int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if log_dir is not NULL. Allowed values: 1 <= x
>Returns: This method does not return an object. Instead the trained or fine-tuned model is saved to disk.
```

Method clone(): The objects of this class are cloneable with this method.

Usage:

```
.AIFEFunnelTransformer$clone(deep = FALSE)
```

Arguments:

deep Whether to make a deep clone.

Note

The model uses a configuration with truncate_seq = TRUE to avoid implementation problems with tensorflow.

This model uses a WordPiece tokenizer like BERT and can be trained with whole word masking. The transformer library may display a warning, which can be ignored.

References

Dai, Z., Lai, G., Yang, Y. & Le, Q. V. (2020). Funnel-Transformer: Filtering out Sequential Redundancy for Efficient Language Processing. [doi:10.48550/arXiv.2006.03236](https://doi.org/10.48550/arXiv.2006.03236)

Hugging Face documentation

- https://huggingface.co/docs/transformers/model_doc/funnel#funnel-transformer
- https://huggingface.co/docs/transformers/model_doc/funnel#transformers.FunnelModel
- https://huggingface.co/docs/transformers/model_doc/funnel#transformers.TFFunnelModel

See Also

Other R6 classes for transformers: `.AIFEBaseTransformer`, `.AIFEBertTransformer`, `.AIFELongformerTransformer`, `.AIFEMpnetTransformer`, `.AIFERobertaTransformer`

`.AIFELongformerTransformer`

Child R6 class for creation and training of Longformer transformers

Description

This class has the following methods:

- `create`: creates a new transformer based on Longformer.
- `train`: trains and fine-tunes a Longformer model.

Create

New models can be created using the `.AIFELongformerTransformer$create` method.

Train

To train the model, pass the directory of the model to the method `.AIFELongformerTransformer$train`.

Pre-Trained models which can be fine-tuned with this function are available at <https://huggingface.co/>.

Training of this model makes use of dynamic masking.

Super class

`aifedducation::AIFEBaseTransformer -> .AIFELongformerTransformer`

Methods

Public methods:

- `.AIFELongformerTransformer$new()`
- `.AIFELongformerTransformer$create()`
- `.AIFELongformerTransformer$train()`
- `.AIFELongformerTransformer$clone()`

Method `new()`: Creates a new transformer based on Longformer and sets the title.

Usage:

```
.AIFELongformerTransformer$new(init_trace = TRUE)
```

Arguments:

`init_trace` bool option to show prints. If TRUE (by default) - messages will be shown, otherwise (FALSE) - hidden.

Returns: This method returns nothing.

Method `create()`: This method creates a transformer configuration based on the Longformer base architecture and a vocabulary based on Byte-Pair Encoding (BPE) tokenizer using the python `transformers` and `tokenizers` libraries.

This method adds the following '*dependent*' parameters to the base class's inherited `params` list:

- `add_prefix_space`
- `trim_offsets`
- `num_hidden_layer`
- `attention_window`

Usage:

```
.AIFELongformerTransformer$create(
  model_dir,
  text_dataset,
  vocab_size = 30522,
  add_prefix_space = FALSE,
  trim_offsets = TRUE,
  max_position_embeddings = 512,
  hidden_size = 768,
  num_hidden_layer = 12,
  num_attention_heads = 12,
  intermediate_size = 3072,
  hidden_act = "GELU",
  hidden_dropout_prob = 0.1,
  attention_probs_dropout_prob = 0.1,
  attention_window = 512,
  sustain_track = TRUE,
  sustain_iso_code = NULL,
  sustain_region = NULL,
  sustain_interval = 15,
  trace = TRUE,
  pytorch_safetensors = TRUE,
  log_dir = NULL,
  log_write_interval = 2
)
```

Arguments:

`model_dir` string Path to the directory where the model should be saved. Allowed values: any

`text_dataset` LargeDataSetForText [LargeDataSetForText](#) Object storing textual data.

`vocab_size` int Size of the vocabulary. Allowed values: $1000 \leq x \leq 5e+05$

`add_prefix_space` bool TRUE if an additional space should be inserted to the leading words.

`trim_offsets` bool TRUE trims the whitespaces from the produced offsets.

`max_position_embeddings` int Number of maximum position embeddings. This parameter also determines the maximum length of a sequence which can be processed with the model. Allowed values: $10 \leq x \leq 4048$

`hidden_size` int Number of neurons in each layer. This parameter determines the dimensionality of the resulting text embedding. Allowed values: $1 \leq x \leq 2048$

`num_hidden_layer` int Number of hidden layers. Allowed values: $1 \leq x$

`num_attention_heads` int determining the number of attention heads for a self-attention layer. Only relevant if `attention_type='multihead'` Allowed values: $0 \leq x$

`intermediate_size` int determining the size of the projection layer within a each transformer encoder. Allowed values: $1 \leq x$

`hidden_act` string Name of the activation function. Allowed values: 'gelu', 'relu', 'silu', 'gelu_new'

`hidden_dropout_prob` double Ratio of dropout. Allowed values: $0 \leq x \leq 0.6$

`attention_probs_dropout_prob` double Ratio of dropout for attention probabilities. Allowed values: $0 \leq x \leq 0.6$

`attention_window` int Size of the window around each token for attention mechanism in every layer. Allowed values: $2 \leq x$

`sustain_track` bool If TRUE energy consumption is tracked during training via the python library 'codecarbon'.

`sustain_iso_code` string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes. Allowed values: any

`sustain_region` string Region within a country. Only available for USA and Canada See the documentation of codecarbon for more information. [https://mlco2.github.io/codecarbon\(parameters.html\)](https://mlco2.github.io/codecarbon(parameters.html)) Allowed values: any

`sustain_interval` int Interval in seconds for measuring power usage. Allowed values: $1 \leq x$

`trace` bool TRUE if information about the estimation phase should be printed to the console.

`pytorch_safetensors` bool

- TRUE: a 'pytorch' model is saved in safetensors format.
- FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format (.bin).

`log_dir` string Path to the directory where the log files should be saved. If no logging is desired set this argument to NULL. Allowed values: any

`log_write_interval` int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if `log_dir` is not NULL. Allowed values: $1 \leq x$

Returns: This method does not return an object. Instead, it saves the configuration and vocabulary of the new model to disk.

Method train(): This method can be used to train or fine-tune a transformer based on Longformer Transformer architecture with the help of the python libraries `transformers`, `datasets`, and `tokenizers`.

Usage:

```
.AIFELongformerTransformer$train(
  output_dir,
  model_dir_path,
  text_dataset,
  p_mask = 0.15,
  val_size = 0.1,
  n_epoch = 1,
  batch_size = 12,
  chunk_size = 250,
  full_sequences_only = FALSE,
  min_seq_len = 50,
  learning_rate = 0.03,
  sustain_track = TRUE,
  sustain_iso_code = NULL,
  sustain_region = NULL,
  sustain_interval = 15,
  trace = TRUE,
  pytorch_trace = 1,
  pytorch_safetensors = TRUE,
  log_dir = NULL,
  log_write_interval = 2
)
```

Arguments:

`output_dir` string Path to the directory where the model should be saved. Allowed values: any

`model_dir_path` string Path to the directory where the original model is stored. Allowed values: any

`text_dataset` LargeDataSetForText [LargeDataSetForText](#) Object storing textual data.

`p_mask` double Ratio that determines the number of words/tokens used for masking. Allowed values: $0 < x < 1$

`val_size` double between 0 and 1, indicating the proportion of cases which should be used for the validation sample during the estimation of the model. The remaining cases are part of the training data. Allowed values: $0 < x < 1$

`n_epoch` int Number of training epochs. Allowed values: $1 \leq x$

`batch_size` int Size of the batches for training. Allowed values: $1 \leq x$

`chunk_size` int Maximum length of every sequence. Must be equal or less the global maximum size allowed by the model. Allowed values: $100 \leq x$

`full_sequences_only` bool TRUE for using only chunks with a sequence length equal to `chunk_size`.

`min_seq_len` int Only relevant if `full_sequences_only = FALSE`. Value determines the minimal sequence length included in training process. Allowed values: $10 \leq x$

`learning_rate` double Initial learning rate for the training. Allowed values: $0 < x \leq 1$

`sustain_track` bool If TRUE energy consumption is tracked during training via the python library 'codecarbon'.

`sustain_iso_code` string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes. Allowed values: any

```

sustain_region string Region within a country. Only available for USA and Canada See the
documentation of codecarbon for more information. https://mlco2.github.io/codecarbon/
parameters.html Allowed values: any
sustain_interval int Interval in seconds for measuring power usage. Allowed values: 1 <=
x
trace bool TRUE if information about the estimation phase should be printed to the console.
pytorch_trace int ml_trace=0 does not print any information about the training process
from pytorch on the console. Allowed values: 0 <= x <= 1
pytorch_safetensors bool
  • TRUE: a 'pytorch' model is saved in safetensors format.
  • FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format
(.bin).
log_dir string Path to the directory where the log files should be saved. If no logging is
desired set this argument to NULL. Allowed values: any
log_write_interval int Time in seconds determining the interval in which the logger should
try to update the log files. Only relevant if log_dir is not NULL. Allowed values: 1 <= x

```

Returns: This method does not return an object. Instead the trained or fine-tuned model is saved to disk.

Method clone(): The objects of this class are cloneable with this method.

Usage:

```
.AIFELongformerTransformer$clone(deep = FALSE)
```

Arguments:

deep Whether to make a deep clone.

References

Beltagy, I., Peters, M. E., & Cohan, A. (2020). Longformer: The Long-Document Transformer.
[doi:10.48550/arXiv.2004.05150](https://doi.org/10.48550/arXiv.2004.05150)

Hugging Face Documentation

- https://huggingface.co/docs/transformers/model_doc/longformer
- https://huggingface.co/docs/transformers/model_doc/longformer#transformers.LongformerModel
- https://huggingface.co/docs/transformers/model_doc/longformer#transformers.TFLongformerModel

See Also

Other R6 classes for transformers: [.AIFEBaseTransformer](#), [.AIFEBertTransformer](#), [.AIFEFunnelTransformer](#),
[.AIFEMpnetTransformer](#), [.AIFERobertaTransformer](#)

.AIFEModernBertTransformer

Child R6 class for creation and training of ModernBERT transformers

Description

This class has the following methods:

- `create`: creates a new transformer based on ModernBERT.
- `train`: trains and fine-tunes a ModernBERT model.

Create

New models can be created using the `.AIFEModernBertTransformer$create` method.

Train

To train the model, pass the directory of the model to the method `.AIFEModernBertTransformer$train`.

Pre-Trained models that can be fine-tuned using this method are available at <https://huggingface.co/>.

The model is trained using dynamic masking, as opposed to the original paper, which used static masking.

Super class

`aifedducation::AIFEBaseTransformer -> .AIFEModernBertTransformer`

Methods

Public methods:

- `.AIFEModernBertTransformer$new()`
- `.AIFEModernBertTransformer$create()`
- `.AIFEModernBertTransformer$train()`
- `.AIFEModernBertTransformer$clone()`

Method `new()`: Creates a new transformer based on ModernBERT and sets the title.

Usage:

`.AIFEModernBertTransformer$new(init_trace = TRUE)`

Arguments:

`init_trace` bool option to show prints. If TRUE (by default) - messages will be shown, otherwise (FALSE) - hidden.

Returns: This method returns nothing.

Method create(): This method creates a transformer configuration based on the ModernBERT base architecture and a vocabulary based on WordPiece by using the python libraries transformers and tokenizers.

This method adds the following '*dependent*' parameters to the base class's inherited params list:

- vocab_do_lower_case
- num_hidden_layer

Usage:

```
.AIFEModernBertTransformer$create(
    model_dir,
    text_dataset,
    vocab_size = 30522,
    vocab_do_lower_case = FALSE,
    max_position_embeddings = 512,
    hidden_size = 768,
    num_hidden_layer = 12,
    num_attention_heads = 12,
    intermediate_size = 3072,
    hidden_act = "GELU",
    hidden_dropout_prob = 0.1,
    attention_probs_dropout_prob = 0.1,
    sustain_track = FALSE,
    sustain_iso_code = NULL,
    sustain_region = NULL,
    sustain_interval = 15,
    trace = TRUE,
    pytorch_safetensors = TRUE,
    log_dir = NULL,
    log_write_interval = 2
)
```

Arguments:

`model_dir` string Path to the directory where the model should be saved. Allowed values:
any

`text_dataset` LargeDataSetForText [LargeDataSetForText](#) Object storing textual data.

`vocab_size` int Size of the vocabulary. Allowed values: $1000 \leq x \leq 5e+05$

`vocab_do_lower_case` bool TRUE if all words/tokens should be lower case.

`max_position_embeddings` int Number of maximum position embeddings. This parameter also determines the maximum length of a sequence which can be processed with the model.
Allowed values: $10 \leq x \leq 4048$

`hidden_size` int Number of neurons in each layer. This parameter determines the dimensionality of the resulting text embedding. Allowed values: $1 \leq x \leq 2048$

`num_hidden_layer` int Number of hidden layers. Allowed values: $1 \leq x$

`num_attention_heads` int determining the number of attention heads for a self-attention layer.
Only relevant if `attention_type='multihead'` Allowed values: $0 \leq x$

`intermediate_size` int determining the size of the projection layer within a each transformer encoder. Allowed values: $1 \leq x$

```

hidden_act string Name of the activation function. Allowed values: 'gelu', 'relu', 'silu',
    'gelu_new'
hidden_dropout_prob double Ratio of dropout. Allowed values: 0 <= x <= 0.6
attention_probs_dropout_prob double Ratio of dropout for attention probabilities. Al-
    lowed values: 0 <= x <= 0.6
sustain_track bool If TRUE energy consumption is tracked during training via the python
    library 'codecarbon'.
sustain_iso_code string ISO code (Alpha-3-Code) for the country. This variable must be
    set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List\_of\_ISO\_3166\_country\_codes. Allowed values: any
sustain_region string Region within a country. Only available for USA and Canada See the
    documentation of codecarbon for more information. https://mlco2.github.io/codecarbon/parameters.html Allowed values: any
sustain_interval int Interval in seconds for measuring power usage. Allowed values: 1 <=
    x
trace bool TRUE if information about the estimation phase should be printed to the console.
pytorch_safetensors bool
    • TRUE: a 'pytorch' model is saved in safetensors format.
    • FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format
        (.bin).
log_dir string Path to the directory where the log files should be saved. If no logging is
    desired set this argument to NULL. Allowed values: any
log_write_interval int Time in seconds determining the interval in which the logger should
    try to update the log files. Only relevant if log_dir is not NULL. Allowed values: 1 <= x
>Returns: This method does not return an object. Instead, it saves the configuration and vocab-
    uary of the new model to disk.
```

Method train(): This method can be used to train or fine-tune a transformer based on ModernBERT architecture with the help of the python libraries `transformers`, `datasets`, and `tokenizers`.

Usage:

```
.AIFEModernBertTransformer$train(
  output_dir,
  model_dir_path,
  text_dataset,
  p_mask = 0.15,
  whole_word = TRUE,
  val_size = 0.1,
  n_epoch = 1,
  batch_size = 12,
  chunk_size = 250,
  full_sequences_only = FALSE,
  min_seq_len = 50,
  learning_rate = 0.003,
  sustain_track = FALSE,
  sustain_iso_code = NULL,
  sustain_region = NULL,
```

```

    sustain_interval = 15,
    trace = TRUE,
    pytorch_trace = 1,
    pytorch_safetensors = TRUE,
    log_dir = NULL,
    log_write_interval = 2
)
Arguments:
output_dir string Path to the directory where the model should be saved. Allowed values:
any
model_dir_path string Path to the directory where the original model is stored. Allowed
values: any
text_dataset LargeDataSetForText LargeDataSetForText Object storing textual data.
p_mask double Ratio that determines the number of words/tokens used for masking. Allowed
values:  $0 < x < 1$ 
whole_word bool * TRUE: whole word masking should be applied.
• FALSE: token masking is used.
val_size double between 0 and 1, indicating the proportion of cases which should be used for
the validation sample during the estimation of the model. The remaining cases are part of
the training data. Allowed values:  $0 < x < 1$ 
n_epoch int Number of training epochs. Allowed values:  $1 \leq x$ 
batch_size int Size of the batches for training. Allowed values:  $1 \leq x$ 
chunk_size int Maximum length of every sequence. Must be equal or less the global maxi-
mum size allowed by the model. Allowed values:  $100 \leq x$ 
full_sequences_only bool TRUE for using only chunks with a sequence length equal to
chunk_size.
min_seq_len bool TRUE for using only chunks with a sequence length equal to chunk_size.
learning_rate double Initial learning rate for the training. Allowed values:  $0 < x \leq 1$ 
sustain_track bool If TRUE energy consumption is tracked during training via the python
library 'codecarbon'.
sustain_iso_code string ISO code (Alpha-3-Code) for the country. This variable must be
set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List\_of\_ISO\_3166\_country\_codes. Allowed values: any
sustain_region string Region within a country. Only available for USA and Canada See the
documentation of codecarbon for more information. https://mlco2.github.io/codecarbon/parameters.html Allowed values: any
sustain_interval int Interval in seconds for measuring power usage. Allowed values:  $1 \leq
x$ 
trace bool TRUE if information about the estimation phase should be printed to the console.
pytorch_trace int ml_trace=0 does not print any information about the training process
from pytorch on the console. Allowed values:  $0 \leq x \leq 1$ 
pytorch_safetensors bool
• TRUE: a 'pytorch' model is saved in safetensors format.
• FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format
(.bin).

```

`log_dir` string Path to the directory where the log files should be saved. If no logging is desired set this argument to NULL. Allowed values: any

`log_write_interval` int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if `log_dir` is not NULL. Allowed values: $1 \leq x$

Returns: This method does not return an object. Instead the trained or fine-tuned model is saved to disk.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
.AIFEModernBertTransformer$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

Note

This model uses a WordPiece tokenizer like BERT and can be trained with whole word masking. The transformer library may display a warning, which can be ignored.

References

Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In J. Burstein, C. Doran, & T. Solorio (Eds.), Proceedings of the 2019 Conference of the North (pp. 4171–4186). Association for Computational Linguistics. [doi:10.18653/v1/N191423](https://doi.org/10.18653/v1/N191423)

Hugging Face documentation

- https://huggingface.co/docs/transformers/model_doc/bert
- https://huggingface.co/docs/transformers/model_doc/bert#transformers.BertForMaskedLM
- https://huggingface.co/docs/transformers/model_doc/bert#transformers.TFBertForMaskedLM

See Also

Other Transformers for developers: `.AIFETrConfig`, `.AIFETrModel`, `.AIFETrModelMLM`, `.AIFETrObj`, `.AIFETrTokenizer`, `.aife_transformer.check_type()`, `.aife_transformer.load_model()`, `.aife_transformer.load_aife_transformer.load_model_mlm()`, `.aife_transformer.load_tokenizer()`

.AIFEMpnetTransformer Child R6 class for creation and training of MPNet transformers

Description

This class has the following methods:

- `create`: creates a new transformer based on MPNet.
- `train`: trains and fine-tunes a MPNet model.

Create

New models can be created using the `.AIFEMpnetTransformer$create` method.

Train

To train the model, pass the directory of the model to the method `.AIFEMpnetTransformer$train`.

Super class

`aifedducation::AIFEBaseTransformer -> .AIFEMpnetTransformer`

Public fields

`special_tokens_list` list List for special tokens with the following elements:

- `cls` - CLS token representation (`<s>`)
- `pad` - pad token representation (`<pad>`)
- `sep` - sep token representation (`</s>`)
- `unk` - unk token representation (`<unk>`)
- `mask` - mask token representation (`<mask>`)

Methods

Public methods:

- `.AIFEMpnetTransformer$new()`
- `.AIFEMpnetTransformer$create()`
- `.AIFEMpnetTransformer$train()`
- `.AIFEMpnetTransformer$clone()`

Method `new()`: Creates a new transformer based on MPNet and sets the title.

Usage:

`.AIFEMpnetTransformer$new(init_trace = TRUE)`

Arguments:

`init_trace` bool option to show prints. If `TRUE` (by default) - messages will be shown, otherwise (`FALSE`) - hidden.

Returns: This method returns nothing.

Method `create()`: This method creates a transformer configuration based on the MPNet base architecture.

This method adds the following '*dependent*' parameters to the base class's inherited `params` list:

- `vocab_do_lower_case`
- `num_hidden_layer`

Usage:

```
.AIFEMpnetTransformer$create(
    model_dir,
    text_dataset,
    vocab_size = 30522,
    vocab_do_lower_case = FALSE,
    max_position_embeddings = 512,
    hidden_size = 768,
    num_hidden_layer = 12,
    num_attention_heads = 12,
    intermediate_size = 3072,
    hidden_act = "GELU",
    hidden_dropout_prob = 0.1,
    attention_probs_dropout_prob = 0.1,
    sustain_track = FALSE,
    sustain_iso_code = NULL,
    sustain_region = NULL,
    sustain_interval = 15,
    trace = TRUE,
    pytorch_safetensors = TRUE,
    log_dir = NULL,
    log_write_interval = 2
)
```

Arguments:

`model_dir` string Path to the directory where the model should be saved. Allowed values:
any

`text_dataset` LargeDataSetForText `LargeDataSetForText` Object storing textual data.

`vocab_size` int Size of the vocabulary. Allowed values: $1000 \leq x \leq 5e+05$

`vocab_do_lower_case` bool TRUE if all words/tokens should be lower case.

`max_position_embeddings` int Number of maximum position embeddings. This parameter
also determines the maximum length of a sequence which can be processed with the model.
Allowed values: $10 \leq x \leq 4048$

`hidden_size` int Number of neurons in each layer. This parameter determines the dimension-
ality of the resulting text embedding. Allowed values: $1 \leq x \leq 2048$

`num_hidden_layer` int Number of hidden layers. Allowed values: $1 \leq x$

`num_attention_heads` int determining the number of attention heads for a self-attention layer.
Only relevant if `attention_type='multihead'` Allowed values: $0 \leq x$

`intermediate_size` int determining the size of the projection layer within a each transformer
encoder. Allowed values: $1 \leq x$

`hidden_act` string Name of the activation function. Allowed values: 'gelu', 'relu', 'silu',
'gelu_new'

`hidden_dropout_prob` double Ratio of dropout. Allowed values: $0 \leq x \leq 0.6$

`attention_probs_dropout_prob` double Ratio of dropout for attention probabilities. Al-
lowed values: $0 \leq x \leq 0.6$

`sustain_track` bool If TRUE energy consumption is tracked during training via the python
library 'codecarbon'.

sustain_iso_code string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes. Allowed values: any

sustain_region string Region within a country. Only available for USA and Canada See the documentation of codecarbon for more information. [https://mlco2.github.io/codecarbon\(parameters.html](https://mlco2.github.io/codecarbon(parameters.html)) Allowed values: any

sustain_interval int Interval in seconds for measuring power usage. Allowed values: $1 \leq x$

trace bool TRUE if information about the estimation phase should be printed to the console.

pytorch_safetensors bool

- TRUE: a 'pytorch' model is saved in safetensors format.
- FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format (.bin).

log_dir string Path to the directory where the log files should be saved. If no logging is desired set this argument to NULL. Allowed values: any

log_write_interval int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if log_dir is not NULL. Allowed values: $1 \leq x$

Returns: This method does not return an object. Instead, it saves the configuration and vocabulary of the new model to disk.

Method train(): This method can be used to train or fine-tune a transformer based on MPNet architecture with the help of the python libraries transformers, datasets, and tokenizers.

This method adds the following '*dependent*' parameter to the base class's inherited params list:

- p_perm

Usage:

```
.AIFEMpnetTransformer$train(
  output_dir,
  model_dir_path,
  text_dataset,
  p_mask = 0.15,
  p_perm = 0.15,
  whole_word = TRUE,
  val_size = 0.1,
  n_epoch = 1,
  batch_size = 12,
  chunk_size = 250,
  full_sequences_only = FALSE,
  min_seq_len = 50,
  learning_rate = 0.003,
  sustain_track = FALSE,
  sustain_iso_code = NULL,
  sustain_region = NULL,
  sustain_interval = 15,
  trace = TRUE,
  pytorch_trace = 1,
  pytorch_safetensors = TRUE,
```

```

    log_dir = NULL,
    log_write_interval = 2
)
Arguments:
output_dir string Path to the directory where the model should be saved. Allowed values:
any
model_dir_path string Path to the directory where the original model is stored. Allowed
values: any
text_dataset LargeDataSetForText LargeDataSetForText Object storing textual data.
p_mask double Ratio that determines the number of words/tokens used for masking. Allowed
values:  $0 < x < 1$ 
p_perm double Ratio that determines the number of words/tokens used for permutation.
whole_word bool * TRUE: whole word masking should be applied.
• FALSE: token masking is used.
val_size double between 0 and 1, indicating the proportion of cases which should be used for
the validation sample during the estimation of the model. The remaining cases are part of
the training data. Allowed values:  $0 < x < 1$ 
n_epoch int Number of training epochs. Allowed values:  $1 \leq x$ 
batch_size int Size of the batches for training. Allowed values:  $1 \leq x$ 
chunk_size int Maximum length of every sequence. Must be equal or less the global maxi-
mum size allowed by the model. Allowed values:  $100 \leq x$ 
full_sequences_only bool TRUE for using only chunks with a sequence length equal to
chunk_size.
min_seq_len int Only relevant if full_sequences_only = FALSE. Value determines the min-
imal sequence length included in training process. Allowed values:  $10 \leq x$ 
learning_rate double Initial learning rate for the training. Allowed values:  $0 < x \leq 1$ 
sustain_track bool If TRUE energy consumption is tracked during training via the python
library 'codecarbon'.
sustain_iso_code string ISO code (Alpha-3-Code) for the country. This variable must be
set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List\_of\_ISO\_3166\_country\_codes. Allowed values: any
sustain_region string Region within a country. Only available for USA and Canada See the
documentation of codecarbon for more information. https://mlco2.github.io/codecarbon/parameters.html Allowed values: any
sustain_interval int Interval in seconds for measuring power usage. Allowed values:  $1 \leq
x$ 
trace bool TRUE if information about the estimation phase should be printed to the console.
pytorch_trace int ml_trace=0 does not print any information about the training process
from pytorch on the console. Allowed values:  $0 \leq x \leq 1$ 
pytorch_safetensors bool
• TRUE: a 'pytorch' model is saved in safetensors format.
• FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format
(.bin).
log_dir string Path to the directory where the log files should be saved. If no logging is
desired set this argument to NULL. Allowed values: any

```

`log_write_interval int` Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if `log_dir` is not NULL. Allowed values: $1 \leq x$

Returns: This method does not return an object. Instead the trained or fine-tuned model is saved to disk.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

`.AIFEMpnetTransformer$clone(deep = FALSE)`

Arguments:

`deep` Whether to make a deep clone.

Note

Using this class with tensorflow is not supported. Supported framework is pytorch.

References

Song,K., Tan, X., Qin, T., Lu, J. & Liu, T.-Y. (2020). MPNet: Masked and Permuted Pre-training for Language Understanding. [doi:10.48550/arXiv.2004.09297](https://doi.org/10.48550/arXiv.2004.09297)

Hugging Face documentation

- https://huggingface.co/docs/transformers/model_doc/mpnet
- https://huggingface.co/docs/transformers/model_doc/mpnet#transformers.MPNetForMaskedLM
- https://huggingface.co/docs/transformers/model_doc/mpnet#transformers.TFPMPNetForMaskedLM

See Also

Other R6 classes for transformers: [.AIFEBaseTransformer](#), [.AIFEBertTransformer](#), [.AIFEFunnelTransformer](#), [.AIFELongformerTransformer](#), [.AIFERobertaTransformer](#)

.AIFERobertaTransformer

Child R6 class for creation and training of RoBERTa transformers

Description

This class has the following methods:

- `create`: creates a new transformer based on RoBERTa.
- `train`: trains and fine-tunes a RoBERTa model.

Create

New models can be created using the `.AIFERobertaTransformer$create` method.

Train

To train the model, pass the directory of the model to the method `.AIFERobertaTransformer$train`.

Pre-Trained models which can be fine-tuned with this function are available at <https://huggingface.co/>.

Training of this model makes use of dynamic masking.

Super class

```
aifedducation:: .AIFEBaseTransformer -> .AIFERobertaTransformer
```

Methods

Public methods:

- `.AIFERobertaTransformer$new()`
- `.AIFERobertaTransformer$create()`
- `.AIFERobertaTransformer$train()`
- `.AIFERobertaTransformer$clone()`

Method `new()`: Creates a new transformer based on RoBERTa and sets the title.

Usage:

```
.AIFERobertaTransformer$new(init_trace = TRUE)
```

Arguments:

`init_trace` bool option to show prints. If `TRUE` (by default) - messages will be shown, otherwise (`FALSE`) - hidden.

Returns: This method returns nothing.

Method `create()`: This method creates a transformer configuration based on the RoBERTa base architecture and a vocabulary based on Byte-Pair Encoding (BPE) tokenizer using the python `transformers` and `tokenizers` libraries.

This method adds the following '*dependent*' parameters to the base class' inherited `params` list:

- `add_prefix_space`
- `trim_offsets`
- `num_hidden_layer`

Usage:

```
.AIFERobertaTransformer$create(
  model_dir,
  text_dataset,
  vocab_size = 30522,
  add_prefix_space = FALSE,
  trim_offsets = TRUE,
  max_position_embeddings = 512,
  hidden_size = 768,
  num_hidden_layer = 12,
  num_attention_heads = 12,
  intermediate_size = 3072,
```

```

hidden_act = "GELU",
hidden_dropout_prob = 0.1,
attention_probs_dropout_prob = 0.1,
sustain_track = TRUE,
sustain_iso_code = NULL,
sustain_region = NULL,
sustain_interval = 15,
trace = TRUE,
pytorch_safetensors = TRUE,
log_dir = NULL,
log_write_interval = 2
)

```

Arguments:

`model_dir` string Path to the directory where the model should be saved. Allowed values:
any

`text_dataset` LargeDataSetForText `LargeDataSetForText` Object storing textual data.

`vocab_size` int Size of the vocabulary. Allowed values: $1000 \leq x \leq 5e+05$

`add_prefix_space` bool TRUE if an additional space should be inserted to the leading words.

`trim_offsets` bool TRUE trims the whitespaces from the produced offsets.

`max_position_embeddings` int Number of maximum position embeddings. This parameter also determines the maximum length of a sequence which can be processed with the model.
Allowed values: $10 \leq x \leq 4048$

`hidden_size` int Number of neurons in each layer. This parameter determines the dimensionality of the resulting text embedding. Allowed values: $1 \leq x \leq 2048$

`num_hidden_layer` int Number of hidden layers. Allowed values: $1 \leq x$

`num_attention_heads` int determining the number of attention heads for a self-attention layer.
Only relevant if `attention_type='multihead'` Allowed values: $0 \leq x$

`intermediate_size` int determining the size of the projection layer within a each transformer encoder. Allowed values: $1 \leq x$

`hidden_act` string Name of the activation function. Allowed values: 'gelu', 'relu', 'silu', 'gelu_new'

`hidden_dropout_prob` double Ratio of dropout. Allowed values: $0 \leq x \leq 0.6$

`attention_probs_dropout_prob` double Ratio of dropout for attention probabilities. Allowed values: $0 \leq x \leq 0.6$

`sustain_track` bool If TRUE energy consumption is tracked during training via the python library 'codecarbon'.

`sustain_iso_code` string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes. Allowed values: any

`sustain_region` string Region within a country. Only available for USA and Canada See the documentation of codecarbon for more information. <https://mlco2.github.io/codecarbon/parameters.html> Allowed values: any

`sustain_interval` int Interval in seconds for measuring power usage. Allowed values: $1 \leq x$

`trace` bool TRUE if information about the estimation phase should be printed to the console.

`pytorch_safetensors` bool
 • TRUE: a 'pytorch' model is saved in safetensors format.
 • FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format (.bin).

`log_dir` string Path to the directory where the log files should be saved. If no logging is desired set this argument to NULL. Allowed values: any

`log_write_interval` int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if `log_dir` is not NULL. Allowed values: $1 \leq x$

Returns: This method does not return an object. Instead, it saves the configuration and vocabulary of the new model to disk.

Method train(): This method can be used to train or fine-tune a transformer based on RoBERTa Transformer architecture with the help of the python libraries `transformers`, `datasets`, and `tokenizers`.

Usage:

```
.AIFERobertaTransformer$train(  
    output_dir,  
    model_dir_path,  
    text_dataset,  
    p_mask = 0.15,  
    val_size = 0.1,  
    n_epoch = 1,  
    batch_size = 12,  
    chunk_size = 250,  
    full_sequences_only = FALSE,  
    min_seq_len = 50,  
    learning_rate = 0.03,  
    sustain_track = TRUE,  
    sustain_iso_code = NULL,  
    sustain_region = NULL,  
    sustain_interval = 15,  
    trace = TRUE,  
    pytorch_trace = 1,  
    pytorch_safetensors = TRUE,  
    log_dir = NULL,  
    log_write_interval = 2  
)
```

Arguments:

`output_dir` string Path to the directory where the model should be saved. Allowed values: any

`model_dir_path` string Path to the directory where the original model is stored. Allowed values: any

`text_dataset` LargeDataSetForText [LargeDataSetForText](#) Object storing textual data.

`p_mask` double Ratio that determines the number of words/tokens used for masking. Allowed values: $0 < x < 1$

`val_size` double between 0 and 1, indicating the proportion of cases which should be used for the validation sample during the estimation of the model. The remaining cases are part of the training data. Allowed values: $0 < x < 1$

`n_epoch` int Number of training epochs. Allowed values: $1 \leq x$

`batch_size` int Size of the batches for training. Allowed values: $1 \leq x$

`chunk_size` int Maximum length of every sequence. Must be equal or less the global maximum size allowed by the model. Allowed values: $100 \leq x$

`full_sequences_only` bool TRUE for using only chunks with a sequence length equal to `chunk_size`.

`min_seq_len` int Only relevant if `full_sequences_only` = FALSE. Value determines the minimal sequence length included in training process. Allowed values: $10 \leq x$

`learning_rate` double Initial learning rate for the training. Allowed values: $0 < x \leq 1$

`sustain_track` bool If TRUE energy consumption is tracked during training via the python library 'codecarbon'.

`sustain_iso_code` string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes. Allowed values: any

`sustain_region` string Region within a country. Only available for USA and Canada See the documentation of codecarbon for more information. <https://mlco2.github.io/codecarbon/parameters.html> Allowed values: any

`sustain_interval` int Interval in seconds for measuring power usage. Allowed values: $1 \leq x$

`trace` bool TRUE if information about the estimation phase should be printed to the console.

`pytorch_trace` int ml_trace=0 does not print any information about the training process from pytorch on the console. Allowed values: $0 \leq x \leq 1$

`pytorch_safetensors` bool

- TRUE: a 'pytorch' model is saved in safetensors format.
- FALSE (or 'safetensors' is not available): model is saved in the standard pytorch format (.bin).

`log_dir` string Path to the directory where the log files should be saved. If no logging is desired set this argument to NULL. Allowed values: any

`log_write_interval` int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if `log_dir` is not NULL. Allowed values: $1 \leq x$

Returns: This method does not return an object. Instead the trained or fine-tuned model is saved to disk.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
.AIFERobertaTransformer$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

References

Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. [doi:10.48550/arXiv.1907.11692](https://doi.org/10.48550/arXiv.1907.11692)

Hugging Face Documentation

- https://huggingface.co/docs/transformers/model_doc/roberta
- https://huggingface.co/docs/transformers/model_doc/roberta#transformers.RobertaModel
- https://huggingface.co/docs/transformers/model_doc/roberta#transformers.TFRobertaModel

See Also

Other R6 classes for transformers: [.AIFEBaseTransformer](#), [.AIFEBertTransformer](#), [.AIFEFunnelTransformer](#), [.AIFELongformerTransformer](#), [.AIFEMpnetTransformer](#)

add_missing_args *Add missing arguments to a list of arguments*

Description

This function is designed for taking the output of `summarize_args_for_long_task` as input. It adds the missing arguments. In general these are arguments that rely on objects of class R6 which can not be exported to a new R session.

Usage

```
add_missing_args(args, path_args, meta_args)
```

Arguments

- | | |
|-----------|--|
| args | Named list List for arguments for the method of a specific class. |
| path_args | Named list List of paths where the objects are stored on disk. |
| meta_args | Named list List containing arguments that are necessary in order to add the missing objects correctly. |

Value

Returns a named list of all arguments that a method of a specific class requires.

See Also

Other Utils Studio Developers: [create_data_embeddings_description\(\)](#), [long_load_target_data\(\)](#), [summarize_args_for_long_task\(\)](#)

<code>AIFEBASEMODEL</code>	<i>Base class for models using neural nets</i>
----------------------------	--

Description

Abstract class for all models that do not rely on the python library 'transformers'. Objects of this class containing fields and methods used in several other classes in 'AI for Education'.

This class is **not** designed for a direct application and should only be used by developers.

Value

A new object of this class.

Public fields

`model ('pytorch_model')`

Field for storing a 'pytorch' model after loading.

`model_config ('list()')`

List for storing information about the configuration of the model.

`last_training ('list()')`

List for storing the history, the configuration, and the results of the last training. This information will be overwritten if a new training is started.

- `last_training$start_time`: Time point when training started.

- `last_training$learning_time`: Duration of the training process.

- `last_training$finish_time`: Time when the last training finished.

- `last_training$history`: History of the last training.

- `last_training$data`: Object of class table storing the initial frequencies of the passed data.

- `last_training$config`: List storing the configuration used for the last training.

Methods

Public methods:

- [AIFEBASEMODEL\\$get_model_info\(\)](#)
- [AIFEBASEMODEL\\$set_publication_info\(\)](#)
- [AIFEBASEMODEL\\$get_publication_info\(\)](#)
- [AIFEBASEMODEL\\$set_model_license\(\)](#)
- [AIFEBASEMODEL\\$get_model_license\(\)](#)
- [AIFEBASEMODEL\\$set_documentation_license\(\)](#)
- [AIFEBASEMODEL\\$get_documentation_license\(\)](#)
- [AIFEBASEMODEL\\$set_model_description\(\)](#)
- [AIFEBASEMODEL\\$get_model_description\(\)](#)
- [AIFEBASEMODEL\\$save\(\)](#)

- `AIFEBaseModel$load()`
- `AIFEBaseModel$get_package_versions()`
- `AIFEBaseModel$get_sustainability_data()`
- `AIFEBaseModel$get_ml_framework()`
- `AIFEBaseModel$count_parameter()`
- `AIFEBaseModel$is_configured()`
- `AIFEBaseModel$is_trained()`
- `AIFEBaseModel$get_private()`
- `AIFEBaseModel$get_all_fields()`
- `AIFEBaseModel$clone()`

Method `get_model_info()`: Method for requesting the model information.

Usage:

`AIFEBaseModel$get_model_info()`

Returns: list of all relevant model information.

Method `set_publication_info()`: Method for setting publication information of the model.

Usage:

`AIFEBaseModel$set_publication_info(authors, citation, url = NULL)`

Arguments:

`authors` List of authors.

`citation` Free text citation.

`url` URL of a corresponding homepage.

Returns: Function does not return a value. It is used for setting the private members for publication information.

Method `get_publication_info()`: Method for requesting the bibliographic information of the model.

Usage:

`AIFEBaseModel$get_publication_info()`

Returns: list with all saved bibliographic information.

Method `set_model_license()`: Method for setting the license of the model.

Usage:

`AIFEBaseModel$set_model_license(license = "CC BY")`

Arguments:

`license` string containing the abbreviation of the license or the license text.

Returns: Function does not return a value. It is used for setting the private member for the software license of the model.

Method `get_model_license()`: Method for getting the license of the model.

Usage:

`AIFEBaseModel$get_model_license()`

Arguments:

`license` string containing the abbreviation of the license or the license text.

Returns: string representing the license for the model.

Method `set_documentation_license()`: Method for setting the license of the model's documentation.

Usage:

```
AIFEBaseModel$set_documentation_license(license = "CC BY")
```

Arguments:

`license` string containing the abbreviation of the license or the license text.

Returns: Function does not return a value. It is used for setting the private member for the documentation license of the model.

Method `get_documentation_license()`: Method for getting the license of the model's documentation.

Usage:

```
AIFEBaseModel$get_documentation_license()
```

Arguments:

`license` string containing the abbreviation of the license or the license text.

Returns: Returns the license as a string.

Method `set_model_description()`: Method for setting a description of the model.

Usage:

```
AIFEBaseModel$set_model_description(  
  eng = NULL,  
  native = NULL,  
  abstract_eng = NULL,  
  abstract_native = NULL,  
  keywords_eng = NULL,  
  keywords_native = NULL  
)
```

Arguments:

`eng` string A text describing the training, its theoretical and empirical background, and output in English.

`native` string A text describing the training , its theoretical and empirical background, and output in the native language of the model.

`abstract_eng` string A text providing a summary of the description in English.

`abstract_native` string A text providing a summary of the description in the native language of the model.

`keywords_eng` vector of keyword in English.

`keywords_native` vector of keyword in the native language of the model.

Returns: Function does not return a value. It is used for setting the private members for the description of the model.

Method `get_model_description()`: Method for requesting the model description.

Usage:

```
AIFEBaseModel$get_model_description()
```

Returns: list with the description of the classifier in English and the native language.

Method `save()`: Method for saving a model.

Usage:

```
AIFEBaseModel$save(dir_path, folder_name)
```

Arguments:

`dir_path` string Path of the directory where the model should be saved.

`folder_name` string Name of the folder that should be created within the directory.

Returns: Function does not return a value. It saves the model to disk.

Method `load()`: Method for importing a model.

Usage:

```
AIFEBaseModel$load(dir_path)
```

Arguments:

`dir_path` string Path of the directory where the model is saved.

Returns: Function does not return a value. It is used to load the weights of a model.

Method `get_package_versions()`: Method for requesting a summary of the R and python packages' versions used for creating the model.

Usage:

```
AIFEBaseModel$get_package_versions()
```

Returns: Returns a list containing the versions of the relevant R and python packages.

Method `get_sustainability_data()`: Method for requesting a summary of tracked energy consumption during training and an estimate of the resulting CO2 equivalents in kg.

Usage:

```
AIFEBaseModel$get_sustainability_data()
```

Returns: Returns a list containing the tracked energy consumption, CO2 equivalents in kg, information on the tracker used, and technical information on the training infrastructure.

Method `get_ml_framework()`: Method for requesting the machine learning framework used for the model.

Usage:

```
AIFEBaseModel$get_ml_framework()
```

Returns: Returns a string describing the machine learning framework used for the classifier.

Method `count_parameter()`: Method for counting the trainable parameters of a model.

Usage:

```
AIFEBaseModel$count_parameter()
```

Returns: Returns the number of trainable parameters of the model.

Method `is_configured()`: Method for checking if the model was successfully configured. An object can only be used if this value is TRUE.

Usage:

```
AIFEBaseModel$is_configured()
```

Returns: bool TRUE if the model is fully configured. FALSE if not.

Method `is_trained()`: Check if the [TEFeatureExtractor](#) is trained.

Usage:

```
AIFEBaseModel$is_trained()
```

Returns: Returns TRUE if the object is trained and FALSE if not.

Method `get_private()`: Method for requesting all private fields and methods. Used for loading and updating an object.

Usage:

```
AIFEBaseModel$get_private()
```

Returns: Returns a list with all private fields and methods.

Method `get_all_fields()`: Return all fields.

Usage:

```
AIFEBaseModel$get_all_fields()
```

Returns: Method returns a list containing all public and private fields of the object.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
AIFEBaseModel$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

See Also

Other R6 Classes for Developers: [ClassifiersBasedOnTextEmbeddings](#), [LargeDataSetBase](#), [ModelsBasedOnTextEmbeddings](#), [TEClassifiersBasedOnProtoNet](#), [TEClassifiersBasedOnRegular](#)

AIFETrType

Transformer types

Description

This list contains transformer types. Elements of the list can be used in [aife_transformer.make\(\)](#) function as input parameter type.

It has the following elements:

- bert = 'bert'
- roberta = 'roberta'
- funnel = 'funnel'
- longformer = 'longformer'
- mpnet = 'mpnet'
- modernbert = 'modernbert'

Elements can be used like AIFETrType\$bert, AIFETrType\$modernbert, AIFETrType\$funnel, etc.

Usage

AIFETrType

Format

An object of class list of length 6.

See Also

Other Transformer: [aife_transformer.make\(\)](#)

aife_transformer.load_model_mlm
Load a MLM-model

Description

Loads a MLM-model for a transformer with the passed type.

Usage

`aife_transformer.load_model_mlm(type, model_dir, from_tf, load_safe)`

Arguments

<code>type</code>	string A type of the transformer. Allowed types are bert, roberta, funnel, longformer, mpnet, modernbert. See AIFETrType list.
<code>model_dir</code>	string Path to the directory where the original model is stored. Allowed values: any
<code>from_tf</code>	bool Whether to load the model weights from a TensorFlow checkpoint save file.
<code>load_safe</code>	bool Whether or not to use safetensors checkpoints.

Value

If success - a MLM-model, otherwise - an error (passed type is invalid).

See Also

Other Transformers for developers: [.AIFEModernBertTransformer](#), [.AIFETrConfig](#), [.AIFETrModel](#), [.AIFETrModelMLM](#), [.AIFETrObj](#), [.AIFETrTokenizer](#), [.aife_transformer.check_type\(\)](#), [.aife_transformer.load_model_config\(\)](#), [.aife_transformer.load_tokenizer\(\)](#)

`aife_transformer.load_tokenizer`
Load a tokenizer

Description

Loads a tokenizer for a transformer with the passed type.

Usage

```
aife_transformer.load_tokenizer(type, model_dir)
```

Arguments

<code>type</code>	string A type of the transformer. Allowed types are bert, roberta, funnel, longformer, mpnet, modernbert. See AIFETrType list.
<code>model_dir</code>	string Path to the directory where the original model is stored. Allowed values: any

Value

If success - a tokenizer, otherwise - an error (passed type is invalid).

See Also

Other Transformers for developers: [.AIFEModernBertTransformer](#), [.AIFETrConfig](#), [.AIFETrModel](#), [.AIFETrModelMLM](#), [.AIFETrObj](#), [.AIFETrTokenizer](#), [.aife_transformer.check_type\(\)](#), [.aife_transformer.load_model_config\(\)](#), [.aife_transformer.load_model_mlm\(\)](#)

aife_transformer.make *Make a transformer*

Description

Creates a new transformer with the passed type. See p.3 Transformer Maker in [Transformers for Developers](#) for details.

See [.AIFEBaseTransformer](#) class for details.

Usage

```
aife_transformer.make(type, init_trace = TRUE)
```

Arguments

type	string A type of the new transformer. Allowed types are bert, roberta, funnel, longformer, mpnet, modernbert. See AIFETrType list.
init_trace	bool option to show prints. If TRUE (by default) - messages will be shown, otherwise (FALSE) - hidden.

Value

If success - a new transformer, otherwise - an error (passed type is invalid).

See Also

Other Transformer: [AIFETrType](#)

auto_n_cores *Number of cores for multiple tasks*

Description

Function for getting the number of cores that should be used for parallel processing of tasks. The number of cores is set to 75 % of the available cores. If the environment variable CI is set to "true" or if the process is running on cran 2 is returned.

Usage

```
auto_n_cores()
```

Value

Returns int as the number of cores.

See Also

Other Utils Developers: [create_object\(\)](#), [create_synthetic_units_from_matrix\(\)](#), [generate_id\(\)](#), [get_n_chunks\(\)](#), [get_synthetic_cases_from_matrix\(\)](#), [matrix_to_array_c\(\)](#), [tensor_to_matrix_c\(\)](#), [to_categorical_c\(\)](#)

build_documentation_for_model
Generate documentation for a classifier class

Description

Function for generating the documentation of a model.

Usage

```
build_documentation_for_model(
  model_name,
  cls_type = NULL,
  core_type = NULL,
  input_type = "text_embeddings"
)
```

Arguments

model_name	string	Name of the model.
cls_type	string	Type of classification
core_type	string	Name of the core type.
input_type	bool	Name of the input type necessary for training and predicting.

Value

Returns a string containing the description written in rmarkdown.

Note

Function is designed to be used with roxygen2 in the regular documentation.

See Also

Other Utils Documentation: [build_layer_stack_documentation_for_vignette\(\)](#), [get_desc_for_core_model_architecture\(\)](#), [get_dict_cls_type\(\)](#), [get_dict_core_models\(\)](#), [get_dict_input_types\(\)](#), [get_layer_dict\(\)](#), [get_layer_documentation\(\)](#), [get_parameter_documentation\(\)](#)

build_layer_stack_documentation_for_vignette

Generate documentation of all layers for an vignette or article

Description

Function for generating the whole documentation for an article used on the packages home page.

Usage

```
build_layer_stack_documentation_for_vignette()
```

Value

Returns a string containing the description written in rmarkdown.

Note

Function is designed to be used with inline r code in rmarkdown vignettes/articles.

See Also

Other Utils Documentation: [build_documentation_for_model\(\)](#), [get_desc_for_core_model_architecture\(\)](#), [get_dict_cls_type\(\)](#), [get_dict_core_models\(\)](#), [get_dict_input_types\(\)](#), [get_layer_dict\(\)](#), [get_layer_documentation\(\)](#), [get_parameter_documentation\(\)](#)

calc_standard_classification_measures

Calculate recall, precision, and f1-scores

Description

Function for calculating recall, precision, and f1-scores.

Usage

```
calc_standard_classification_measures(true_values, predicted_values)
```

Arguments

`true_values` factor containing the true labels/categories.

`predicted_values` factor containing the predicted labels/categories.

Value

Returns a matrix which contains the cases categories in the rows and the measures (precision, recall, f1) in the columns.

See Also

Other performance measures: [cohens_kappa\(\)](#), [fleiss_kappa\(\)](#), [get_coder_metrics\(\)](#), [gwet_ac\(\)](#), [kendalls_w\(\)](#), [kripp_alpha\(\)](#)

calc_tokenizer_statistics

Estimate tokenizer statistics

Description

Function for estimating the tokenizer statistics described by Kaya & Tantug (2024).

Usage

```
calc_tokenizer_statistics(dataset, step = "creation")
```

Arguments

- | | |
|---------|--|
| dataset | Object of class datasets.arrow_dataset.Dataset. The data set must contain a column "length" containing the number of tokens for every sequence and a column "word_ids" containing the word ids within every sequence. |
| step | string indicating to which step the statistics belong. Recommended values are <ul style="list-style-type: none"> • "creation" for the creation of the tokenizer. • "initial_training" for the first training of the transformer. • "fine_tuning" for all following trainings of the transformer. • "training" for a training run of the transformer. |

Value

Returns a list with the following entries:

- n_sequences: Number of sequences
- n_words: Number for words in whole corpus
- n_tokens: Number of tokens in the whole corpus
- mu_t: eqn(n_tokens/n_sequences)
- mu_w: eqn(n_words/n_sequences)
- mu_g: eqn(n_tokens/n_words)

References

Kaya, Y. B., & Tantug, A. C. (2024). Effect of tokenization granularity for Turkish large language models. Intelligent Systems with Applications, 21, 200335. <https://doi.org/10.1016/j.iswa.2024.200335>

cat_message	<i>Print message (cat())</i>
-------------	------------------------------

Description

Prints a message `msg` if `trace` parameter is `TRUE` with current date with `cat()` function.

Usage

```
cat_message(msg, trace)
```

Arguments

msg	string Message that should be printed.
trace	bool Silent printing (<code>FALSE</code>) or not (<code>TRUE</code>).

Value

This function returns nothing.

See Also

Other Utils Log Developers: [clean_pytorch_log_transformers\(\)](#), [output_message\(\)](#), [print_message\(\)](#), [read_log\(\)](#), [read_loss_log\(\)](#), [reset_log\(\)](#), [reset_loss_log\(\)](#), [write_log\(\)](#)

check_adjust_n_samples_on_CI	<i>Set sample size for argument combinations</i>
------------------------------	--

Description

Function adjust the number of samples depending on the test environment. On continuous integration it is limited to a random sample of combinations.

Usage

```
check_adjust_n_samples_on_CI(n_samples_requested, n_CI = 50)
```

Arguments

n_samples_requested	int Number of samples if the test do not run on continuous integration.
n_CI	int Number of samples if the test run on continuous integration.

Value

Returns an `int` depending on the test environment.

See Also

Other Utils TestThat Developers: [generate_args_for_tests\(\)](#), [generate_embeddings\(\)](#), [generate_tensors\(\)](#), [get_current_args_for_print\(\)](#), [get_fixed_test_tensor\(\)](#), [get_test_data_for_classifiers\(\)](#), [random_bool_on_CI\(\)](#)

`check_aif_py_modules` *Check if all necessary python modules are available*

Description

This function checks if all python modules necessary for the package 'aifedducation' to work are available.

Usage

```
check_aif_py_modules(trace = TRUE)
```

Arguments

<code>trace</code>	bool <code>TRUE</code> if a list with all modules and their availability should be printed to the console.
--------------------	--

Value

The function prints a table with all relevant packages and shows which modules are available or unavailable.

If all relevant modules are available, the functions returns `TRUE`. In all other cases it returns `FALSE`

See Also

Other Installation and Configuration: [install_aifedducation\(\)](#), [install_aifedducation_studio\(\)](#), [install_py_modules\(\)](#), [prepare_session\(\)](#), [set_transformers_logger\(\)](#), [update_aifedducation\(\)](#)

check_all_args	<i>Check arguments automatically</i>
----------------	--------------------------------------

Description

This function performs checks for every provided argument. It can only check arguments that are defined in the central parameter dictionary. See [get_param_dict](#) for more details.

Usage

```
check_all_args(args)
```

Arguments

args	Named list containing the arguments and their values.
------	---

Value

Function does nothing return. It raises an error the arguments are not valid.

See Also

Other Utils Checks Developers: [check_class_and_type\(\)](#)

check_class_and_type	<i>Check class and type</i>
----------------------	-----------------------------

Description

Function for checking if an object is of a specific type or class.

Usage

```
check_class_and_type(  
    object,  
    object_name = NULL,  
    type_classes = "bool",  
    allow_NULL = FALSE,  
    min = NULL,  
    max = NULL,  
    allowed_values = NULL  
)
```

Arguments

object	Any R object.
object_name	string Name of the object. This is helpful for debugging.
type_classes	vector of strings containing the type or classes which the object should belong to.
allow_NULL	bool If TRUE allow the object to be NULL.
min	double or int Minimal value for the object.
max	double or int Maximal value for the object.
allowed_values	vector of strings determining the allowed values. If all strings are allowed set this argument to NULL.

Value

Function does nothing return. It raises an error if the object is not of the specified type.

Note

parameter min, max, and allowed_values do not apply if type_classes is a class.
allowed_values does only apply if type_classes is string.

See Also

Other Utils Checks Developers: [check_all_args\(\)](#)

ClassifiersBasedOnTextEmbeddings

Abstract class for all classifiers that use numerical representations of texts instead of words.

Description

Base class for classifiers relying on [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) generated with a [TextEmbeddingModel](#).

Objects of this class containing fields and methods used in several other classes in 'AI for Education'.

This class is **not** designed for a direct application and should only be used by developers.

Value

A new object of this class.

Super classes

[aifedducation::AIFEBaseModel](#) -> [aifedducation::ModelsBasedOnTextEmbeddings](#) -> ClassifiersBasedOnTextEmbe

Public fields

`feature_extractor ('list()')`

List for storing information and objects about the feature_extractor.

`reliability ('list()')`

List for storing central reliability measures of the last training.

- `reliability$test_metric`: Array containing the reliability measures for the test data for every fold and step (in case of pseudo-labeling).
- `reliability$test_metric_mean`: Array containing the reliability measures for the test data. The values represent the mean values for every fold.
- `reliability$raw_iota_objects`: List containing all iota_object generated with the package iotarelr for every fold at the end of the last training for the test data.
- `reliability$raw_iota_objects$iota_objects_end`: List of objects with class iotarelr_iota2 containing the estimated iota reliability of the second generation for the final model for every fold for the test data.
- `reliability$raw_iota_objects$iota_objects_end_free`: List of objects with class iotarelr_iota2 containing the estimated iota reliability of the second generation for the final model for every fold for the test data. Please note that the model is estimated without forcing the Assignment Error Matrix to be in line with the assumption of weak superiority.
- `reliability$iota_object_end`: Object of class iotarelr_iota2 as a mean of the individual objects for every fold for the test data.
- `reliability$iota_object_end_free`: Object of class iotarelr_iota2 as a mean of the individual objects for every fold. Please note that the model is estimated without forcing the Assignment Error Matrix to be in line with the assumption of weak superiority.
- `reliability$standard_measures_end`: Object of class list containing the final measures for precision, recall, and f1 for every fold.
- `reliability$standard_measures_mean`: matrix containing the mean measures for precision, recall, and f1.

Methods

Public methods:

- `ClassifiersBasedOnTextEmbeddings$predict()`
- `ClassifiersBasedOnTextEmbeddings$check_embedding_model()`
- `ClassifiersBasedOnTextEmbeddings$check_feature_extractor_object_type()`
- `ClassifiersBasedOnTextEmbeddings$requires_compression()`
- `ClassifiersBasedOnTextEmbeddings$save()`
- `ClassifiersBasedOnTextEmbeddings$load_from_disk()`
- `ClassifiersBasedOnTextEmbeddings$adjust_target_levels()`
- `ClassifiersBasedOnTextEmbeddings$plot_training_history()`
- `ClassifiersBasedOnTextEmbeddings$plot_coding_stream()`
- `ClassifiersBasedOnTextEmbeddings$clone()`

Method predict(): Method for predicting new data with a trained neural net.

Usage:

```
ClassifiersBasedOnTextEmbeddings$predict(
  newdata,
  batch_size = 32,
  ml_trace = 1
)
```

Arguments:

newdata Object of class [TextEmbeddingModel](#) or [LargeDataSetForTextEmbeddings](#) for which predictions should be made. In addition, this method allows to use objects of class [array](#) and [datasets.arrow_dataset.Dataset](#). However, these should be used only by developers.
batch_size int Size of batches.
ml_trace int `ml_trace=0` does not print any information on the process from the machine learning framework.

Returns: Returns a `data.frame` containing the predictions and the probabilities of the different labels for each case.

Method `check_embedding_model()`: Method for checking if the provided text embeddings are created with the same [TextEmbeddingModel](#) as the classifier.

Usage:

```
ClassifiersBasedOnTextEmbeddings$check_embedding_model(
  text_embeddings,
  require_compressed = FALSE
)
```

Arguments:

text_embeddings Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#).
require_compressed TRUE if a compressed version of the embeddings are necessary. Compressed embeddings are created by an object of class [TEFeatureExtractor](#).

Returns: TRUE if the underlying [TextEmbeddingModel](#) is the same. FALSE if the models differ.

Method `check_feature_extractor_object_type()`: Method for checking an object of class [TEFeatureExtractor](#).

Usage:

```
ClassifiersBasedOnTextEmbeddings$check_feature_extractor_object_type(
  feature_extractor
)
```

Arguments:

feature_extractor Object of class [TEFeatureExtractor](#)

Returns: This method does nothing returns. It raises an error if

- the object is NULL
- the object does not rely on the same machine learning framework as the classifier
- the object is not trained.

Method `requires_compression()`: Method for checking if provided text embeddings must be compressed via a [TEFeatureExtractor](#) before processing.

Usage:

```
ClassifiersBasedOnTextEmbeddings$requires_compression(text_embeddings)
```

Arguments:

text_embeddings Object of class [EmbeddedText](#), [LargeDataSetForTextEmbeddings](#), array or `datasets.arrow_dataset.Dataset`.

Returns: Return TRUE if a compression is necessary and FALSE if not.

Method `save()`: Method for saving a model.

Usage:

```
ClassifiersBasedOnTextEmbeddings$save(dir_path, folder_name)
```

Arguments:

dir_path string Path of the directory where the model should be saved.

folder_name string Name of the folder that should be created within the directory.

Returns: Function does not return a value. It saves the model to disk.

Method `load_from_disk()`: loads an object from disk and updates the object to the current version of the package.

Usage:

```
ClassifiersBasedOnTextEmbeddings$load_from_disk(dir_path)
```

Arguments:

dir_path Path where the object set is stored.

Returns: Method does not return anything. It loads an object from disk.

Method `adjust_target_levels()`: Method transforms the levels of a factor into numbers corresponding to the models definition.

Usage:

```
ClassifiersBasedOnTextEmbeddings$adjust_target_levels(data_targets)
```

Arguments:

data_targets factor containing the labels for cases stored in embeddings. Factor must be named and has to use the same names as used in in the embeddings.

Returns: Method returns a factor containing the numerical representation of categories/classes.

Method `plot_training_history()`: Method for requesting a plot of the training history. This method requires the R package 'ggplot2' to work.

Usage:

```
ClassifiersBasedOnTextEmbeddings$plot_training_history(
  final_training = FALSE,
  pl_step = NULL,
  measure = "loss",
  y_min = NULL,
  y_max = NULL,
  add_min_max = TRUE,
  text_size = 10
)
```

Arguments:

`final_training` bool If FALSE the values of the performance estimation are used. If TRUE only the epochs of the final training are used.

`pl_step` int Number of the step during pseudo labeling to plot. Only relevant if the model was trained with active pseudo labeling.

`measure` string Measure to plot. Allowed values:

- "avg_iota" = Average Iota
- "loss" = Loss
- "accuracy" = Accuracy
- "balanced_accuracy" = Balanced Accuracy

`y_min` Minimal value for the y-axis. Set to NULL for an automatic adjustment.

`y_max` Maximal value for the y-axis. Set to NULL for an automatic adjustment.

`add_min_max` bool If TRUE the minimal and maximal values during performance estimation are part of the plot. If FALSE only the mean values are shown. Parameter is ignored if `final_training=TRUE`.

`text_size` Size of the text.

Returns: Returns a plot of class ggplot visualizing the training process.

Method `plot_coding_stream()`: Method for requesting a plot the coding stream. The plot shows how the cases of different categories/classes are assigned to a the available classes/categories. The visualization is helpful for analyzing the consequences of coding errors.

Usage:

```
ClassifiersBasedOnTextEmbeddings$plot_coding_stream(
  label_categories_size = 3,
  key_size = 0.5,
  text_size = 10
)
```

Arguments:

`label_categories_size` double determining the size of the label for each true and assigned category within the plot.

`key_size` double determining the size of the legend.

`text_size` double determining the size of the text within the legend.

Returns: Returns a plot of class ggplot visualizing the training process.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
ClassifiersBasedOnTextEmbeddings$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

See Also

Other R6 Classes for Developers: [AIFEBaseModel](#), [LargeDataSetBase](#), [ModelsBasedOnTextEmbeddings](#), [TEClassifiersBasedOnProtoNet](#), [TEClassifiersBasedOnRegular](#)

class_vector_to_py_dataset

Convert class vector to arrow data set

Description

Function converts a vector of class indices into an arrow data set.

Usage

```
class_vector_to_py_dataset(vector)
```

Arguments

vector vector of class indices.

Value

Returns a data set of class datasets.arrow_dataset.Dataset containing the class indices.

See Also

Other Utils Python Data Management Developers: [data.frame_to_py_dataset\(\)](#), [get_batches_index\(\)](#), [prepare_r_array_for_dataset\(\)](#), [py_dataset_to_embeddings\(\)](#), [reduce_to_unique\(\)](#), [tensor_list_to_numpy\(\)](#), [tensor_to_numpy\(\)](#)

clean_pytorch_log_transformers

Clean pytorch log of transformers

Description

Function for preparing and cleaning the log created by an object of class Trainer from the python library 'transformer's.

Usage

```
clean_pytorch_log_transformers(log)
```

Arguments

log data.frame containing the log.

Value

Returns a data.frame containing epochs, loss, and val_loss.

See Also

Other Utils Log Developers: [cat_message\(\)](#), [output_message\(\)](#), [print_message\(\)](#), [read_log\(\)](#), [read_loss_log\(\)](#), [reset_log\(\)](#), [reset_loss_log\(\)](#), [write_log\(\)](#)

cohens_kappa	<i>Calculate Cohen's Kappa</i>
--------------	--------------------------------

Description

This function calculates different version of Cohen's Kappa.

Usage

```
cohens_kappa(rater_one, rater_two)
```

Arguments

rater_one	factor rating of the first coder.
rater_two	factor ratings of the second coder.

Value

Returns a list containing the results for Cohen' Kappa if no weights are applied (`kappa_unweighted`), if weights are applied and the weights increase linear (`kappa_linear`), and if weights are applied and the weights increase quadratic (`kappa_squared`).

References

- Cohen, J (1968). Weighted kappa: Nominal scale agreement with provision for scaled disagreement or partial credit. *Psychological Bulletin*, 70(4), 213–220. [doi:10.1037/h0026256](https://doi.org/10.1037/h0026256)
- Cohen, J (1960). A Coefficient of Agreement for Nominal Scales. *Educational and Psychological Measurement*, 20(1), 37–46. [doi:10.1177/001316446002000104](https://doi.org/10.1177/001316446002000104)

See Also

Other performance measures: [calc_standard_classification_measures\(\)](#), [fleiss_kappa\(\)](#), [get_coder_metrics\(\)](#), [gwet_ac\(\)](#), [kendalls_w\(\)](#), [kripp_alpha\(\)](#)

create_dir	<i>Create directory if not exists</i>
------------	---------------------------------------

Description

Check whether the passed dir_path directory exists. If not, creates a new directory and prints a msg message if trace is TRUE.

Usage

```
create_dir(dir_path, trace, msg = "Creating Directory", msg_fun = TRUE)
```

Arguments

dir_path	string A new directory path that should be created.
trace	bool Whether a msg message should be printed.
msg	string A message that should be printed if trace is TRUE.
msg_fun	func Function used for printing the message.

Value

TRUE or FALSE depending on whether the shiny app is active.

See Also

Other Utils File Management Developers: [get_file_extension\(\)](#)

create_object	<i>Create object</i>
---------------	----------------------

Description

Support function for creating objects.

Usage

```
create_object(class)
```

Arguments

class	string Name of the class to be created.#'
-------	---

Value

Returns an object of the requested class.#'

See Also

Other Utils Developers: [auto_n_cores\(\)](#), [create_synthetic_units_from_matrix\(\)](#), [generate_id\(\)](#), [get_n_chunks\(\)](#), [get_synthetic_cases_from_matrix\(\)](#), [matrix_to_array_c\(\)](#), [tensor_to_matrix_c\(\)](#), [to_categorical_c\(\)](#)

create_synthetic_units_from_matrix
Create synthetic units

Description

Function for creating synthetic cases in order to balance the data for training with [TEClassifierRegulator](#) or [TEClassifierProtoNet](#). This is an auxiliary function for use with [get_synthetic_cases_from_matrix](#) to allow parallel computations.

Usage

```
create_synthetic_units_from_matrix(
  matrix_form,
  target,
  required_cases,
  k,
  method,
  cat,
  k_s,
  max_k
)
```

Arguments

<code>matrix_form</code>	Named <code>matrix</code> containing the text embeddings in matrix form. In most cases this object is taken from EmbeddedText\$embeddings .
<code>target</code>	Named factor containing the labels/categories of the corresponding cases.
<code>required_cases</code>	<code>int</code> Number of cases necessary to fill the gap between the frequency of the class under investigation and the major class.
<code>k</code>	<code>int</code> The number of nearest neighbors during sampling process.
<code>method</code>	vector containing strings of the requested methods for generating new cases. Currently "knnor" from this package is available.
<code>cat</code>	<code>string</code> The category for which new cases should be created.
<code>k_s</code>	<code>int</code> Number of ks in the complete generation process.
<code>max_k</code>	<code>int</code> The maximum number of nearest neighbors during sampling process.

Value

Returns a `list` which contains the text embeddings of the new synthetic cases as a named `data.frame` and their labels as a named factor.

See Also

Other Utils Developers: [auto_n_cores\(\)](#), [create_object\(\)](#), [generate_id\(\)](#), [get_n_chunks\(\)](#), [get_synthetic_cases_from_matrix\(\)](#), [matrix_to_array_c\(\)](#), [tensor_to_matrix_c\(\)](#), [to_categorical_c\(\)](#)

data.frame_to_py_dataset

Convert data.frame to arrow data set

Description

Function for converting a data.frame into a pyarrow data set.

Usage

```
data.frame_to_py_dataset(data_frame)
```

Arguments

data_frame Object of class data.frame.

Value

Returns the data.frame as a pyarrow data set of class datasets.arrow_dataset.Dataset.

See Also

Other Utils Python Data Management Developers: [class_vector_to_py_dataset\(\)](#), [get_batches_index\(\)](#), [prepare_r_array_for_dataset\(\)](#), [py_dataset_to_embeddings\(\)](#), [reduce_to_unique\(\)](#), [tensor_list_to_numpy\(\)](#), [tensor_to_numpy\(\)](#)

DataManagerClassifier *Data manager for classification tasks*

Description

Abstract class for managing the data and samples during training a classifier. DataManagerClassifier is used with all classifiers based on text embeddings.

Value

Objects of this class are used for ensuring the correct data management for training different types of classifiers. They are also used for data augmentation by creating synthetic cases with different techniques.

Public fields

`config ('list')`
 Field for storing configuration of the [DataManagerClassifier](#).

`state ('list')`
 Field for storing the current state of the [DataManagerClassifier](#).

`datasets ('list')`
 Field for storing the data sets used during training. All elements of the list are data sets of class `datasets.arrow_dataset.Dataset`. The following data sets are available:

- `data_labeled`: all cases which have a label.
- `data_unlabeled`: all cases which have no label.
- `data_labeled_synthetic`: all synthetic cases with their corresponding labels.
- `data_labeled_pseudo`: subset of `data_unlabeled` if pseudo labels were estimated by a classifier.

`name_idx ('named vector')`
 Field for storing the pairs of indexes and names of every case. The pairs for labeled and unlabeled data are separated.

`samples ('list')`
 Field for storing the assignment of every cases to a train, validation or test data set depending on the concrete fold. Only the indexes and not the names are stored. In addition, the list contains the assignment for the final training which excludes a test data set. If the [DataManagerClassifier](#) uses `i` folds the sample for the final training can be requested with `i+1`.

Methods

Public methods:

- [DataManagerClassifier\\$new\(\)](#)
- [DataManagerClassifier\\$get_config\(\)](#)
- [DataManagerClassifier\\$get_labeled_data\(\)](#)
- [DataManagerClassifier\\$get_unlabeled_data\(\)](#)
- [DataManagerClassifier\\$get_samples\(\)](#)
- [DataManagerClassifier\\$set_state\(\)](#)
- [DataManagerClassifier\\$get_n_folds\(\)](#)
- [DataManagerClassifier\\$get_n_classes\(\)](#)
- [DataManagerClassifier\\$get_statistics\(\)](#)
- [DataManagerClassifier\\$contains_unlabeled_data\(\)](#)
- [DataManagerClassifier\\$get_dataset\(\)](#)
- [DataManagerClassifier\\$get_val_dataset\(\)](#)
- [DataManagerClassifier\\$get_test_dataset\(\)](#)
- [DataManagerClassifier\\$create_synthetic\(\)](#)
- [DataManagerClassifier\\$add_replace_pseudo_data\(\)](#)
- [DataManagerClassifier\\$clone\(\)](#)

Method `new()`: Creating a new instance of this class.

Usage:

```
DataManagerClassifier$new(
  data_embeddings,
  data_targets,
  folds = 5,
  val_size = 0.25,
  pad_value = -100,
  class_levels,
  one_hot_encoding = TRUE,
  add_matrix_map = TRUE,
  sc_methods = "knnor",
  sc_min_k = 1,
  sc_max_k = 10,
  trace = TRUE,
  n_cores = auto_n_cores()
)
```

Arguments:

`data_embeddings` EmbeddedText, LargeDataSetForTextEmbeddings Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#).

`data_targets` factor containing the labels for cases stored in embeddings. Factor must be named and has to use the same names as used in the embeddings.

`folds` int determining the number of cross-fold samples. Allowed values: $1 \leq x$

`val_size` double between 0 and 1, indicating the proportion of cases which should be used for the validation sample during the estimation of the model. The remaining cases are part of the training data. Allowed values: $0 < x < 1$

`pad_value` int Value indicating padding. This value should not be in the range of regular values for computations. Thus it is not recommended to change this value. Default is -100. Allowed values: $x \leq -100$

`class_levels` vector containing the levels (categories or classes) within the target data. Please note that order matters. For ordinal data please ensure that the levels are sorted correctly with later levels indicating a higher category/class. For nominal data the order does not matter.

`one_hot_encoding` bool If TRUE all labels are converted to one hot encoding.

`add_matrix_map` bool If TRUE all embeddings are transformed into a two dimensional matrix. The number of rows equals the number of cases. The number of columns equals times*features.

`sc_methods` string containing the method for generating synthetic cases. Allowed values: 'knnor'

`sc_min_k` int determining the minimal number of k which is used for creating synthetic units. Allowed values: $1 \leq x$

`sc_max_k` int determining the maximal number of k which is used for creating synthetic units. Allowed values: $1 \leq x$

`trace` bool TRUE if information about the estimation phase should be printed to the console.

`n_cores` int Number of cores which should be used during the calculation of synthetic cases. Only relevant if `use_sc=TRUE`. Allowed values: $1 \leq x$

Returns: Method returns an initialized object of class [DataManagerClassifier](#).

Method `get_config()`: Method for requesting the configuration of the [DataManagerClassifier](#).

Usage:

```
DataManagerClassifier$get_config()
```

Returns: Returns a list storing the configuration of the [DataManagerClassifier](#).

Method `get_labeled_data()`: Method for requesting the complete labeled data set.

Usage:

```
DataManagerClassifier$get_labeled_data()
```

Returns: Returns an object of class `datasets.arrow_dataset.Dataset` containing all cases with labels.

Method `get_unlabeled_data()`: Method for requesting the complete unlabeled data set.

Usage:

```
DataManagerClassifier$get_unlabeled_data()
```

Returns: Returns an object of class `datasets.arrow_dataset.Dataset` containing all cases without labels.

Method `get_samples()`: Method for requesting the assignments to train, validation, and test data sets for every fold and the final training.

Usage:

```
DataManagerClassifier$get_samples()
```

Returns: Returns a list storing the assignments to a train, validation, and test data set for every fold. In the case of the sample for the final training the test data set is always empty (NULL).

Method `set_state()`: Method for setting the current state of the [DataManagerClassifier](#).

Usage:

```
DataManagerClassifier$set_state(iteration, step = NULL)
```

Arguments:

`iteration` int determining the current iteration of the training. That is iteration determines the fold to use for training, validation, and testing. If i is the number of fold $i+1$ request the sample for the final training. For requesting the sample for the final training iteration can take a string "final".

`step` int determining the step for estimating and using pseudo labels during training. Only relevant if training is requested with pseudo labels.

Returns: Method does not return anything. It is used for setting the internal state of the [DataManager](#).

Method `get_n_folds()`: Method for requesting the number of folds the [DataManagerClassifier](#) can use with the current data.

Usage:

```
DataManagerClassifier$get_n_folds()
```

Returns: Returns the number of folds the [DataManagerClassifier](#) uses.

Method `get_n_classes()`: Method for requesting the number of classes.

Usage:

```
DataManagerClassifier$get_n_classes()
```

Returns: Returns the number classes.

Method `get_statistics()`: Method for requesting descriptive sample statistics.

Usage:

```
DataManagerClassifier$get_statistics()
```

Returns: Returns a table describing the absolute frequencies of the labeled and unlabeled data. The rows contain the length of the sequences while the columns contain the labels.

Method `contains_unlabeled_data()`: Method for checking if the dataset contains cases without labels.

Usage:

```
DataManagerClassifier$contains_unlabeled_data()
```

Returns: Returns TRUE if the dataset contains cases without labels. Returns FALSE if all cases have labels.

Method `get_dataset()`: Method for requesting a data set for training depending in the current state of the `DataManagerClassifier`.

Usage:

```
DataManagerClassifier$get_dataset(  
  inc_labeled = TRUE,  
  inc_unlabeled = FALSE,  
  inc_synthetic = FALSE,  
  inc_pseudo_data = FALSE  
)
```

Arguments:

`inc_labeled` bool If TRUE the data set includes all cases which have labels.

`inc_unlabeled` bool If TRUE the data set includes all cases which have no labels.

`inc_synthetic` bool If TRUE the data set includes all synthetic cases with their corresponding labels.

`inc_pseudo_data` bool If TRUE the data set includes all cases which have pseudo labels.

Returns: Returns an object of class `datasets.arrow_dataset.Dataset` containing the requested kind of data along with all requested transformations for training. Please note that this method returns a data sets that is designed for training only. The corresponding validation data set is requested with `get_val_dataset` and the corresponding test data set with `get_test_dataset`.

Method `get_val_dataset()`: Method for requesting a data set for validation depending in the current state of the `DataManagerClassifier`.

Usage:

```
DataManagerClassifier$get_val_dataset()
```

Returns: Returns an object of class `datasets.arrow_dataset.Dataset` containing the requested kind of data along with all requested transformations for validation. The corresponding data set for training can be requested with `get_dataset` and the corresponding data set for testing with `get_test_dataset`.

Method `get_test_dataset()`: Method for requesting a data set for testing depending in the current state of the DataManagerClassifier.

Usage:

```
DataManagerClassifier$get_test_dataset()
```

Returns: Returns an object of class `datasets.arrow_dataset.Dataset` containing the requested kind of data along with all requested transformations for validation. The corresponding data set for training can be requested with `get_dataset` and the corresponding data set for validation with `get_val_dataset`.

Method `create_synthetic()`: Method for generating synthetic data used during training. The process uses all labeled data belonging to the current state of the DataManagerClassifier.

Usage:

```
DataManagerClassifier$create_synthetic(trace = TRUE, inc_pseudo_data = FALSE)
```

Arguments:

`trace` bool If TRUE information on the process are printed to the console.

`inc_pseudo_data` bool If TRUE data with pseudo labels are used in addition to the labeled data for generating synthetic cases.

Returns: This method does nothing return. It generates a new data set for synthetic cases which are stored as an object of class `datasets.arrow_dataset.Dataset` in the field `datasets$data_labeled_synthetic`. Please note that a call of this method will override an existing data set in the corresponding field.

Method `add_replace_pseudo_data()`: Method for adding data with pseudo labels generated by a classifier

Usage:

```
DataManagerClassifier$add_replace_pseudo_data(inputs, labels)
```

Arguments:

`inputs` array or matrix representing the input data.

`labels` factor containing the corresponding pseudo labels.

Returns: This method does nothing return. It generates a new data set for synthetic cases which are stored as an object of class `datasets.arrow_dataset.Dataset` in the field `datasets$data_labeled_pseudo`. Please note that a call of this method will override an existing data set in the corresponding field.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
DataManagerClassifier$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

EmbeddedText*Abstract class for small data sets containing text embeddings*

Description

Object of class R6 which stores the text embeddings generated by an object of class [TextEmbeddingModel](#). The text embeddings are stored within memory/RAM. In the case of a high number of documents the data may not fit into memory/RAM. Thus, please use this object only for a small sample of texts. In general, it is recommended to use an object of class [LargeDataSetForTextEmbeddings](#) which can deal with any number of texts.

Value

Returns an object of class [EmbeddedText](#). These objects are used for storing and managing the text embeddings created with objects of class [TextEmbeddingModel](#). Objects of class [EmbeddedText](#) serve as input for objects of class [TEClassifierRegular](#), [TEClassifierProtoNet](#), and [TEFeatureExtractor](#). The main aim of this class is to provide a structured link between embedding models and classifiers. Since objects of this class save information on the text embedding model that created the text embedding it ensures that only embedding generated with same embedding model are combined. Furthermore, the stored information allows objects to check if embeddings of the correct text embedding model are used for training and predicting.

Public fields

`embeddings ('data.frame()')`

data.frame containing the text embeddings for all chunks. Documents are in the rows. Embedding dimensions are in the columns.

Methods**Public methods:**

- [EmbeddedText\\$configure\(\)](#)
- [EmbeddedText\\$save\(\)](#)
- [EmbeddedText\\$is_configured\(\)](#)
- [EmbeddedText\\$load_from_disk\(\)](#)
- [EmbeddedText\\$get_model_info\(\)](#)
- [EmbeddedText\\$get_model_label\(\)](#)
- [EmbeddedText\\$get_times\(\)](#)
- [EmbeddedText\\$get_features\(\)](#)
- [EmbeddedText\\$get_original_features\(\)](#)
- [EmbeddedText\\$get_pad_value\(\)](#)
- [EmbeddedText\\$is_compressed\(\)](#)
- [EmbeddedText\\$add_feature_extractor_info\(\)](#)
- [EmbeddedText\\$get_feature_extractor_info\(\)](#)
- [EmbeddedText\\$convert_to_LargeDataSetForTextEmbeddings\(\)](#)

- `EmbeddedText$n_rows()`
- `EmbeddedText$get_all_fields()`
- `EmbeddedText$set_package_versions()`
- `EmbeddedText$get_package_versions()`
- `EmbeddedText$clone()`

Method `configure()`: Creates a new object representing text embeddings.

Usage:

```
EmbeddedText$configure(
  model_name = NA,
  model_label = NA,
  model_date = NA,
  model_method = NA,
  model_version = NA,
  model_language = NA,
  param_seq_length = NA,
  param_chunks = NULL,
  param_features = NULL,
  param_overlap = NULL,
  param_emb_layer_min = NULL,
  param_emb_layer_max = NULL,
  param_emb_pool_type = NULL,
  param_aggregation = NULL,
  param_pad_value = -100,
  embeddings
)
```

Arguments:

`model_name` `string` Name of the model that generates this embedding.
`model_label` `string` Label of the model that generates this embedding.
`model_date` `string` Date when the embedding generating model was created.
`model_method` `string` Method of the underlying embedding model.
`model_version` `string` Version of the model that generated this embedding.
`model_language` `string` Language of the model that generated this embedding.
`param_seq_length` `int` Maximum number of tokens that processes the generating model for a chunk.
`param_chunks` `int` Maximum number of chunks which are supported by the generating model.
`param_features` `int` Number of dimensions of the text embeddings.
`param_overlap` `int` Number of tokens that were added at the beginning of the sequence for the next chunk by this model. #'
`param_emb_layer_min` `int or string` determining the first layer to be included in the creation of embeddings.
`param_emb_layer_max` `int or string` determining the last layer to be included in the creation of embeddings.
`param_emb_pool_type` `string` determining the method for pooling the token embeddings within each layer.

`param_aggregation string` Aggregation method of the hidden states. Deprecated. Only included for backward compatibility.

`param_pad_value int` Value indicating padding. This value should no be in the range of regular values for computations. Thus it is not recommended to chance this value. Default is -100. Allowed values: $x \leq -100$

`embeddings data.frame` containing the text embeddings.

Returns: Returns an object of class `EmbeddedText` which stores the text embeddings produced by an objects of class `TextEmbeddingModel`.

Method `save()`: Saves a data set to disk.

Usage:

`EmbeddedText$save(dir_path, folder_name, create_dir = TRUE)`

Arguments:

`dir_path` Path where to store the data set.

`folder_name string` Name of the folder for storing the data set.

`create_dir bool` If True the directory will be created if it does not exist.

Returns: Method does not return anything. It write the data set to disk.

Method `is_configured()`: Method for checking if the model was successfully configured. An object can only be used if this value is TRUE.

Usage:

`EmbeddedText$is_configured()`

Returns: bool TRUE if the model is fully configured. FALSE if not.

Method `load_from_disk()`: loads an object of class `EmbeddedText` from disk and updates the object to the current version of the package.

Usage:

`EmbeddedText$load_from_disk(dir_path)`

Arguments:

`dir_path` Path where the data set set is stored.

Returns: Method does not return anything. It loads an object from disk.

Method `get_model_info()`: Method for retrieving information about the model that generated this embedding.

Usage:

`EmbeddedText$get_model_info()`

Returns: list contains all saved information about the underlying text embedding model.

Method `get_model_label()`: Method for retrieving the label of the model that generated this embedding.

Usage:

`EmbeddedText$get_model_label()`

Returns: string Label of the corresponding text embedding model

Method `get_times()`: Number of chunks/times of the text embeddings.

Usage:

```
EmbeddedText$get_times()
```

Returns: Returns an int describing the number of chunks/times of the text embeddings.

Method `get_features()`: Number of actual features/dimensions of the text embeddings. In the case a [feature extractor](#) was used the number of features is smaller as the original number of features. To receive the original number of features (the number of features before applying a [feature extractor](#)) you can use the method `get_original_features` of this class.

Usage:

```
EmbeddedText$get_features()
```

Returns: Returns an int describing the number of features/dimensions of the text embeddings.

Method `get_original_features()`: Number of original features/dimensions of the text embeddings.

Usage:

```
EmbeddedText$get_original_features()
```

Returns: Returns an int describing the number of features/dimensions if no [feature extractor](#) is used or before a [feature extractor](#)) is applied.

Method `get_pad_value()`: Value for indicating padding.

Usage:

```
EmbeddedText$get_pad_value()
```

Returns: Returns an int describing the value used for padding.

Method `is_compressed()`: Checks if the text embedding were reduced by a [feature extractor](#).

Usage:

```
EmbeddedText$is_compressed()
```

Returns: Returns TRUE if the number of dimensions was reduced by a [feature extractor](#). If not return FALSE.

Method `add_feature_extractor_info()`: Method setting information on the [feature extractor](#) that was used to reduce the number of dimensions of the text embeddings. This information should only be used if a [feature extractor](#) was applied.

Usage:

```
EmbeddedText$add_feature_extractor_info(
  model_name,
  model_label = NA,
  features = NA,
  method = NA,
  noise_factor = NA,
  optimizer = NA
)
```

Arguments:

```

model_name string Name of the underlying TextEmbeddingModel.
model_label string Label of the underlying TextEmbeddingModel.
features int Number of dimension (features) for the compressed text embeddings.
method string Method that the TEFeatureExtractor applies for generating the compressed
text embeddings.
noise_factor double Noise factor of the TEFeatureExtractor.
optimizer string Optimizer used during training the TEFeatureExtractor.

```

Returns: Method does nothing return. It sets information on a [feature extractor](#).

Method `get_feature_extractor_info()`: Method for receiving information on the [feature extractor](#) that was used to reduce the number of dimensions of the text embeddings.

Usage:

```
EmbeddedText$get_feature_extractor_info()
```

Returns: Returns a list with information on the [feature extractor](#). If no [feature extractor](#) was used it returns NULL.

Method `convert_to_LargeDataSetForTextEmbeddings()`: Method for converting this object to an object of class [LargeDataSetForTextEmbeddings](#).

Usage:

```
EmbeddedText$convert_to_LargeDataSetForTextEmbeddings()
```

Returns: Returns an object of class [LargeDataSetForTextEmbeddings](#) which uses memory mapping allowing to work with large data sets.

Method `n_rows()`: Number of rows.

Usage:

```
EmbeddedText$n_rows()
```

Returns: Returns the number of rows of the text embeddings which represent the number of cases.

Method `get_all_fields()`: Return all fields.

Usage:

```
EmbeddedText$get_all_fields()
```

Returns: Method returns a list containing all public and private fields of the object.

Method `set_package_versions()`: Method for setting the package version for 'aifeducation', 'reticulate', 'torch', and 'numpy' to the currently used versions.

Usage:

```
EmbeddedText$set_package_versions()
```

Returns: Method does not return anything. It is used to set the private fields for package versions.

Method `get_package_versions()`: Method for requesting a summary of the R and python packages' versions used for creating the model.

Usage:

```
EmbeddedText$get_package_versions()
```

Returns: Returns a list containing the versions of the relevant R and python packages.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
EmbeddedText$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

See Also

Other Data Management: [LargeDataSetForText](#), [LargeDataSetForTextEmbeddings](#)

fleiss_kappa

Calculate Fleiss' Kappa

Description

This function calculates Fleiss' Kappa.

Usage

```
fleiss_kappa(rater_one, rater_two, additional_raters = NULL)
```

Arguments

<code>rater_one</code>	factor rating of the first coder.
<code>rater_two</code>	factor ratings of the second coder.
<code>additional_raters</code>	<code>list</code> Additional raters with same requirements as <code>rater_one</code> and <code>rater_two</code> . If there are no additional raters set to <code>NULL</code> .

Value

Returns the value for Fleiss' Kappa.

References

Fleiss, J. L. (1971). Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76(5), 378–382. [doi:10.1037/h0031619](https://doi.org/10.1037/h0031619)

See Also

Other performance measures: [calc_standard_classification_measures\(\)](#), [cohens_kappa\(\)](#), [get_coder_metrics\(\)](#), [gweet_ac\(\)](#), [kendalls_w\(\)](#), [kripp_alpha\(\)](#)

generate_args_for_tests

Generate combinations of arguments

Description

Function generates a specific number of combinations for a method. These are used for automating tests of objects.

Usage

```
generate_args_for_tests(  
    object_name,  
    method,  
    var_objects = list(),  
    necessary_objects = list(),  
    var_override = list(sustain_interval = 30, trace = FALSE, epochs = 50, batch_size = 20,  
        ml_trace = 0, n_cores = 2, data_folds = 2, pl_max_steps = 2, pl_max = 1, pl_anchor =  
        1, pl_min = 0, sustain_track = TRUE, sustain_iso_code = "DEU", data_val_size = 0.25,  
        lr_rate = 0.001, lr_warm_up_ratio = 0.01)  
)
```

Arguments

object_name	string Name of the object to generate the arguments for.
method	string Name of the method of the object to generate the arguments for.
var_objects	list of other objects which should be combined with the other arguments.
necessary_objects	list of other objects which are part of every combination.
var_override	Named list containing the arguments which should be set to a specific value for all combinations.

Value

Returns a list with combinations of arguments.

Note

var_objects, necessary_objects, and var_override the names must exactly match the name of the parameter. Otherwise they are not applied. Names of arguments which are not part of a method are ignored. #'

See Also

Other Utils TestThat Developers: [check_adjust_n_samples_on_CI\(\)](#), [generate_embeddings\(\)](#), [generate_tensors\(\)](#), [get_current_args_for_print\(\)](#), [get_fixed_test_tensor\(\)](#), [get_test_data_for_classifier\(\)](#), [random_bool_on_CI\(\)](#)

`generate_embeddings` *Generate test embeddings*

Description

Functions generates a random test embedding that can be used for testing methods and functions. The embeddings have the shape (Batch, Times,Features).

Usage

```
generate_embeddings(times, features, seq_len, pad_value)
```

Arguments

<code>times</code>	<code>int</code> Maximal length of a sequence.
<code>features</code>	<code>int</code> Number of features of the sequence.
<code>seq_len</code>	Numeric vector containing the length of the given cases. The length of this vector determines the value for 'Batch'. Values must be at least 1 and maximal times.
<code>pad_value</code>	<code>int</code> Value used to indicate padding.

Value

Returns an array with dim (length(`seq_len`),`times`,`features`).

Note

To generate a 'PyTorch' object please use [generate_tensors](#).

See Also

Other Utils TestThat Developers: [check_adjust_n_samples_on_CI\(\)](#), [generate_args_for_tests\(\)](#), [generate_tensors\(\)](#), [get_current_args_for_print\(\)](#), [get_fixed_test_tensor\(\)](#), [get_test_data_for_classifier\(\)](#), [random_bool_on_CI\(\)](#)

`generate_id` *Generate ID suffix for objects*

Description

Function for generating an ID suffix for objects of class [TextEmbeddingModel](#), [TEClassifierRegular](#), and [TEClassifierProtoNet](#).

Usage

```
generate_id(length = 16)
```

Arguments

length int determining the length of the id suffix.

Value

Returns a string of the requested length.

See Also

Other Utils Developers: [auto_n_cores\(\)](#), [create_object\(\)](#), [create_synthetic_units_from_matrix\(\)](#), [get_n_chunks\(\)](#), [get_synthetic_cases_from_matrix\(\)](#), [matrix_to_array_c\(\)](#), [tensor_to_matrix_c\(\)](#), [to_categorical_c\(\)](#)

generate_tensors *Generate test tensors*

Description

Functions generates a random test tensor that can be used for testing methods and functions based on 'PyTorch'. The tensors have the shape (Batch, Times, Features).

Usage

```
generate_tensors(times, features, seq_len, pad_value)
```

Arguments

times int Maximal length of a sequence.
features int Number of features of the sequence.
seq_len Numeric vector containing the length of the given cases. The length of this vector determines the value for 'Batch'. Values must be at least 1 and maximal times.
pad_value int Value used to indicate padding.

Value

Returns an object of class Tensor from 'PyTorch'.

Note

To request a R array please use [generate_embeddings](#).

See Also

Other Utils TestThat Developers: [check_adjust_n_samples_on_CI\(\)](#), [generate_args_for_tests\(\)](#), [generate_embeddings\(\)](#), [get_current_args_for_print\(\)](#), [get_fixed_test_tensor\(\)](#), [get_test_data_for_classification\(\)](#), [random_bool_on_CI\(\)](#)

`get_alpha_3_codes` *Country Alpha 3 Codes*

Description

Function for requesting a vector containing the alpha-3 codes for most countries.

Usage

```
get_alpha_3_codes()
```

Value

Returns a vector containing the alpha-3 codes for most countries.

See Also

Other Utils Sustainability Developers: [summarize_tracked_sustainability\(\)](#)

`get_batches_index` *Assign cases to batches*

Description

Function groups cases into batches.

Usage

```
get_batches_index(number_rows, batch_size, zero_based = FALSE)
```

Arguments

<code>number_rows</code>	int representing the number of cases or rows of a matrix or array.
<code>batch_size</code>	int size of a batch.
<code>zero_based</code>	bool If TRUE the indices of the cases within each batch are zero based. One based if FALSE.

Value

Returns a list of batches. Each entry in the list contains a vector of int representing the cases belonging to that batch.

See Also

Other Utils Python Data Management Developers: [class_vector_to_py_dataset\(\)](#), [data.frame_to_py_dataset\(\)](#), [prepare_r_array_for_dataset\(\)](#), [py_dataset_to_embeddings\(\)](#), [reduce_to_unique\(\)](#), [tensor_list_to_numpy\(\)](#), [tensor_to_numpy\(\)](#)

get_called_args *Called arguments*

Description

Function for receiving all arguments that were called by a method or function.

Usage

```
get_called_args(n = 1)
```

Arguments

n int level of the nested environments where to extract the arguments.

Value

Returns a named list of all arguments and their values.

See Also

Other Parameter Dictionary: [get_TEClassifiers_class_names\(\)](#), [get_depr_obj_names\(\)](#), [get_magnitude_values\(\)](#), [get_param_def\(\)](#), [get_param_dict\(\)](#), [get_param_doc_desc\(\)](#)

get_coder_metrics *Calculate reliability measures based on content analysis*

Description

This function calculates different reliability measures which are based on the empirical research method of content analysis.

Usage

```
get_coder_metrics(  
  true_values = NULL,  
  predicted_values = NULL,  
  return_names_only = FALSE  
)
```

Arguments

true_values factor containing the true labels/categories.

predicted_values factor containing the predicted labels/categories.

return_names_only bool If TRUE returns only the names of the resulting vector. Use FALSE to request computation of the values.

Value

If `return_names_only = FALSE` returns a vector with the following reliability measures:

- **`iota_index`**: Iota Index from the Iota Reliability Concept Version 2.
- **`min_iota2`**: Minimal Iota from Iota Reliability Concept Version 2.
- **`avg_iota2`**: Average Iota from Iota Reliability Concept Version 2.
- **`max_iota2`**: Maximum Iota from Iota Reliability Concept Version 2.
- **`min_alpha`**: Minmal Alpha Reliability from Iota Reliability Concept Version 2.
- **`avg_alpha`**: Average Alpha Reliability from Iota Reliability Concept Version 2.
- **`max_alpha`**: Maximum Alpha Reliability from Iota Reliability Concept Version 2.
- **`static_iota_index`**: Static Iota Index from Iota Reliability Concept Version 2.
- **`dynamic_iota_index`**: Dynamic Iota Index Iota Reliability Concept Version 2.
- **`kalpha_nominal`**: Krippendorff's Alpha for nominal variables.
- **`kalpha_ordinal`**: Krippendorff's Alpha for ordinal variables.
- **`kendall`**: Kendall's coefficient of concordance W with correction for ties.
- **`c_kappa_unweighted`**: Cohen's Kappa unweighted.
- **`c_kappa_linear`**: Weighted Cohen's Kappa with linear increasing weights.
- **`c_kappa_squared`**: Weighted Cohen's Kappa with quadratic increasing weights.
- **`kappa_fleiss`**: Fleiss' Kappa for multiple raters without exact estimation.
- **`percentage_agreement`**: Percentage Agreement.
- **`balanced_accuracy`**: Average accuracy within each class.
- **`gwet_ac1_nominal`**: Gwet's Agreement Coefficient 1 (AC1) for nominal data which is unweighted.
- **`gwet_ac2_linear`**: Gwet's Agreement Coefficient 2 (AC2) for ordinal data with linear weights.
- **`gwet_ac2_quadratic`**: Gwet's Agreement Coefficient 2 (AC2) for ordinal data with quadratic weights.

If `return_names_only = TRUE` returns only the names of the vector elements.

See Also

Other performance measures: [calc_standard_classification_measures\(\)](#), [cohens_kappa\(\)](#), [fleiss_kappa\(\)](#), [gwet_ac\(\)](#), [kendalls_w\(\)](#), [kripp_alpha\(\)](#)

get_current_args_for_print
Print arguments

Description

Functions prints the used arguments. The aim of this function is to print the arguments to the console that resulted in a failed test.

Usage

```
get_current_args_for_print(arg_list)
```

Arguments

arg_list Named list of arguments. The list should be generated with [generate_args_for_tests\(\)](#).

Value

Function does nothing return.

See Also

Other Utils TestThat Developers: [check_adjust_n_samples_on_CI\(\)](#), [generate_args_for_tests\(\)](#), [generate_embeddings\(\)](#), [generate_tensors\(\)](#), [get_fixed_test_tensor\(\)](#), [get_test_data_for_classifiers\(\)](#), [random_bool_on_CI\(\)](#)

get_depr_obj_names Get names of deprecated objects

Description

Function returns the names of all objects that are deprecated.

Usage

```
get_depr_obj_names()
```

Value

Returns a vector containing the names.

See Also

Other Parameter Dictionary: [get_TEClassifiers_class_names\(\)](#), [get_called_args\(\)](#), [get_magnitude_values\(\)](#), [get_param_def\(\)](#), [get_param_dict\(\)](#), [get_param_doc_desc\(\)](#)

`get_desc_for_core_model_architecture`
Generate documentation for core models

Description

Function for generating the documentation of a specific core model.

Usage

```
get_desc_for_core_model_architecture(
  name,
  title_format = "bold",
  inc_img = FALSE
)
```

Arguments

<code>name</code>	string Name of the core model.
<code>title_format</code>	string Kind of format of the title.
<code>inc_img</code>	bool Include a visualization of the layer.

Value

Returns a string containing the description written in rmarkdown.

See Also

Other Utils Documentation: [build_documentation_for_model\(\)](#), [build_layer_stack_documentation_for_vignette\(\)](#), [get_dict_cls_type\(\)](#), [get_dict_core_models\(\)](#), [get_dict_input_types\(\)](#), [get_layer_dict\(\)](#), [get_layer_documentation\(\)](#), [get_parameter_documentation\(\)](#)

`get_file_extension` *Get file extension*

Description

Function for requesting the file extension

Usage

```
get_file_extension(file_path)
```

Arguments

<code>file_path</code>	string Path to a file.
------------------------	------------------------

Value

Returns the extension of a file as a string.

See Also

Other Utils File Management Developers: [create_dir\(\)](#)

get_fixed_test_tensor *Generate static test tensor*

Description

Function generates a static test tensor which is always the same.

Usage

```
get_fixed_test_tensor(pad_value)
```

Arguments

pad_value int Value used to indicate padding.

Value

Returns an object of class Tensor which is always the same except padding. Shape (5,3,7).

See Also

Other Utils TestThat Developers: [check_adjust_n_samples_on_CI\(\)](#), [generate_args_for_tests\(\)](#), [generate_embeddings\(\)](#), [generate_tensors\(\)](#), [get_current_args_for_print\(\)](#), [get_test_data_for_classifiers\(\)](#), [random_bool_on_CI\(\)](#)

get_layer_documentation
Generate layer documentation

Description

Function for generating the documentation of a specific layer.

Usage

```
get_layer_documentation(
  layer_name,
  title_format = "bold",
  subtitle_format = "italic",
  inc_img = FALSE,
  inc_params = FALSE,
  inc_references = FALSE
)
```

Arguments

<code>layer_name</code>	string	Name of the layer.
<code>title_format</code>	string	Kind of format of the title.
<code>subtitle_format</code>		string Kind of format for all sub-titles.
<code>inc_img</code>	bool	Include a visualization of the layer.
<code>inc_params</code>	bool	Include a description of every parameter of the layer.
<code>inc_references</code>	bool	Include a list of literature references for the layer.

Value

Returns a string containing the description written in rmarkdown.

See Also

Other Utils Documentation: [build_documentation_for_model\(\)](#), [build_layer_stack_documentation_for_vignette\(\)](#), [get_desc_for_core_model_architecture\(\)](#), [get_dict_cls_type\(\)](#), [get_dict_core_models\(\)](#), [get_dict_input_types\(\)](#), [get_layer_dict\(\)](#), [get_parameter_documentation\(\)](#)

`get_magnitude_values` *Magnitudes of an argument*

Description

Function calculates different magnitude for a numeric argument.

Usage

```
get_magnitude_values(magnitude, n_elements = 9, max, min)
```

Arguments

<code>magnitude</code>	double	Factor using for creating the magnitude.
<code>n_elements</code>	int	Number of values to return.
<code>max</code>	double	The maximal value.
<code>min</code>	double	The minimal value.

Value

Returns a numeric vector with the generated values. The values are calculated with the following formula: max * magnitudeⁱ for i=1,...,n_elements. Only values equal or greater min are returned.

See Also

Other Parameter Dictionary: [get_TEClassifiers_class_names\(\)](#), [get_called_args\(\)](#), [get_depr_obj_names\(\)](#), [get_param_def\(\)](#), [get_param_dict\(\)](#), [get_param_doc_desc\(\)](#)

get_n_chunks

Get the number of chunks/sequences for each case

Description

Function for calculating the number of chunks/sequences for every case.

Usage

```
get_n_chunks(text_embeddings, features, times, pad_value = -100)
```

Arguments

text_embeddings	data.frame or array containing the text embeddings.
features	int Number of features within each sequence.
times	int Number of sequences.
pad_value	int Value indicating padding. This value should no be in the range of regular values for computations. Thus it is not recommended to chance this value. Default is -100. Allowed values: x <= -100

Value

Namedvector of integers representing the number of chunks/sequences for every case.

See Also

Other Utils Developers: [auto_n_cores\(\)](#), [create_object\(\)](#), [create_synthetic_units_from_matrix\(\)](#), [generate_id\(\)](#), [get_synthetic_cases_from_matrix\(\)](#), [matrix_to_array_c\(\)](#), [tensor_to_matrix_c\(\)](#), [to_categorical_c\(\)](#)

`get_parameter_documentation`
Generate layer documentation

Description

Function for generating the documentation of a specific layer.

Usage

```
get_parameter_documentation(  
    param_name,  
    param_dict,  
    as_list = TRUE,  
    inc_param_name = TRUE  
)
```

Arguments

<code>param_name</code>	string Name of the parameter.
<code>param_dict</code>	list storing the parameter description.
<code>as_list</code>	bool If TRUE returns the element as part of a list.
<code>inc_param_name</code>	bool If TRUE the documentation includes the name of the parameter.

Value

Returns a string containing the description written in rmarkdown.

See Also

Other Utils Documentation: [build_documentation_for_model\(\)](#), [build_layer_stack_documentation_for_vignette\(\)](#), [get_desc_for_core_model_architecture\(\)](#), [get_dict_cls_type\(\)](#), [get_dict_core_models\(\)](#), [get_dict_input_types\(\)](#), [get_layer_dict\(\)](#), [get_layer_documentation\(\)](#)

`get_param_def` *Definition of an argument*

Description

Function returns the definition of an argument. Please note that only definitions of arguments can be requested which are used for transformers or classifier models.

Usage

```
get_param_def(param_name)
```

Arguments

param_name string Name of the parameter to request its definition.

Value

Returns a list with the definition of the argument. See [get_param_dict](#) for more details.

See Also

Other Parameter Dictionary: [get_TEClassifiers_class_names\(\)](#), [get_called_args\(\)](#), [get_depr_obj_names\(\)](#), [get_magnitude_values\(\)](#), [get_param_dict\(\)](#), [get_param_doc_desc\(\)](#)

get_param_dict *Get dictionary of all parameters*

Description

Function provides a list containing important characteristics of the parameter used in the models. The list does contain only the definition of arguments for transformer models and all classifiers. The arguments of other functions in this package are documented separately.

The aim of this list is to automatize argument checking and widget generation for *AI for Education - Studio*.

Usage

`get_param_dict()`

Value

Returns a named list. The names correspond to specific arguments. The list contains a list for every argument with the following components:

- type: The type of allowed values.
- allow_null: A bool indicating if the argument can be set to NULL.
- min: The minimal value the argument can be. Set to NULL if not relevant. Set to -Inf if there is no minimum.
- max: The maximal value the argument can be. Set to NULL if not relevant. Set to Inf if there is no Minimum.
- desc: A string which includes the description of the argument written in markdown. This string is for the documentation the parameter.
- values_desc: A named list containing a description of every possible value. The names must exactly match the strings in allowed_values. Descriptions should be written in markdown.
- allowed_values: vector of allowed values. This is only relevant if the argument is not numeric. During the checking of the arguments it is checked if the provided values can be found in this vector. If all values are allowed set to NULL.

- `default_value`: The default value of the argument. If there is no default set to NULL.
- `default_historic`: Historic default value. This can be necessary for backward compatibility.
- `gui_box`: string Name of the box in AI for Education - Studio where the argument appears. If it should not appear set to NULL.
- `gui_label`: string Label of the controlling widget in AI for Education - Studio.

See Also

Other Parameter Dictionary: [get_TEClassifiers_class_names\(\)](#), [get_called_args\(\)](#), [get_depr_obj_names\(\)](#), [get_magnitude_values\(\)](#), [get_param_def\(\)](#), [get_param_doc_desc\(\)](#)

`get_param_doc_desc` *Description of an argument*

Description

Function provides the description of an argument in markdown. Its aim is to be used for documenting the parameter of functions.

Usage

`get_param_doc_desc(param_name)`

Arguments

`param_name` string Name of the parameter to request its definition.

Value

Returns a string which contains the description of the argument in markdown. The concrete format depends on the type of the argument.

See Also

Other Parameter Dictionary: [get_TEClassifiers_class_names\(\)](#), [get_called_args\(\)](#), [get_depr_obj_names\(\)](#), [get_magnitude_values\(\)](#), [get_param_def\(\)](#), [get_param_dict\(\)](#)

get_py_package_version

Get versions of a specific python package

Description

Function for requesting the version of a specific python package.

Usage

```
get_py_package_version(package_name)
```

Arguments

package_name string Name of the package.

Value

Returns the version as string or NA if the package does not exist.

See Also

Other Utils Python Developers: [get_py_package_versions\(\)](#), [load_all_py_scripts\(\)](#), [load_py_scripts\(\)](#), [run_py_file\(\)](#)

get_py_package_versions

Get versions of python components

Description

Function for requesting a summary of the versions of all critical python components.

Usage

```
get_py_package_versions()
```

Value

Returns a list that contains the version number of python and the versions of critical python packages. If a package is not available version is set to NA.

See Also

Other Utils Python Developers: [get_py_package_version\(\)](#), [load_all_py_scripts\(\)](#), [load_py_scripts\(\)](#), [run_py_file\(\)](#)

```
get_synthetic_cases_from_matrix
    Create synthetic cases for balancing training data
```

Description

This function creates synthetic cases for balancing the training with classifier models.

Usage

```
get_synthetic_cases_from_matrix(
  matrix_form,
  times,
  features,
  target,
  sequence_length,
  method = c("knnor"),
  min_k = 1,
  max_k = 6
)
```

Arguments

matrix_form	Named <code>matrix</code> containing the text embeddings in a matrix form.
times	<code>int</code> for the number of sequences/times.
features	<code>int</code> for the number of features within each sequence.
target	Named <code>factor</code> containing the labels of the corresponding embeddings.
sequence_length	<code>int</code> Length of the text embedding sequences.
method	<code>vector</code> containing strings of the requested methods for generating new cases. Currently "knnor" from this package is available.
min_k	<code>int</code> The minimal number of nearest neighbors during sampling process.
max_k	<code>int</code> The maximum number of nearest neighbors during sampling process.

Value

`list` with the following components:

- `synthetic_embeddings`: Named `data.frame` containing the text embeddings of the synthetic cases.
- `synthetic_targets`: Named `factor` containing the labels of the corresponding synthetic cases.
- `n_synthetic_units`: `table` showing the number of synthetic cases for every label/category.

See Also

Other Utils Developers: [auto_n_cores\(\)](#), [create_object\(\)](#), [create_synthetic_units_from_matrix\(\)](#), [generate_id\(\)](#), [get_n_chunks\(\)](#), [matrix_to_array_c\(\)](#), [tensor_to_matrix_c\(\)](#), [to_categorical_c\(\)](#)

get_TEClassifiers_class_names

Get names of classifiers

Description

Function returns the names of all classifiers which are child classes of a specific super class.

Usage

```
get_TEClassifiers_class_names(super_class = NULL)
```

Arguments

super_class string Name of the super class the classifiers should be child of. To request the names of all classifiers set this argument to NULL.

Value

Returns a vector containing the names of the classifiers.

See Also

Other Parameter Dictionary: [get_called_args\(\)](#), [get_depr_obj_names\(\)](#), [get_magnitude_values\(\)](#), [get_param_def\(\)](#), [get_param_dict\(\)](#), [get_param_doc_desc\(\)](#)

get_test_data_for_classifiers

Get test data

Description

Function returns example data for testing the package

Usage

```
get_test_data_for_classifiers(class_range = c(2, 3), path_test_embeddings)
```

Arguments

class_range vector containing the number of classes.
path_test_embeddings string Path to the location where the test data is stored.

Value

Returns a list with test data.

See Also

Other Utils TestThat Developers: [check_adjust_n_samples_on_CI\(\)](#), [generate_args_for_tests\(\)](#), [generate_embeddings\(\)](#), [generate_tensors\(\)](#), [get_current_args_for_print\(\)](#), [get_fixed_test_tensor\(\)](#), [random_bool_on_CI\(\)](#)

gwet_ac

*Calculate Gwet's AC1 and AC2***Description**

This function calculates Gwets Agreement Coefficients.

Usage

```
gwet_ac(rater_one, rater_two, additional_raters = NULL)
```

Arguments

rater_one	factor rating of the first coder.
rater_two	factor ratings of the second coder.
additional_raters	list Additional raters with same requirements as rater_one and rater_two. If there are no additional raters set to NULL.

Value

Returns a list with the following entries

- ac1: Gwet's Agreement Coefficient 1 (AC1) for nominal data which is unweighted.
- ac2_linear: Gwet's Agreement Coefficient 2 (AC2) for ordinal data with linear weights.
- ac2_quadratic: Gwet's Agreement Coefficient 2 (AC2) for ordinal data with quadratic weights.

Note

Weights are calculated as described in Gwet (2021).

Missing values are supported.

References

Gwet, K. L. (2021). Handbook of inter-rater reliability: The definitive guide to measuring the extent of agreement among raters (Fifth edition, volume 1). AgreeStat Analytics.

See Also

Other performance measures: [calc_standard_classification_measures\(\)](#), [cohens_kappa\(\)](#), [fleiss_kappa\(\)](#), [get_coder_metrics\(\)](#), [kendalls_w\(\)](#), [kripp_alpha\(\)](#)

imdb_movie_reviews *Standford Movie Review Dataset*

Description

A [data.frame](#) consisting of a subset of 100 negative and 200 positive movie reviews from the dataset provided by Maas et al. (2011). The [data.frame](#) consists of three columns. The first column 'text' stores the movie review. The second stores the labels (0 = negative, 1 = positive). The last column stores the id. The purpose of the data is for illustration in vignettes and for tests.

Usage

`imdb_movie_reviews`

Format

`data.frame`

References

Maas, A. L., Daly, R. E., Pham, P. T., Huang, D., Ng, A. Y., & Potts, C. (2011). Learning Word Vectors for Sentiment Analysis. In D. Lin, Y. Matsumoto, & R. Mihalcea (Eds.), Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies (pp. 142–150). Association for Computational Linguistics. <doi: 10.5555/2002472.2002491>

install_aifedducation *Install aifedducation on a machine*

Description

Function for installing 'aifedducation' on a machine.

Using a virtual environment (`use_conda=FALSE`) If 'python' is already installed the installed version is used. In the case that the required version of 'python' is different from the existing version the new version is installed. In all other cases python will be installed on the system.

#' Using a conda environment (`use_conda=TRUE`) If 'miniconda' is already existing on the machine no installation of 'miniconda' is applied. In this case the system checks for update and updates 'miniconda' to the newest version. If 'miniconda' is not found on the system it will be installed.

Usage

```
install_aifeduation(
  install_aifeduation_studio = TRUE,
  python_version = "3.12",
  cuda_version = "12.4",
  use_conda = FALSE
)
```

Arguments

<code>install_aifeduation_studio</code>	
	bool If TRUE all necessary R packages are installed for using AI for Education Studio.
<code>python_version</code>	string Python version to use/install.
<code>cuda_version</code>	string determining the requested version of cuda.
<code>use_conda</code>	bool If TRUE installation installs 'miniconda' and uses 'conda' as package manager. If FALSE installation installs python and uses virtual environments for package management.

Value

Function does nothing return. It installs python, optional R packages, and necessary 'python' packages on a machine.

Note

On MAC OS torch will be installed without support for cuda.

See Also

Other Installation and Configuration: [check_aif_py_modules\(\)](#), [install_aifeduation_studio\(\)](#), [install_py_modules\(\)](#), [prepare_session\(\)](#), [set_transformers_logger\(\)](#), [update_aifeduation\(\)](#)

install_aifeduation_studio

Install 'AI for Education - Studio' on a machine

Description

Function installs/updates all relevant R packages necessary to run the shiny app "AI for Education - Studio".

Usage

```
install_aifeduation_studio()
```

Value

Function does nothing return. It installs/updates R packages.

See Also

Other Installation and Configuration: [check_aif_py_modules\(\)](#), [install_aifedducation\(\)](#), [install_py_modules\(\)](#), [prepare_session\(\)](#), [set_transformers_logger\(\)](#), [update_aifedducation\(\)](#)

install_py_modules *Installing necessary python modules to an environment*

Description

Function for installing the necessary python modules.

Usage

```
install_py_modules(  
    envname = "aifedducation",  
    transformer_version = "<=4.52.4",  
    tokenizers_version = "<=0.21.1",  
    pandas_version = "<=2.3.0",  
    datasets_version = "<=3.6.0",  
    codecarbon_version = "<=3.0.2",  
    safetensors_version = "<=0.5.3",  
    torcheval_version = "<=0.0.7",  
    accelerate_version = "<=1.8.1",  
    pytorch_cuda_version = "12.6",  
    python_version = "3.12",  
    remove_first = FALSE,  
    use_conda = FALSE  
)
```

Arguments

envname	string	Name of the environment where the packages should be installed.
transformer_version	string	determining the desired version of the python library 'transformers'.
tokenizers_version	string	determining the desired version of the python library 'tokenizers'.
pandas_version	string	determining the desired version of the python library 'pandas'.
datasets_version	string	determining the desired version of the python library 'datasets'.
codecarbon_version	string	determining the desired version of the python library 'codecarbon'.

```

safetensors_version
    string determining the desired version of the python library 'safetensors'.

torcheval_version
    string determining the desired version of the python library 'torcheval'.

accelerate_version
    string determining the desired version of the python library 'accelerate'.

pytorch_cuda_version
    string determining the desired version of 'cuda' for 'PyTorch'. To install 'PyTorch' without cuda set to NULL.

python_version  string Python version to use.

remove_first    bool If TRUE removes the environment completely before recreating the environment and installing the packages. If FALSE the packages are installed in the existing environment without any prior changes.

use_conda       bool If TRUE uses 'conda' for package management. If FALSE uses virtual environments for package management.

```

Value

Returns no values or objects. Function is used for installing the necessary python libraries in a conda environment.

Note

Function tries to identify the type of operating system. In the case that MAC OS is detected 'PyTorch' is installed without support for cuda.

See Also

Other Installation and Configuration: [check_aif_py_modules\(\)](#), [install_aifedducation\(\)](#), [install_aifedducation_stu](#)
[prepare_session\(\)](#), [set_transformers_logger\(\)](#), [update_aifedducation\(\)](#)

kendalls_w

Calculate Kendall's coefficient of concordance w

Description

This function calculates Kendall's coefficient of concordance w with and without correction.

Usage

```
kendalls_w(rater_one, rater_two, additional_raters = NULL)
```

Arguments

rater_one factor rating of the first coder.
 rater_two factor ratings of the second coder.
 additional_raters
 list Additional raters with same requirements as rater_one and rater_two.
 If there are no additional raters set to NULL.

Value

Returns a list containing the results for Kendall's coefficient of concordance w with and without correction.

See Also

Other performance measures: [calc_standard_classification_measures\(\)](#), [cohens_kappa\(\)](#), [fleiss_kappa\(\)](#), [get_coder_metrics\(\)](#), [gwt_ac\(\)](#), [kripp_alpha\(\)](#)

knnor

K-Nearest Neighbor OveRampling approach (KNNOR)

Description

K-Nearest Neighbor OveRampling approach (KNNOR)

Usage

```
knnor(dataset, k, aug_num, cycles_number_limit = 100L)
```

Arguments

dataset list containing the following fields:
 • embeddings: an 2-D array (matrix) with size batch x times*features
 • labels: an 1-D array (vector) of integers with batch elements
 k unsigned integer number of nearest neighbors
 aug_num unsigned integer number of datapoints to be augmented
 cycles_number_limit
 unsigned integer number of maximum try cycles

Value

Returns artificial points (2-D array (matrix) with size aug_numxtimes*features')

References

Islam, A., Belhaouari, S. B., Rehman, A. U. & Bensmail, H. (2022). KNNOR: An oversampling technique for imbalanced datasets. Applied Soft Computing, 115, 108288. <https://doi.org/10.1016/j.asoc.2021.108288>

`knnor_is_same_class` *Validate a new point*

Description

Function written in C++ for validating a new point (KNNOR-Validation)

Usage

```
knnor_is_same_class(new_point, dataset, labels, k)
```

Arguments

<code>new_point</code>	1-D array (vector) new data point to be validated before adding (with <code>times*features</code> elements)
<code>dataset</code>	2-D array (matrix) current embeddings (with size <code>batch x times*features</code>)
<code>labels</code>	1-D array (vector) of integers with batch elements
<code>k</code>	unsigned integer number of nearest neighbors

Value

Returns TRUE if a new point can be added, otherwise - FALSE

`kripp_alpha` *Calculate Krippendorff's Alpha*

Description

This function calculates different Krippendorff's Alpha for nominal and ordinal variables.

Usage

```
kripp_alpha(rater_one, rater_two, additional_raters = NULL)
```

Arguments

<code>rater_one</code>	factor rating of the first coder.
<code>rater_two</code>	factor ratings of the second coder.
<code>additional_raters</code>	list Additional raters with same requirements as <code>rater_one</code> and <code>rater_two</code> . If there are no additional raters set to NULL.

Value

Returns a list containing the results for Krippendorff's Alpha for nominal and ordinal data.

Note

Missing values are supported.

References

Krippendorff, K. (2019). Content Analysis: An Introduction to Its Methodology (4th Ed.). SAGE

See Also

Other performance measures: [calc_standard_classification_measures\(\)](#), [cohens_kappa\(\)](#), [fleiss_kappa\(\)](#), [get_coder_metrics\(\)](#), [gwt_ac\(\)](#), [kendalls_w\(\)](#)

LargeDataSetBase

Abstract base class for large data sets

Description

This object contains public and private methods which may be useful for every large data sets. Objects of this class are not intended to be used directly.

Value

Returns a new object of this class.

Methods**Public methods:**

- [LargeDataSetBase\\$n_cols\(\)](#)
- [LargeDataSetBase\\$n_rows\(\)](#)
- [LargeDataSetBase\\$get_colnames\(\)](#)
- [LargeDataSetBase\\$get_dataset\(\)](#)
- [LargeDataSetBase\\$reduce_to_unique_ids\(\)](#)
- [LargeDataSetBase\\$select\(\)](#)
- [LargeDataSetBase\\$get_ids\(\)](#)
- [LargeDataSetBase\\$save\(\)](#)
- [LargeDataSetBase\\$load_from_disk\(\)](#)
- [LargeDataSetBase\\$load\(\)](#)
- [LargeDataSetBase\\$set_package_versions\(\)](#)
- [LargeDataSetBase\\$get_package_versions\(\)](#)
- [LargeDataSetBase\\$get_all_fields\(\)](#)
- [LargeDataSetBase\\$clone\(\)](#)

Method n_cols(): Number of columns in the data set.

Usage:

`LargeDataSetBase$n_cols()`

Returns: int describing the number of columns in the data set.

Method `n_rows()`: Number of rows in the data set.

Usage:

`LargeDataSetBase$n_rows()`

Returns: int describing the number of rows in the data set.

Method `get_colnames()`: Get names of the columns in the data set.

Usage:

`LargeDataSetBase$get_colnames()`

Returns: vector containing the names of the columns as strings.

Method `get_dataset()`: Get data set.

Usage:

`LargeDataSetBase$get_dataset()`

Returns: Returns the data set of this object as an object of class `datasets.arrow_dataset.Dataset`.

Method `reduce_to_unique_ids()`: Reduces the data set to a data set containing only unique ids. In the case an id exists multiple times in the data set the first case remains in the data set. The other cases are dropped.

Attention Calling this method will change the data set in place.

Usage:

`LargeDataSetBase$reduce_to_unique_ids()`

Returns: Method does not return anything. It changes the data set of this object in place.

Method `select()`: Returns a data set which contains only the cases belonging to the specific indices.

Usage:

`LargeDataSetBase$select(indicies)`

Arguments:

`indicies` vector of int for selecting rows in the data set. **Attention** The indices are zero-based.

Returns: Returns a data set of class `datasets.arrow_dataset.Dataset` with the selected rows.

Method `get_ids()`: Get ids

Usage:

`LargeDataSetBase$get_ids()`

Returns: Returns a vector containing the ids of every row as strings.

Method `save()`: Saves a data set to disk.

Usage:

```
LargeDataSetBase$save(dir_path, folder_name, create_dir = TRUE)
```

Arguments:

dir_path Path where to store the data set.

folder_name string Name of the folder for storing the data set.

create_dir bool If True the directory will be created if it does not exist.

Returns: Method does not return anything. It write the data set to disk.

Method load_from_disk(): loads an object of class [LargeDataSetBase](#) from disk 'and updates the object to the current version of the package.

Usage:

```
LargeDataSetBase$load_from_disk(dir_path)
```

Arguments:

dir_path Path where the data set set is stored.

Returns: Method does not return anything. It loads an object from disk.

Method load(): Loads a data set from disk.

Usage:

```
LargeDataSetBase$load(dir_path)
```

Arguments:

dir_path Path where the data set is stored.

Returns: Method does not return anything. It loads a data set from disk.

Method set_package_versions(): Method for setting the package version for 'aifedducation', 'reticulate', 'torch', and 'numpy' to the currently used versions.

Usage:

```
LargeDataSetBase$set_package_versions()
```

Returns: Method does not return anything. It is used to set the private fields fo package versions.

Method get_package_versions(): Method for requesting a summary of the R and python packages' versions used for creating the model.

Usage:

```
LargeDataSetBase$get_package_versions()
```

Returns: Returns a list containing the versions of the relevant R and python packages.

Method get_all_fields(): Return all fields.

Usage:

```
LargeDataSetBase$get_all_fields()
```

Returns: Method returns a list containing all public and private fields of the object.

Method clone(): The objects of this class are cloneable with this method.

Usage:

```
LargeDataSetBase$clone(deep = FALSE)
```

Arguments:

deep Whether to make a deep clone.

See Also

Other R6 Classes for Developers: [AIFEBaseModel](#), [ClassifiersBasedOnTextEmbeddings](#), [ModelsBasedOnTextEmbedding](#)
[TEClassifiersBasedOnProtoNet](#), [TEClassifiersBasedOnRegular](#)

LargeDataSetForText *Abstract class for large data sets containing raw texts*

Description

This object stores raw texts. The data of this objects is not stored in memory directly. By using memory mapping these objects allow to work with data sets which do not fit into memory/RAM.

Value

Returns a new object of this class.

Super class

[aifedducation::LargeDataSetBase](#) -> LargeDataSetForText

Methods**Public methods:**

- [LargeDataSetForText\\$new\(\)](#)
- [LargeDataSetForText\\$add_from_files_txt\(\)](#)
- [LargeDataSetForText\\$add_from_files_pdf\(\)](#)
- [LargeDataSetForText\\$add_from_files_xlsx\(\)](#)
- [LargeDataSetForText\\$add_from_data.frame\(\)](#)
- [LargeDataSetForText\\$get_private\(\)](#)
- [LargeDataSetForText\\$clone\(\)](#)

Method `new()`: Method for creation of [LargeDataSetForText](#) instance. It can be initialized with `init_data` parameter if passed (Uses `add_from_data.frame()` method if `init_data` is `data.frame`).

Usage:

```
LargeDataSetForText$new(init_data = NULL)
```

Arguments:

`init_data` Initial `data.frame` for dataset.

Returns: A new instance of this class initialized with `init_data` if passed.

Method `add_from_files_txt()`: Method for adding raw texts saved within .txt files to the data set. Please note the the directory should contain one folder for each .txt file. In order to create an informative data set every folder can contain the following additional files:

- `bib_entry.txt`: containing a text version of the bibliographic information of the raw text.

- license.txt: containing a statement about the license to use the raw text such as "CC BY".
 - url_license.txt: containing the url/link to the license in the internet.
 - text_license.txt: containing the license in raw text.
 - url_source.txt: containing the url/link to the source in the internet.
- The id of every .txt file is the file name without file extension. Please be aware to provide unique file names. Id and raw texts are mandatory, bibliographic and license information are optional.

Usage:

```
LargeDataSetForText$add_from_files_txt(
  dir_path,
  batch_size = 500,
  log_file = NULL,
  log_write_interval = 2,
  log_top_value = 0,
  log_top_total = 1,
  log_top_message = NA,
  clean_text = TRUE,
  trace = TRUE
)
```

Arguments:

dir_path Path to the directory where the files are stored.

batch_size int determining the number of files to process at once.

log_file string Path to the file where the log should be saved. If no logging is desired set this argument to NULL.

log_write_interval int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if log_file is not NULL.

log_top_value int indicating the current iteration of the process.

log_top_total int determining the maximal number of iterations.

log_top_message string providing additional information of the process.

clean_text bool If TRUE the text is modified to improve the quality of the following analysis:

- Some special symbols are removed.
- All spaces at the beginning and the end of a row are removed.
- Multiple spaces are reduced to single space.
- All rows with a number from 1 to 999 at the beginning or at the end are removed (header and footer).
- List of content is removed.
- Hyphenation is made undone.
- Line breaks within a paragraph are removed.
- Multiple line breaks are reduced to a single line break.

trace bool If TRUE information on the progress is printed to the console.

Returns: The method does not return anything. It adds new raw texts to the data set.

Method add_from_files_pdf(): Method for adding raw texts saved within .pdf files to the data set. Please note the the directory should contain one folder for each .pdf file. In order to create an informative data set every folder can contain the following additional files:

- bib_entry.txt: containing a text version of the bibliographic information of the raw text.
 - license.txt: containing a statement about the license to use the raw text such as "CC BY".
 - url_license.txt: containing the url/link to the license in the internet.
 - text_license.txt: containing the license in raw text.
 - url_source.txt: containing the url/link to the source in the internet.
- The id of every .pdf file is the file name without file extension. Please be aware to provide unique file names. Id and raw texts are mandatory, bibliographic and license information are optional.

Usage:

```
LargeDataSetForText$add_from_files_pdf(
  dir_path,
  batch_size = 500,
  log_file = NULL,
  log_write_interval = 2,
  log_top_value = 0,
  log_top_total = 1,
  log_top_message = NA,
  clean_text = TRUE,
  trace = TRUE
)
```

Arguments:

`dir_path` Path to the directory where the files are stored.

`batch_size` int determining the number of files to process at once.

`log_file` string Path to the file where the log should be saved. If no logging is desired set this argument to NULL.

`log_write_interval` int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if `log_file` is not NULL.

`log_top_value` int indicating the current iteration of the process.

`log_top_total` int determining the maximal number of iterations.

`log_top_message` string providing additional information of the process.

`clean_text` bool If TRUE the text is modified to improve the quality of the following analysis:

- Some special symbols are removed.
- All spaces at the beginning and the end of a row are removed.
- Multiple spaces are reduced to single space.
- All rows with a number from 1 to 999 at the beginning or at the end are removed (header and footer).
- List of content is removed.
- Hyphenation is made undone.
- Line breaks within a paragraph are removed.
- Multiple line breaks are reduced to a single line break.

`trace` bool If TRUE information on the progress is printed to the console.

Returns: The method does not return anything. It adds new raw texts to the data set.

Method add_from_files_xlsx(): Method for adding raw texts saved within .xlsx files to the data set. The method assumes that the texts are saved in the rows and that the columns store the id and the raw texts in the columns. In addition, a column for the bibliography information and the license can be added. The column names for these rows must be specified with the following arguments. They must be the same for all .xlsx files in the chosen directory. Id and raw texts are mandatory, bibliographic, license, license's url, license's text, and source's url are optional. Additional columns are dropped.

Usage:

```
LargeDataSetForText$add_from_files_xlsx(
  dir_path,
  trace = TRUE,
  id_column = "id",
  text_column = "text",
  bib_entry_column = "bib_entry",
  license_column = "license",
  url_license_column = "url_license",
  text_license_column = "text_license",
  url_source_column = "url_source",
  log_file = NULL,
  log_write_interval = 2,
  log_top_value = 0,
  log_top_total = 1,
  log_top_message = NA
)
```

Arguments:

dir_path Path to the directory where the files are stored.
trace bool If TRUE prints information on the progress to the console.
id_column string Name of the column storing the ids for the texts.
text_column string Name of the column storing the raw text.
bib_entry_column string Name of the column storing the bibliographic information of the texts.
license_column string Name of the column storing information about the licenses.
url_license_column string Name of the column storing information about the url to the license in the internet.
text_license_column string Name of the column storing the license as text.
url_source_column string Name of the column storing information about the url to the source in the internet.
log_file string Path to the file where the log should be saved. If no logging is desired set this argument to NULL.
log_write_interval int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if **log_file** is not NULL.
log_top_value int indicating the current iteration of the process.
log_top_total int determining the maximal number of iterations.
log_top_message string providing additional information of the process.

Returns: The method does not return anything. It adds new raw texts to the data set.

Method `add_from_data.frame()`: Method for adding raw texts from a `data.frame`

Usage:

```
LargeDataSetForText$add_from_data.frame(data_frame)
```

Arguments:

`data_frame` Object of class `data.frame` with at least the following columns "id", "text", "bib_entry", "license", "url_license", "text_license", and "url_source". If "id" and/or "text" is missing an error occurs. If the other columns are not present in the `data.frame` they are added with empty values(NA). Additional columns are dropped.

Returns: The method does not return anything. It adds new raw texts to the data set.

Method `get_private()`: Method for requesting all private fields and methods. Used for loading and updating an object.

Usage:

```
LargeDataSetForText$get_private()
```

Returns: Returns a list with all private fields and methods.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
LargeDataSetForText$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

See Also

Other Data Management: [EmbeddedText](#), [LargeDataSetForTextEmbeddings](#)

LargeDataSetForTextEmbeddings

Abstract class for large data sets containing text embeddings

Description

This object stores text embeddings which are usually produced by an object of class [TextEmbeddingModel](#). The data of this objects is not stored in memory directly. By using memory mapping these objects allow to work with data sets which do not fit into memory/RAM.

[LargeDataSetForTextEmbeddings](#) are used for storing and managing the text embeddings created with objects of class [TextEmbeddingModel](#). Objects of class [LargeDataSetForTextEmbeddings](#) serve as input for objects of class [ClassifiersBasedOnTextEmbeddings](#) and [TEFeatureExtractor](#). The main aim of this class is to provide a structured link between embedding models and classifiers. Since objects of this class save information on the text embedding model that created the text embedding it ensures that only embeddings generated with same embedding model are combined. Furthermore, the stored information allows objects to check if embeddings of the correct text embedding model are used for training and predicting.

This class is not designed for a direct use.

Value

Returns a new object of this class.

Super class

aifedducation: :LargeDataSetBase -> LargeDataSetForTextEmbeddings

Methods**Public methods:**

- `LargeDataSetForTextEmbeddings$configure()`
- `LargeDataSetForTextEmbeddings$is_configured()`
- `LargeDataSetForTextEmbeddings$get_text_embedding_model_name()`
- `LargeDataSetForTextEmbeddings$get_model_info()`
- `LargeDataSetForTextEmbeddings$load_from_disk()`
- `LargeDataSetForTextEmbeddings$get_model_label()`
- `LargeDataSetForTextEmbeddings$add_feature_extractor_info()`
- `LargeDataSetForTextEmbeddings$get_feature_extractor_info()`
- `LargeDataSetForTextEmbeddings$is_compressed()`
- `LargeDataSetForTextEmbeddings$get_times()`
- `LargeDataSetForTextEmbeddings$get_features()`
- `LargeDataSetForTextEmbeddings$get_original_features()`
- `LargeDataSetForTextEmbeddings$get_pad_value()`
- `LargeDataSetForTextEmbeddings$add_embeddings_from_array()`
- `LargeDataSetForTextEmbeddings$add_embeddings_from_EMBEDDED_TEXT()`
- `LargeDataSetForTextEmbeddings$add_embeddings_from_LargeDataSetForTextEmbeddings()`
- `LargeDataSetForTextEmbeddings$convert_to_EMBEDDED_TEXT()`
- `LargeDataSetForTextEmbeddings$clone()`

Method `configure():` Creates a new object representing text embeddings.

Usage:

```
LargeDataSetForTextEmbeddings$configure(  
  model_name = NA,  
  model_label = NA,  
  model_date = NA,  
  model_method = NA,  
  model_version = NA,  
  model_language = NA,  
  param_seq_length = NA,  
  param_chunks = NULL,  
  param_features = NULL,  
  param_overlap = NULL,  
  param_emb_layer_min = NULL,  
  param_emb_layer_max = NULL,  
  param_emb_pool_type = NULL,
```

```

    param_pad_value = -100,
    param_aggregation = NULL
)
Arguments:
model_name string Name of the model that generates this embedding.
model_label string Label of the model that generates this embedding.
model_date string Date when the embedding generating model was created.
model_method string Method of the underlying embedding model.
model_version string Version of the model that generated this embedding.
model_language string Language of the model that generated this embedding.
param_seq_length int Maximum number of tokens that processes the generating model for a
chunk.
param_chunks int Maximum number of chunks which are supported by the generating model.
param_features int Number of dimensions of the text embeddings.
param_overlap int Number of tokens that were added at the beginning of the sequence for
the next chunk by this model.
param_emb_layer_min int or string determining the first layer to be included in the creation
of embeddings.
param_emb_layer_max int or string determining the last layer to be included in the creation
of embeddings.
param_emb_pool_type string determining the method for pooling the token embeddings within
each layer.
param_pad_value int Value indicating padding. This value should no be in the range of
regular values for computations. Thus it is not recommended to chance this value. De-
fault is -100. Allowed values: x <= -100
param_aggregation string Aggregation method of the hidden states. Deprecated. Only in-
cluded for backward compatibility.

```

Returns: The method returns a new object of this class.

Method is_configured(): Method for checking if the model was successfully configured. An object can only be used if this value is TRUE.

Usage:

```
LargeDataSetForTextEmbeddings$is_configured()
```

Returns: bool TRUE if the model is fully configured. FALSE if not.

Method get_text_embedding_model_name(): Method for requesting the name (unique id) of the underlying text embedding model.

Usage:

```
LargeDataSetForTextEmbeddings$get_text_embedding_model_name()
```

Returns: Returns a string describing name of the text embedding model.

Method get_model_info(): Method for retrieving information about the model that generated this embedding.

Usage:

```
LargeDataSetForTextEmbeddings$get_model_info()
```

Returns: list containing all saved information about the underlying text embedding model.

Method `load_from_disk()`: loads an object of class `LargeDataSetForTextEmbeddings` from disk and updates the object to the current version of the package.

Usage:

```
LargeDataSetForTextEmbeddings$load_from_disk(dir_path)
```

Arguments:

`dir_path` Path where the data set set is stored.

Returns: Method does not return anything. It loads an object from disk.

Method `get_model_label()`: Method for retrieving the label of the model that generated this embedding.

Usage:

```
LargeDataSetForTextEmbeddings$get_model_label()
```

Returns: string Label of the corresponding text embedding model

Method `add_feature_extractor_info()`: Method setting information on the `TEFeatureExtractor` that was used to reduce the number of dimensions of the text embeddings. This information should only be used if a `TEFeatureExtractor` was applied.

Usage:

```
LargeDataSetForTextEmbeddings$add_feature_extractor_info(
  model_name,
  model_label = NA,
  features = NA,
  method = NA,
  noise_factor = NA,
  optimizer = NA
)
```

Arguments:

`model_name` string Name of the underlying `TextEmbeddingModel`.

`model_label` string Label of the underlying `TextEmbeddingModel`.

`features` int Number of dimension (features) for the **compressed** text embeddings.

`method` string Method that the `TEFeatureExtractor` applies for genereating the compressed text embeddings.

`noise_factor` double Noise factor of the `TEFeatureExtractor`.

`optimizer` string Optimizer used during training the `TEFeatureExtractor`.

Returns: Method does nothing return. It sets information on a `TEFeatureExtractor`.

Method `get_feature_extractor_info()`: Method for receiving information on the `TEFeatureExtractor` that was used to reduce the number of dimensions of the text embeddings.

Usage:

```
LargeDataSetForTextEmbeddings$get_feature_extractor_info()
```

Returns: Returns a list with information on the [TEFeatureExtractor](#). If no [TEFeatureExtractor](#) was used it returns NULL.

Method `is_compressed()`: Checks if the text embedding were reduced by a [TEFeatureExtractor](#).

Usage:

```
LargeDataSetForTextEmbeddings$is_compressed()
```

Returns: Returns TRUE if the number of dimensions was reduced by a [TEFeatureExtractor](#). If not return FALSE.

Method `get_times()`: Number of chunks/times of the text embeddings.

Usage:

```
LargeDataSetForTextEmbeddings$get_times()
```

Returns: Returns an int describing the number of chunks/times of the text embeddings.

Method `get_features()`: Number of actual features/dimensions of the text embeddings. In the case a [TEFeatureExtractor](#) was used the number of features is smaller as the original number of features. To receive the original number of features (the number of features before applying a [TEFeatureExtractor](#)) you can use the method `get_original_features` of this class.

Usage:

```
LargeDataSetForTextEmbeddings$get_features()
```

Returns: Returns an int describing the number of features/dimensions of the text embeddings.

Method `get_original_features()`: Number of original features/dimensions of the text embeddings.

Usage:

```
LargeDataSetForTextEmbeddings$get_original_features()
```

Returns: Returns an int describing the number of features/dimensions if no [TEFeatureExtractor](#) is used or before a [TEFeatureExtractor](#) is applied.

Method `get_pad_value()`: Value for indicating padding.

Usage:

```
LargeDataSetForTextEmbeddings$get_pad_value()
```

Returns: Returns an int describing the value used for padding.

Method `add_embeddings_from_array()`: Method for adding new data to the data set from an array. Please note that the method does not check if cases already exist in the data set. To reduce the data set to unique cases call the method `reduce_to_unique_ids`.

Usage:

```
LargeDataSetForTextEmbeddings$add_embeddings_from_array(embedding_array)
```

Arguments:

`embedding_array` array containing the text embeddings.

Returns: The method does not return anything. It adds new data to the data set.

Method `add_embeddings_from_EMBEDDEDTEXT()`: Method for adding new data to the data set from an [EmbeddedText](#). Please note that the method does not check if cases already exist in the data set. To reduce the data set to unique cases call the method `reduce_to_unique_ids`.

Usage:

```
LargeDataSetForTextEmbeddings$add_embeddings_from_EMBEDDEDTEXT(EMBEDDEDTEXT)
```

Arguments:

`EMBEDDEDTEXT` Object of class [EmbeddedText](#).

Returns: The method does not return anything. It adds new data to the data set.

Method `add_embeddings_from_LargeDataSetForTextEmbeddings()`: Method for adding new data to the data set from an [LargeDataSetForTextEmbeddings](#). Please note that the method does not check if cases already exist in the data set. To reduce the data set to unique cases call the method `reduce_to_unique_ids`.

Usage:

```
LargeDataSetForTextEmbeddings$add_embeddings_from_LargeDataSetForTextEmbeddings(  
  dataset  
)
```

Arguments:

`dataset` Object of class [LargeDataSetForTextEmbeddings](#).

Returns: The method does not return anything. It adds new data to the data set.

Method `convert_to_EMBEDDEDTEXT()`: Method for converting this object to an object of class [EmbeddedText](#).

Attention This object uses memory mapping to allow the usage of data sets that do not fit into memory. By calling this method the data set will be loaded and stored into memory/RAM. This may lead to an out-of-memory error.

Usage:

```
LargeDataSetForTextEmbeddings$convert_to_EMBEDDEDTEXT()
```

Returns: `LargeDataSetForTextEmbeddings` an object of class [EmbeddedText](#) which is stored in the memory/RAM.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
LargeDataSetForTextEmbeddings$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

See Also

Other Data Management: [EmbeddedText](#), [LargeDataSetForText](#)

`load_all_py_scripts` *Load and re-load all python scripts*

Description

Function loads or re-loads all python scripts within the package 'aifedducation'.

Usage

```
load_all_py_scripts()
```

Value

Function does nothing return. It loads the requested scripts.

See Also

Other Utils Python Developers: [get_py_package_version\(\)](#), [get_py_package_versions\(\)](#), [load_py_scripts\(\)](#), [run_py_file\(\)](#)

`load_from_disk` *Loading objects created with 'aifedducation'*

Description

Function for loading objects created with 'aifedducation'.

Usage

```
load_from_disk(dir_path)
```

Arguments

`dir_path` string Path to the directory where the model is stored.

Value

Returns an object of class [TEClassifierRegular](#), [TEClassifierProtoNet](#), [TEFeatureExtractor](#), [TextEmbeddingModel](#), [LargeDataSetForTextEmbeddings](#), [LargeDataSetForText](#) or [EmbeddedText](#).

See Also

Other Saving and Loading: [save_to_disk\(\)](#)

load_py_scripts	<i>Load and re-load python scripts</i>
-----------------	--

Description

Function loads or re-loads python scripts within the package 'aifedducation'.

Usage

```
load_py_scripts(files)
```

Arguments

files vector containing the file names of the scripts that should be loaded.

Value

Function does nothing return. It loads the requested scripts.

See Also

Other Utils Python Developers: [get_py_package_version\(\)](#), [get_py_package_versions\(\)](#), [load_all_py_scripts\(\)](#), [run_py_file\(\)](#)

long_load_target_data	<i>Load target data for long running tasks</i>
-----------------------	--

Description

Function loads the target data for a long running task.

Usage

```
long_load_target_data(file_path, selectet_column)
```

Arguments

file_path string Path to the file storing the target data.

selectet_column string Name of the column containing the target data.

Details

This function assumes that the target data is stored as a columns with the cases in the rows and the categories in the columns. The ids of the cases must be stored in a column called "id".

Value

Returns a named factor containing the target data.

See Also

Other Utils Studio Developers: [add_missing_args\(\)](#), [create_data_embeddings_description\(\)](#), [summarize_args_for_long_task\(\)](#)

matrix_to_array_c *Reshape matrix to array*

Description

Function written in C++ for reshaping a matrix containing sequential data into an array for use with keras.

Usage

```
matrix_to_array_c(matrix, times, features)
```

Arguments

matrix	matrix containing the sequential data.
times	size_t Number of sequences.
features	size_t Number of features within each sequence.

Value

Returns an array. The first dimension corresponds to the cases, the second to the times, and the third to the features.

See Also

Other Utils Developers: [auto_n_cores\(\)](#), [create_object\(\)](#), [create_synthetic_units_from_matrix\(\)](#), [generate_id\(\)](#), [get_n_chunks\(\)](#), [get_synthetic_cases_from_matrix\(\)](#), [tensor_to_matrix_c\(\)](#), [to_categorical_c\(\)](#)

ModelsBasedOnTextEmbeddings

Base class for models using neural nets

Description

Abstract class for all models that do not rely on the python library 'transformers'. All models of this class require text embeddings as input. These are provided as objects of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#).

Objects of this class containing fields and methods used in several other classes in 'AI for Education'.

This class is **not** designed for a direct application and should only be used by developers.

Value

A new object of this class.

Super class

[aifedducation::AIFEBaseModel](#) -> ModelsBasedOnTextEmbeddings

Methods

Public methods:

- [ModelsBasedOnTextEmbeddings\\$get_text_embedding_model\(\)](#)
- [ModelsBasedOnTextEmbeddings\\$get_text_embedding_model_name\(\)](#)
- [ModelsBasedOnTextEmbeddings\\$check_embedding_model\(\)](#)
- [ModelsBasedOnTextEmbeddings\\$load_from_disk\(\)](#)
- [ModelsBasedOnTextEmbeddings\\$plot_training_history\(\)](#)
- [ModelsBasedOnTextEmbeddings\\$clone\(\)](#)

Method `get_text_embedding_model()`: Method for requesting the text embedding model information.

Usage:

`ModelsBasedOnTextEmbeddings$get_text_embedding_model()`

Returns: list of all relevant model information on the text embedding model underlying the model.

Method `get_text_embedding_model_name()`: Method for requesting the name (unique id) of the underlying text embedding model.

Usage:

`ModelsBasedOnTextEmbeddings$get_text_embedding_model_name()`

Returns: Returns a string describing name of the text embedding model.

Method `check_embedding_model()`: Method for checking if the provided text embeddings are created with the same [TextEmbeddingModel](#) as the model.

Usage:

```
ModelsBasedOnTextEmbeddings$check_embedding_model(text_embeddings)
```

Arguments:

`text_embeddings` Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#).

Returns: TRUE if the underlying [TextEmbeddingModel](#) are the same. FALSE if the models differ.

Method `load_from_disk()`: loads an object from disk and updates the object to the current version of the package.

Usage:

```
ModelsBasedOnTextEmbeddings$load_from_disk(dir_path)
```

Arguments:

`dir_path` Path where the object set is stored.

Returns: Method does not return anything. It loads an object from disk.

Method `plot_training_history()`: Method for requesting a plot of the training history. This method requires the *R* package 'ggplot2' to work.

Usage:

```
ModelsBasedOnTextEmbeddings$plot_training_history(
  final_training = FALSE,
  pl_step = NULL,
  measure = "loss",
  y_min = NULL,
  y_max = NULL,
  add_min_max = TRUE,
  text_size = 10
)
```

Arguments:

`final_training` bool If FALSE the values of the performance estimation are used. If TRUE only the epochs of the final training are used.

`pl_step` int Number of the step during pseudo labeling to plot. Only relevant if the model was trained with active pseudo labeling.

`measure` Measure to plot.

`y_min` Minimal value for the y-axis. Set to NULL for an automatic adjustment.

`y_max` Maximal value for the y-axis. Set to NULL for an automatic adjustment.

`add_min_max` bool If TRUE the minimal and maximal values during performance estimation are part of the plot. If FALSE only the mean values are shown. Parameter is ignored if `final_training`=TRUE.

`text_size` Size of the text.

Returns: Returns a plot of class ggplot visualizing the training process. Prepare history data of objects Function for preparing the history data of a model in order to be plotted in AI for Education - Studio.

final bool If TRUE the history data of the final training is used for the data set. pl_step int If use_pl=TRUE select the step within pseudo labeling for which the data should be prepared. Returns a named list with the training history data of the model. The reported measures depend on the provided model.

Utils Studio Developers internal

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
ModelsBasedOnTextEmbeddings$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

See Also

Other R6 Classes for Developers: [AIFEBaseModel](#), [ClassifiersBasedOnTextEmbeddings](#), [LargeDataSetBase](#), [TEClassifiersBasedOnProtoNet](#), [TEClassifiersBasedOnRegular](#)

output_message

Print message

Description

Prints a message `msg` if `trace` parameter is TRUE with current date with `message()` or `cat()` function.

Usage

```
output_message(msg, trace, msg_fun)
```

Arguments

<code>msg</code>	string Message that should be printed.
<code>trace</code>	bool Silent printing (FALSE) or not (TRUE).
<code>msg_fun</code>	bool value that determines what function should be used. TRUE for <code>message()</code> , FALSE for <code>cat()</code> .

Value

This function returns nothing.

See Also

Other Utils Log Developers: [cat_message\(\)](#), [clean_pytorch_log_transformers\(\)](#), [print_message\(\)](#), [read_log\(\)](#), [read_loss_log\(\)](#), [reset_log\(\)](#), [reset_loss_log\(\)](#), [write_log\(\)](#)

prepare_r_array_for_dataset*Convert R array for arrow data set***Description**

Function converts a R array into a numpy array that can be added to an arrow data set. The array should represent embeddings.

Usage

```
prepare_r_array_for_dataset(r_array)
```

Arguments

r_array	array representing embeddings.
---------	--------------------------------

Value

Returns a numpy array.

See Also

Other Utils Python Data Management Developers: [class_vector_to_py_dataset\(\)](#), [data.frame_to_py_dataset\(\)](#), [get_batches_index\(\)](#), [py_dataset_to_embeddings\(\)](#), [reduce_to_unique\(\)](#), [tensor_list_to_numpy\(\)](#), [tensor_to_numpy\(\)](#)

prepare_session*Function for setting up a python environment within R.***Description**

This functions checks for python and a specified environment. If the environment exists it will be activated. If python is already initialized it uses the current environment.

Usage

```
prepare_session(env_type = "auto", envname = "aifedducation")
```

Arguments

env_type	string If set to "venv" virtual environment is requested. If set to "conda" a 'conda' environment is requested. If set to "auto" the function tries to activate a virtual environment with the given name. If this environment does not exist it tries to activate a conda environment with the given name. If this fails the default virtual environment is used.
envname	string envname name of the requested environment.

Value

Function does not return anything. It is used for preparing python and R.

See Also

Other Installation and Configuration: [check_aif_py_modules\(\)](#), [install_aifedducation\(\)](#), [install_aifedducation_stu](#)
[install_py_modules\(\)](#), [set_transformers_logger\(\)](#), [update_aifedducation\(\)](#)

print_message

Print message (message())

Description

Prints a message msg if trace parameter is TRUE with current date with message() function.

Usage

```
print_message(msg, trace)
```

Arguments

msg	string Message that should be printed.
trace	bool Silent printing (FALSE) or not (TRUE).

Value

This function returns nothing.

See Also

Other Utils Log Developers: [cat_message\(\)](#), [clean_pytorch_log_transformers\(\)](#), [output_message\(\)](#),
[read_log\(\)](#), [read_loss_log\(\)](#), [reset_log\(\)](#), [reset_loss_log\(\)](#), [write_log\(\)](#)

py_dataset_to_embeddings

Convert arrow data set to an arrow data set

Description

Function for converting an arrow data set into a data set that can be used to store and process embeddings.

Usage

```
py_dataset_to_embeddings(py_dataset)
```

Arguments

py_dataset Object of class `datasets.arrow_dataset.Dataset`.

Value

Returns the data set of class `datasets.arrow_dataset.Dataset` with only two columns ("id", "input"). "id" stores the name of the cases while "input" stores the embeddings.

See Also

Other Utils Python Data Management Developers: [class_vector_to_py_dataset\(\)](#), [data.frame_to_py_dataset\(\)](#), [get_batches_index\(\)](#), [prepare_r_array_for_dataset\(\)](#), [reduce_to_unique\(\)](#), [tensor_list_to_numpy\(\)](#), [tensor_to_numpy\(\)](#)

random_bool_on_CI *Random bool on Continuous Integration*

Description

Function returns randomly TRUE or FALSE if on CI. It returns FALSE if it is not on CI.

Usage

`random_bool_on_CI()`

Value

Returns a bool.

See Also

Other Utils TestThat Developers: [check_adjust_n_samples_on_CI\(\)](#), [generate_args_for_tests\(\)](#), [generate_embeddings\(\)](#), [generate_tensors\(\)](#), [get_current_args_for_print\(\)](#), [get_fixed_test_tensor\(\)](#), [get_test_data_for_classifiers\(\)](#)

read_log*Function for reading a log file in R*

Description

This function reads a log file at the given location. The log file should be created with [write_log](#).

Usage

```
read_log(file_path)
```

Arguments

file_path string Path to the log file.

Value

Returns a matrix containing the log file.

See Also

Other Utils Log Developers: [cat_message\(\)](#), [clean_pytorch_log_transformers\(\)](#), [output_message\(\)](#), [print_message\(\)](#), [read_loss_log\(\)](#), [reset_log\(\)](#), [reset_loss_log\(\)](#), [write_log\(\)](#)

read_loss_log*Function for reading a log file containing a record of the loss during training.*

Description

This function reads a log file that contains values for every epoch for the loss. The values are grouped for training and validation data. The log contains values for test data if test data was available during training.

Usage

```
read_loss_log(path_loss)
```

Arguments

path_loss string Path to the log file.

Details

In general the loss is written by a python function during model's training.

Value

Function returns a matrix that contains two or three row depending on the data inside the loss log. In the case of two rows the first represents the training data and the second the validation data. In the case of three rows the third row represents the values for test data. All Columns represent the epochs.

See Also

Other Utils Log Developers: [cat_message\(\)](#), [clean_pytorch_log_transformers\(\)](#), [output_message\(\)](#), [print_message\(\)](#), [read_log\(\)](#), [reset_log\(\)](#), [reset_loss_log\(\)](#), [write_log\(\)](#)

reduce_to_unique	<i>Reduce to unique cases</i>
------------------	-------------------------------

Description

Function creates an arrow data set that contains only unique cases. That is, duplicates are removed.

Usage

```
reduce_to_unique(dataset_to_reduce, column_name)
```

Arguments

dataset_to_reduce
Object of class `datasets.arrow_dataset.Dataset`.

column_name
string Name of the column whose values should be unique.

Value

Returns a data set of class `datasets.arrow_dataset.Dataset` where the duplicates are removed according to the given column.

See Also

Other Utils Python Data Management Developers: [class_vector_to_py_dataset\(\)](#), [data.frame_to_py_dataset\(\)](#), [get_batches_index\(\)](#), [prepare_r_array_for_dataset\(\)](#), [py_dataset_to_embeddings\(\)](#), [tensor_list_to_numpy\(\)](#), [tensor_to_numpy\(\)](#)

reset_log	<i>Function that resets a log file.</i>
-----------	---

Description

This function writes a log file with default values. The file can be read with [read_log](#).

Usage

```
reset_log(log_path)
```

Arguments

log_path string Path to the log file.

Value

Function does nothing return. It is used to write an "empty" log file.

See Also

Other Utils Log Developers: [cat_message\(\)](#), [clean_pytorch_log_transformers\(\)](#), [output_message\(\)](#), [print_message\(\)](#), [read_log\(\)](#), [read_loss_log\(\)](#), [reset_loss_log\(\)](#), [write_log\(\)](#)

reset_loss_log	<i>Reset log for loss information</i>
----------------	---------------------------------------

Description

This function writes an empty log file for loss information.

Usage

```
reset_loss_log(log_path, epochs)
```

Arguments

log_path string Path to the log file.
epochs int Number of epochs for the complete training process.

Value

Function does nothing return. It writes a log file at the given location. The file is a .csv file that contains three rows. The first row takes the value for the training, the second for the validation, and the third row for the test data. The columns represent epochs.

See Also

Other Utils Log Developers: [cat_message\(\)](#), [clean_pytorch_log_transformers\(\)](#), [output_message\(\)](#), [print_message\(\)](#), [read_log\(\)](#), [read_loss_log\(\)](#), [reset_log\(\)](#), [write_log\(\)](#)

run_py_file

*Run python file***Description**

Used to run python files with `reticulate::py_run_file()` from folder `python`.

Usage

```
run_py_file(py_file_name)
```

Arguments

`py_file_name` string Name of a python file to run. The file must be in the `python` folder of `aifedducation` package.

Value

This function returns nothing.

See Also

Other Utils Python Developers: [get_py_package_version\(\)](#), [get_py_package_versions\(\)](#), [load_all_py_scripts\(\)](#), [load_py_scripts\(\)](#)

save_to_disk

*Saving objects created with 'aifedducation'***Description**

Function for saving objects created with 'aifedducation'.

Usage

```
save_to_disk(object, dir_path, folder_name)
```

Arguments

<code>object</code>	Object of class <code>TEClassifierRegular</code> , <code>TEClassifierProtoNet</code> , <code>TEFeatureExtractor</code> , <code>TextEmbeddingModel</code> , <code>LargeDataSetForTextEmbeddings</code> , <code>LargeDataSetForText</code> or <code>EmbeddedText</code> which should be saved.
<code>dir_path</code>	string Path to the directory where the model is stored.
<code>folder_name</code>	string Name of the folder where the files should be stored.

Value

Function does not return a value. It saves the model to disk.

No return value, called for side effects.

See Also

Other Saving and Loading: [load_from_disk\(\)](#)

set_transformers_logger

Sets the level for logging information of the 'transformers' library

Description

This function changes the level for logging information of the 'transformers' library. It influences the output printed to console for creating and training transformer models as well as [TextEmbeddingModels](#).

Usage

```
set_transformers_logger(level = "ERROR")
```

Arguments

level	string Minimal level that should be printed to console. Four levels are available: INFO, WARNING, ERROR and DEBUG
-------	---

Value

This function does not return anything. It is used for its side effects.

See Also

Other Installation and Configuration: [check_aif_py_modules\(\)](#), [install_aifedducation\(\)](#), [install_aifedducation_stu](#)
[install_py_modules\(\)](#), [prepare_session\(\)](#), [update_aifedducation\(\)](#)

start_aifedducation_studio
Aifedducation Studio

Description

Functions starts a shiny app that represents Aifedducation Studio.

Usage

```
start_aifedducation_studio()
```

Value

This function does nothing return. It is used to start a shiny app.

summarize_args_for_long_task
Summarize arguments from shiny input

Description

This function extracts the input relevant for a specific method of a specific class from shiny input.

In addition, it adds the path to all objects which can not be exported to another R session. These object must be loaded separately in the new session with the function `add_missing_args`. The paths are intended to be used with `shiny::ExtendedTask`. The final preparation of the arguments should be done with

The function can also be used to override the default value of a method or to add value for arguments which are not part of shiny input (use parameter `override_args`).

Usage

```
summarize_args_for_long_task(  

  input,  

  object_class,  

  method = "configure",  

  path_args = list(path_to_embeddings = NULL, path_to_textual_dataset = NULL,  

    path_to_target_data = NULL, path_to_feature_extractor = NULL, destination_path =  

    NULL, folder_name = NULL),  

  override_args = list(),  

  meta_args = list(py_environment_type = get_py_env_type(), py_env_name =  

    get_py_env_name(), target_data_column = input$data_target_column, object_class =  

    input$classifier_type)  

)
```

Arguments

<code>input</code>	Shiny input.
<code>object_class</code>	<code>string</code> Class of the object.
<code>method</code>	<code>string</code> Method of the class for which the arguments should be extracted and prepared.
<code>path_args</code>	<code>list</code> List containing the path to object that can not be exported to another R session. These must be loaded in the session.
<code>override_args</code>	<code>list</code> List containing all arguments that should be set manually. The values override default values of the argument and values which are part of <code>input</code> .
<code>meta_args</code>	<code>list</code> List containing information that are not relevant for the arguments of the method but are necessary to set up the <code>shiny::ExtendedTask</code> correctly.

Value

Returns a named list with the following entries:

- `args`: Named list of all arguments necessary for the method of the class.
- `path_args`: Named list of all paths for loading the objects missing in `args`.
- `meta_args`: Named list of all arguments that are not part of the arguments of the method but which are necessary to set up the `shiny::ExtendedTask` correctly.

Note

Please note that all list are named list of the format (`argument_name=values`).

See Also

Other Utils Studio Developers: [add_missing_args\(\)](#), [create_data_embeddings_description\(\)](#), [long_load_target_data\(\)](#)

TEClassifierParallel *Text embedding classifier with a neural net*

Description

Classification Type

This is a probability classifier that predicts a probability distribution for different classes/categories. This is the standard case most common in literature.

Parallel Core Architecture

This model is based on a parallel architecture. An input is passed to different types of layers separately. At the end the outputs are combined to create the final output of the whole model.

Transformer Encoder Layers

Description

The transformer encoder layers follow the structure of the encoder layers used in transformer models. A single layer is designed as described by Chollet, Kalinowski, and Allaire (2022, p. 373) with the exception that single components of the layers (such as the activation function, the kind of residual connection, the kind of normalization or the kind of attention) can be customized. All parameters with the prefix *tf_* can be used to configure this layer.

Feature Layer

Description

The feature layer is a dense layer that can be used to increase or decrease the number of features of the input data before passing the data into your model. The aim of this layer is to increase or reduce the complexity of the data for your model. The output size of this layer determines the number of features for all following layers. In the special case that the requested number of features equals the number of features of the text embeddings this layer is reduced to a dropout layer with masking capabilities. All parameters with the prefix *feat_* can be used to configure this layer.

Dense Layers

Description

A fully connected layer. The layer is applied to every step of a sequence. All parameters with the prefix *dense_* can be used to configure this layer.

Multiple N-Gram Layers

Description

This type of layer focuses on sub-sequence and performs an 1d convolutional operation. On a word and token level these sub-sequences can be interpreted as n-grams (Jacovi, Shalom & Goldberg 2018). The convolution is done across all features. The number of filters equals the number of features of the input tensor. Thus, the shape of the tensor is retained (Pham, Kruszewski & Boleda 2016).

The layer is able to consider multiple n-grams at the same time. In this case the convolution of the n-grams is done separately and the resulting tensors are concatenated along the feature dimension. The number of filters for every n-gram is set to *num_features/num_n-grams*. Thus, the resulting tensor has the same shape as the input tensor.

Sub-sequences that are masked in the input are also masked in the output.

The output of this layer can be understand as the results of the n-gram filters. Stacking this layer allows the model to perform n-gram detection of n-grams (meta perspective). All parameters with the prefix *ng_conv_* can be used to configure this layer.

Recurrent Layers

Description

A regular recurrent layer either as Gated Recurrent Unit (GRU) or Long Short-Term Memory (LSTM) layer. Uses PyTorchs implementation. All parameters with the prefix *rec_* can be used to configure this layer.

Merge Layer

Description

Layer for combining the output of different layers. All inputs must be sequential data of shape (Batch, Times, Features). First, pooling over time is applied extracting the minimal and/or maximal features. Second, the pooled tensors are combined by calculating their weighted sum. Different attention mechanism can be used to dynamically calculate the corresponding weights. This allows

the model to decide which part of the data is most useful. Finally, pooling over features is applied extracting a specific number of maximal and/or minimal features. A normalization of all input at the beginning of the layer is possible. All parameters with the prefix `merge_` can be used to configure this layer.

Training and Prediction

For the creation and training of a classifier an object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` on the one hand and a `factor` on the other hand are necessary.

The object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` contains the numerical text representations (text embeddings) of the raw texts generated by an object of class `TextEmbeddingModel`. For supporting large data sets it is recommended to use `LargeDataSetForTextEmbeddings` instead of `EmbeddedText`.

The `factor` contains the classes/categories for every text. Missing values (unlabeled cases) are supported and can be used for pseudo labeling.

For predictions an object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` has to be used which was created with the same `TextEmbeddingModel` as for training.

Value

Returns a new object of this class ready for configuration or for loading a saved classifier.

Super classes

```
aifedducation::AIFEBaseModel -> aifedducation::ModelsBasedOnTextEmbeddings -> aifedducation::ClassifiersB
-> aifedducation::TEClassifiersBasedOnRegular -> TEClassifierParallel
```

Methods

Public methods:

- `TEClassifierParallel$configure()`
- `TEClassifierParallel$clone()`

Method `configure()`: Creating a new instance of this class.

Usage:

```
TEClassifierParallel$configure(
  name = NULL,
  label = NULL,
  text_embeddings = NULL,
  feature_extractor = NULL,
  target_levels = NULL,
  shared_feat_layer = TRUE,
  feat_act_fct = "ELU",
  feat_size = 50,
  feat_bias = TRUE,
  feat_dropout = 0,
  feat_parametrizations = "None",
  feat_normalization_type = "LayerNorm",
  ng_conv_act_fct = "ELU",
```

```

ng_conv_n_layers = 1,
ng_conv_ks_min = 2,
ng_conv_ks_max = 4,
ng_conv_bias = FALSE,
ng_conv_dropout = 0.1,
ng_conv_parametrizations = "None",
ng_conv_normalization_type = "LayerNorm",
ng_conv_residual_type = "ResidualGate",
dense_act_fct = "ELU",
dense_n_layers = 1,
dense_dropout = 0.5,
dense_bias = FALSE,
dense_parametrizations = "None",
dense_normalization_type = "LayerNorm",
dense_residual_type = "ResidualGate",
rec_act_fct = "Tanh",
rec_n_layers = 1,
rec_type = "GRU",
rec_bidirectional = FALSE,
rec_dropout = 0.2,
rec_bias = FALSE,
rec_parametrizations = "None",
rec_normalization_type = "LayerNorm",
rec_residual_type = "ResidualGate",
tf_act_fct = "ELU",
tf_dense_dim = 50,
tf_n_layers = 1,
tf_dropout_rate_1 = 0.1,
tf_dropout_rate_2 = 0.5,
tf_attention_type = "MultiHead",
tf_positional_type = "absolute",
tf_num_heads = 1,
tf_bias = FALSE,
tf_parametrizations = "None",
tf_normalization_type = "LayerNorm",
tf_residual_type = "ResidualGate",
merge_attention_type = "multi_head",
merge_num_heads = 1,
merge_normalization_type = "LayerNorm",
merge_pooling_features = 50,
merge_pooling_type = "MinMax"
)

```

Arguments:

name string Name of the new model. Please refer to common name conventions. Free text can be used with parameter label. If set to NULL a unique ID is generated automatically.

Allowed values: any

label string Label for the new model. Here you can use free text. Allowed values: any

text_embeddings EmbeddedText, LargeDataSetForTextEmbeddings Object of class [Em-](#)

`beddedText` or `LargeDataSetForTextEmbeddings`.

`feature_extractor` `TEFeatureExtractor` Object of class `TEFeatureExtractor` which should be used in order to reduce the number of dimensions of the text embeddings. If no feature extractor should be applied set `NULL`.

`target_levels` `vector` containing the levels (categories or classes) within the target data. Please note that order matters. For ordinal data please ensure that the levels are sorted correctly with later levels indicating a higher category/class. For nominal data the order does not matter.

`shared_feat_layer` `bool` If `TRUE` all streams use the same feature layer. If `FALSE` all streams use their own feature layer.

`feat_act_fct` `string` Activation function for all layers. Allowed values: `'ELU'`, `'LeakyReLU'`, `'ReLU'`, `'GELU'`, `'Sigmoid'`, `'Tanh'`, `'PReLU'`

`feat_size` `int` Number of neurons for each dense layer. Allowed values: $2 \leq x$

`feat_bias` `bool` If `TRUE` a bias term is added to all layers. If `FALSE` no bias term is added to the layers.

`feat_dropout` `double` determining the dropout for the dense projection of the feature layer. Allowed values: $0 \leq x \leq 0.6$

`feat_parametrizations` `string` Re-Parametrizations of the weights of layers. Allowed values: `'None'`, `'OrthogonalWeights'`, `'WeightNorm'`, `'SpectralNorm'`

`feat_normalization_type` `string` Type of normalization applied to all layers and stack layers. Allowed values: `'LayerNorm'`, `'None'`

`ng_conv_act_fct` `string` Activation function for all layers. Allowed values: `'ELU'`, `'LeakyReLU'`, `'ReLU'`, `'GELU'`, `'Sigmoid'`, `'Tanh'`, `'PReLU'`

`ng_conv_n_layers` `int` determining how many times the n-gram layers should be added to the network. Allowed values: $0 \leq x$

`ng_conv_ks_min` `int` determining the minimal window size for n-grams. Allowed values: $2 \leq x$

`ng_conv_ks_max` `int` determining the maximal window size for n-grams. Allowed values: $2 \leq x$

`ng_conv_bias` `bool` If `TRUE` a bias term is added to all layers. If `FALSE` no bias term is added to the layers.

`ng_conv_dropout` `double` determining the dropout for n-gram convolution layers. Allowed values: $0 \leq x \leq 0.6$

`ng_conv_parametrizations` `string` Re-Parametrizations of the weights of layers. Allowed values: `'None'`, `'OrthogonalWeights'`, `'WeightNorm'`, `'SpectralNorm'`

`ng_conv_normalization_type` `string` Type of normalization applied to all layers and stack layers. Allowed values: `'LayerNorm'`, `'None'`

`ng_conv_residual_type` `string` Type of residual connection for all layers and stack of layers. Allowed values: `'ResidualGate'`, `'Addition'`, `'None'`

`dense_act_fct` `string` Activation function for all layers. Allowed values: `'ELU'`, `'LeakyReLU'`, `'ReLU'`, `'GELU'`, `'Sigmoid'`, `'Tanh'`, `'PReLU'`

`dense_n_layers` `int` Number of dense layers. Allowed values: $0 \leq x$

`dense_dropout` `double` determining the dropout between dense layers. Allowed values: $0 \leq x \leq 0.6$

`dense_bias` `bool` If `TRUE` a bias term is added to all layers. If `FALSE` no bias term is added to the layers.

```

dense_parametrizations string Re-Parametrizations of the weights of layers. Allowed val-
    ues: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
dense_normalization_type string Type of normalization applied to all layers and stack lay-
    ers. Allowed values: 'LayerNorm', 'None'
dense_residual_type string Type of residual connection for all layers and stack of layers.
    Allowed values: 'ResidualGate', 'Addition', 'None'
rec_act_fct string Activation function for all layers. Allowed values: 'Tanh'
rec_n_layers int Number of recurrent layers. Allowed values: 0 <= x
rec_type string Type of the recurrent layers. rec_type='GRU' for Gated Recurrent Unit and
    rec_type='LSTM' for Long Short-Term Memory. Allowed values: 'GRU', 'LSTM'
rec_bidirectional bool If TRUE a bidirectional version of the recurrent layers is used.
rec_dropout double determining the dropout between recurrent layers. Allowed values: 0 <= x <= 0.6
rec_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the
    layers.
rec_parametrizations string Re-Parametrizations of the weights of layers. Allowed val-
    ues: 'None'
rec_normalization_type string Type of normalization applied to all layers and stack layers.
    Allowed values: 'LayerNorm', 'None'
rec_residual_type string Type of residual connection for all layers and stack of layers.
    Allowed values: 'ResidualGate', 'Addition', 'None'
tf_act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU',
    'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'
tf_dense_dim int determining the size of the projection layer within a each transformer en-
    coder. Allowed values: 1 <= x
tf_n_layers int determining how many times the encoder should be added to the network.
    Allowed values: 0 <= x
tf_dropout_rate_1 double determining the dropout after the attention mechanism within the
    transformer encoder layers. Allowed values: 0 <= x <= 0.6
tf_dropout_rate_2 double determining the dropout for the dense projection within the trans-
    former encoder layers. Allowed values: 0 <= x <= 0.6
tf_attention_type string Choose the attention type. Allowed values: 'Fourier', 'Multi-
    Head'
tf_positional_type string Type of processing positional information. Allowed values: 'ab-
    solute'
tf_num_heads int determining the number of attention heads for a self-attention layer. Only
    relevant if attention_type='multihead' Allowed values: 0 <= x
tf_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the
    layers.
tf_parametrizations string Re-Parametrizations of the weights of layers. Allowed values:
    'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
tf_normalization_type string Type of normalization applied to all layers and stack layers.
    Allowed values: 'LayerNorm', 'None'
tf_residual_type string Type of residual connection for all layers and stack of layers.
    Allowed values: 'ResidualGate', 'Addition', 'None'

```

```

merge_attention_type string Choose the attention type. Allowed values: 'Fourier', 'MultiHead'
merge_num_heads int determining the number of attention heads for a self-attention layer.
    Only relevant if attention_type='multihead' Allowed values: 0 <= x
merge_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'
merge_pooling_features int Number of features to be extracted at the end of the model.
    Allowed values: 1 <= x
merge_pooling_type string Type of extracting intermediate features. Allowed values: 'Max', 'Min', 'MinMax'
```

Returns: Function does nothing return. It modifies the current object.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
TEClassifierParallel$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

See Also

Other Classification: [TEClassifierParallelPrototype](#), [TEClassifierProtoNet](#), [TEClassifierRegular](#), [TEClassifierSequential](#), [TEClassifierSequentialPrototype](#)

TEClassifierParallelPrototype

Text embedding classifier with a ProtoNet

Description

Classification Type

This object is a metric based classifier and represents in implementation of a prototypical network for few-shot learning as described by Snell, Swersky, and Zemel (2017). The network uses a multi way contrastive loss described by Zhang et al. (2019). The network learns to scale the metric as described by Oreshkin, Rodriguez, and Lacoste (2018).

Parallel Core Architecture

This model is based on a parallel architecture. An input is passed to different types of layers separately. At the end the outputs are combined to create the final output of the whole model.

Transformer Encoder Layers

Description

The transformer encoder layers follow the structure of the encoder layers used in transformer models. A single layer is designed as described by Chollet, Kalinowski, and Allaire (2022, p. 373) with the exception that single components of the layers (such as the activation function, the kind

of residual connection, the kind of normalization or the kind of attention) can be customized. All parameters with the prefix *tf_* can be used to configure this layer.

Feature Layer

Description

The feature layer is a dense layer that can be used to increase or decrease the number of features of the input data before passing the data into your model. The aim of this layer is to increase or reduce the complexity of the data for your model. The output size of this layer determines the number of features for all following layers. In the special case that the requested number of features equals the number of features of the text embeddings this layer is reduced to a dropout layer with masking capabilities. All parameters with the prefix *feat_* can be used to configure this layer.

Dense Layers

Description

A fully connected layer. The layer is applied to every step of a sequence. All parameters with the prefix *dense_* can be used to configure this layer.

Multiple N-Gram Layers

Description

This type of layer focuses on sub-sequence and performs an 1d convolutional operation. On a word and token level these sub-sequences can be interpreted as n-grams (Jacovi, Shalom & Goldberg 2018). The convolution is done across all features. The number of filters equals the number of features of the input tensor. Thus, the shape of the tensor is retained (Pham, Kruszewski & Boleda 2016).

The layer is able to consider multiple n-grams at the same time. In this case the convolution of the n-grams is done separately and the resulting tensors are concatenated along the feature dimension. The number of filters for every n-gram is set to *num_features/num_n-grams*. Thus, the resulting tensor has the same shape as the input tensor.

Sub-sequences that are masked in the input are also masked in the output.

The output of this layer can be understand as the results of the n-gram filters. Stacking this layer allows the model to perform n-gram detection of n-grams (meta perspective). All parameters with the prefix *ng_conv_* can be used to configure this layer.

Recurrent Layers

Description

A regular recurrent layer either as Gated Recurrent Unit (GRU) or Long Short-Term Memory (LSTM) layer. Uses PyTorchs implementation. All parameters with the prefix *rec_* can be used to configure this layer.

Merge Layer

Description

Layer for combining the output of different layers. All inputs must be sequential data of shape (Batch, Times, Features). First, pooling over time is applied extracting the minimal and/or maximal features. Second, the pooled tensors are combined by calculating their weighted sum. Different attention mechanism can be used to dynamically calculate the corresponding weights. This allows the model to decide which part of the data is most usefull. Finally, pooling over features is applied extracting a specific number of maximal and/or minimal features. A normalization of all input at

the begining of the layer is possible. All parameters with the prefix `merge_` can be used to configure this layer.

Training and Prediction

For the creation and training of a classifier an object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` on the one hand and a `factor` on the other hand are necessary.

The object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` contains the numerical text representations (text embeddings) of the raw texts generated by an object of class `TextEmbeddingModel`. For supporting large data sets it is recommended to use `LargeDataSetForTextEmbeddings` instead of `EmbeddedText`.

The factor contains the classes/categories for every text. Missing values (unlabeled cases) are supported and can be used for pseudo labeling.

For predictions an object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` has to be used which was created with the same `TextEmbeddingModel` as for training.

Value

Returns a new object of this class ready for configuration or for loading a saved classifier.

Super classes

```
aifedducation::AIFEBaseModel -> aifedducation::ModelsBasedOnTextEmbeddings -> aifedducation::ClassifiersB-> aifedducation::TEClassifiersBasedOnProtoNet -> TEClassifierParallelPrototype
```

Methods

Public methods:

- `TEClassifierParallelPrototype$configure()`
- `TEClassifierParallelPrototype$clone()`

Method `configure()`: Creating a new instance of this class.

Usage:

```
TEClassifierParallelPrototype$configure(
  name = NULL,
  label = NULL,
  text_embeddings = NULL,
  feature_extractor = NULL,
  target_levels = NULL,
  metric_type = "Euclidean",
  shared_feat_layer = TRUE,
  feat_act_fct = "ELU",
  feat_size = 50,
  feat_bias = TRUE,
  feat_dropout = 0,
  feat_parametrizations = "None",
  feat_normalization_type = "LayerNorm",
  ng_conv_act_fct = "ELU",
  ng_conv_n_layers = 1,
```

```

ng_conv_ks_min = 2,
ng_conv_ks_max = 4,
ng_conv_bias = FALSE,
ng_conv_dropout = 0.1,
ng_conv_parametrizations = "None",
ng_conv_normalization_type = "LayerNorm",
ng_conv_residual_type = "ResidualGate",
dense_act_fct = "ELU",
dense_n_layers = 1,
dense_dropout = 0.5,
dense_bias = FALSE,
dense_parametrizations = "None",
dense_normalization_type = "LayerNorm",
dense_residual_type = "ResidualGate",
rec_act_fct = "Tanh",
rec_n_layers = 1,
rec_type = "GRU",
rec_bidirectional = FALSE,
rec_dropout = 0.2,
rec_bias = FALSE,
rec_parametrizations = "None",
rec_normalization_type = "LayerNorm",
rec_residual_type = "ResidualGate",
tf_act_fct = "ELU",
tf_dense_dim = 50,
tf_n_layers = 1,
tf_dropout_rate_1 = 0.1,
tf_dropout_rate_2 = 0.5,
tf_attention_type = "MultiHead",
tf_positional_type = "absolute",
tf_num_heads = 1,
tf_bias = FALSE,
tf_parametrizations = "None",
tf_normalization_type = "LayerNorm",
tf_residual_type = "ResidualGate",
merge_attention_type = "multi_head",
merge_num_heads = 1,
merge_normalization_type = "LayerNorm",
merge_pooling_features = 50,
merge_pooling_type = "MinMax",
embedding_dim = 2
)

```

Arguments:

name string Name of the new model. Please refer to common name conventions. Free text can be used with parameter label. If set to NULL a unique ID is generated automatically.

Allowed values: any

label string Label for the new model. Here you can use free text. Allowed values: any

text_embeddings EmbeddedText, LargeDataSetForTextEmbeddings Object of class [Em-](#)

beddedText or **LargeDataSetForTextEmbeddings**.

feature_extractor `TEFeatureExtractor` Object of class `TEFeatureExtractor` which should be used in order to reduce the number of dimensions of the text embeddings. If no feature extractor should be applied set `NULL`.

target_levels `vector` containing the levels (categories or classes) within the target data. Please note that order matters. For ordinal data please ensure that the levels are sorted correctly with later levels indicating a higher category/class. For nominal data the order does not matter.

metric_type `string` Type of metric used for calculating the distance. Allowed values: 'Euclidean'

shared_feat_layer `bool` If `TRUE` all streams use the same feature layer. If `FALSE` all streams use their own feature layer.

feat_act_fct `string` Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU', 'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'

feat_size `int` Number of neurons for each dense layer. Allowed values: $2 \leq x$

feat_bias `bool` If `TRUE` a bias term is added to all layers. If `FALSE` no bias term is added to the layers.

feat_dropout `double` determining the dropout for the dense projection of the feature layer. Allowed values: $0 \leq x \leq 0.6$

feat_parametrizations `string` Re-Parametrizations of the weights of layers. Allowed values: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'

feat_normalization_type `string` Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'

ng_conv_act_fct `string` Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU', 'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'

ng_conv_n_layers `int` determining how many times the n-gram layers should be added to the network. Allowed values: $0 \leq x$

ng_conv_ks_min `int` determining the minimal window size for n-grams. Allowed values: $2 \leq x$

ng_conv_ks_max `int` determining the maximal window size for n-grams. Allowed values: $2 \leq x$

ng_conv_bias `bool` If `TRUE` a bias term is added to all layers. If `FALSE` no bias term is added to the layers.

ng_conv_dropout `double` determining the dropout for n-gram convolution layers. Allowed values: $0 \leq x \leq 0.6$

ng_conv_parametrizations `string` Re-Parametrizations of the weights of layers. Allowed values: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'

ng_conv_normalization_type `string` Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'

ng_conv_residual_type `string` Type of residual connection for all layers and stack of layers. Allowed values: 'ResidualGate', 'Addition', 'None'

dense_act_fct `string` Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU', 'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'

dense_n_layers `int` Number of dense layers. Allowed values: $0 \leq x$

dense_dropout `double` determining the dropout between dense layers. Allowed values: $0 \leq x \leq 0.6$

dense_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.
dense_parametrizations string Re-Parametrizations of the weights of layers. Allowed values: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
dense_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'
dense_residual_type string Type of residual connection for all layers and stack of layers. Allowed values: 'ResidualGate', 'Addition', 'None'
rec_act_fct string Activation function for all layers. Allowed values: 'Tanh'
rec_n_layers int Number of recurrent layers. Allowed values: $0 \leq x$
rec_type string Type of the recurrent layers. `rec_type='GRU'` for Gated Recurrent Unit and `rec_type='LSTM'` for Long Short-Term Memory. Allowed values: 'GRU', 'LSTM'
rec_bidirectional bool If TRUE a bidirectional version of the recurrent layers is used.
rec_dropout double determining the dropout between recurrent layers. Allowed values: $0 \leq x \leq 0.6$
rec_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.
rec_parametrizations string Re-Parametrizations of the weights of layers. Allowed values: 'None'
rec_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'
rec_residual_type string Type of residual connection for all layers and stack of layers. Allowed values: 'ResidualGate', 'Addition', 'None'
tf_act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU', 'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'
tf_dense_dim int determining the size of the projection layer within a each transformer encoder. Allowed values: $1 \leq x$
tf_n_layers int determining how many times the encoder should be added to the network. Allowed values: $0 \leq x$
tf_dropout_rate_1 double determining the dropout after the attention mechanism within the transformer encoder layers. Allowed values: $0 \leq x \leq 0.6$
tf_dropout_rate_2 double determining the dropout for the dense projection within the transformer encoder layers. Allowed values: $0 \leq x \leq 0.6$
tf_attention_type string Choose the attention type. Allowed values: 'Fourier', 'Multi-Head'
tf_positional_type string Type of processing positional information. Allowed values: 'absolute'
tf_num_heads int determining the number of attention heads for a self-attention layer. Only relevant if `attention_type='multihead'` Allowed values: $0 \leq x$
tf_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.
tf_parametrizations string Re-Parametrizations of the weights of layers. Allowed values: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
tf_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'

```

tf_residual_type string Type of residual connection for all layers and stack of layers.
    Allowed values: 'ResidualGate', 'Addition', 'None'

merge_attention_type string Choose the attention type. Allowed values: 'Fourier', 'MultiHead'

merge_num_heads int determining the number of attention heads for a self-attention layer.
    Only relevant if attention_type='multihead' Allowed values:  $0 \leq x$ 

merge_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'

merge_pooling_features int Number of features to be extracted at the end of the model.
    Allowed values:  $1 \leq x$ 

merge_pooling_type string Type of extracting intermediate features. Allowed values: 'Max',
    'Min', 'MinMax'

embedding_dim int determining the number of dimensions for the embedding. Allowed values:  $2 \leq x$ 

```

Returns: Function does nothing return. It modifies the current object.

Method clone(): The objects of this class are cloneable with this method.

Usage:

```
TEClassifierParallelPrototype$clone(deep = FALSE)
```

Arguments:

deep Whether to make a deep clone.

References

- Oreshkin, B. N., Rodriguez, P. & Lacoste, A. (2018). TADAM: Task dependent adaptive metric for improved few-shot learning. <https://doi.org/10.48550/arXiv.1805.10123>
- Snell, J., Swersky, K. & Zemel, R. S. (2017). Prototypical Networks for Few-shot Learning. <https://doi.org/10.48550/arXiv.1703.05175>
- Zhang, X., Nie, J., Zong, L., Yu, H. & Liang, W. (2019). One Shot Learning with Margin. In Q. Yang, Z.-H. Zhou, Z. Gong, M.-L. Zhang & S.-J. Huang (Eds.), Lecture Notes in Computer Science. Advances in Knowledge Discovery and Data Mining (Vol. 11440, pp. 305–317). Springer International Publishing. https://doi.org/10.1007/978-3-030-16145-3_24

See Also

Other Classification: [TEClassifierParallel](#), [TEClassifierProtoNet](#), [TEClassifierRegular](#), [TEClassifierSequential](#), [TEClassifierSequentialPrototype](#)

TEClassifierProtoNet *Text embedding classifier with a ProtoNet*

Description

Abstract class for neural nets with 'pytorch'.

This class is **deprecated**. Please use an Object of class [TEClassifierSequentialPrototype](#) instead.

This object represents in implementation of a prototypical network for few-shot learning as described by Snell, Swersky, and Zemel (2017). The network uses a multi way contrastive loss described by Zhang et al. (2019). The network learns to scale the metric as described by Oreshkin, Rodriguez, and Lacoste (2018)

Value

Objects of this class are used for assigning texts to classes/categories. For the creation and training of a classifier an object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) and a factor are necessary. The object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) contains the numerical text representations (text embeddings) of the raw texts generated by an object of class [TextEmbeddingModel](#). The factor contains the classes/categories for every text. Missing values (unlabeled cases) are supported. For predictions an object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) has to be used which was created with the same [TextEmbeddingModel](#) as for training.

Super classes

```
aifedducation::AIFEBaseModel -> aifedducation::ModelsBasedOnTextEmbeddings -> aifedducation::ClassifiersB
-> aifedducation::TEClassifiersBasedOnProtoNet -> TEClassifierProtoNet
```

Methods

Public methods:

- [TEClassifierProtoNet\\$new\(\)](#)
- [TEClassifierProtoNet\\$configure\(\)](#)
- [TEClassifierProtoNet\\$embed\(\)](#)
- [TEClassifierProtoNet\\$plot_embeddings\(\)](#)
- [TEClassifierProtoNet\\$clone\(\)](#)

Method [new\(\)](#): Creating a new instance of this class.

Usage:

`TEClassifierProtoNet$new()`

Returns: Returns an object of class [TEClassifierProtoNet](#) which is ready for configuration.

Method [configure\(\)](#): Creating a new instance of this class.

Usage:

```
TEClassifierProtoNet$configure(
  name = NULL,
  label = NULL,
  text_embeddings = NULL,
  feature_extractor = NULL,
  target_levels = NULL,
  dense_size = 4,
  dense_layers = 0,
  rec_size = 4,
  rec_layers = 2,
  rec_type = "GRU",
  rec_bidirectional = FALSE,
  embedding_dim = 2,
  self_attention_heads = 0,
  intermediate_size = NULL,
  attention_type = "Fourier",
  add_pos_embedding = TRUE,
  act_fct = "ELU",
  parametrizations = "None",
  rec_dropout = 0.1,
  repeat_encoder = 1,
  dense_dropout = 0.4,
  encoder_dropout = 0.1
)
```

Arguments:

name string Name of the new model. Please refer to common name conventions. Free text can be used with parameter **label**. If set to **NULL** a unique ID is generated automatically.

Allowed values: any

label string Label for the new model. Here you can use free text. Allowed values: any

text_embeddings EmbeddedText, LargeDataSetForTextEmbeddings Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#).

feature_extractor TFeatureExtractor Object of class [TFeatureExtractor](#) which should be used in order to reduce the number of dimensions of the text embeddings. If no feature extractor should be applied set **NULL**.

target_levels vector containing the levels (categories or classes) within the target data. Please note that order matters. For ordinal data please ensure that the levels are sorted correctly with later levels indicating a higher category/class. For nominal data the order does not matter.

dense_size int Number of neurons for each dense layer. Allowed values: $1 \leq x$

dense_layers int Number of dense layers. Allowed values: $0 \leq x$

rec_size int Number of neurons for each recurrent layer. Allowed values: $1 \leq x$

rec_layers int Number of recurrent layers. Allowed values: $0 \leq x$

rec_type string Type of the recurrent layers. **rec_type='GRU'** for Gated Recurrent Unit and **rec_type='LSTM'** for Long Short-Term Memory. Allowed values: 'GRU', 'LSTM'

rec_bidirectional bool If **TRUE** a bidirectional version of the recurrent layers is used.

embedding_dim int determining the number of dimensions for the embedding. Allowed values: $2 \leq x$

```

self_attention_heads int determining the number of attention heads for a self-attention
layer. Only relevant if attention_type='multihead' Allowed values: 0 <= x
intermediate_size int determining the size of the projection layer within a each transformer
encoder. Allowed values: 1 <= x
attention_type string Choose the attention type. Allowed values: 'Fourier', 'MultiHead'
add_pos_embedding bool TRUE if positional embedding should be used.
act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU',
'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'
parametrizations string Re-Parametrizations of the weights of layers. Allowed values:
'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
rec_dropout double determining the dropout between recurrent layers. Allowed values: 0 <= x <= 0.6
repeat_encoder int determining how many times the encoder should be added to the net-
work. Allowed values: 0 <= x
dense_dropout double determining the dropout between dense layers. Allowed values: 0 <= x <= 0.6
encoder_dropout double determining the dropout for the dense projection within the trans-
former encoder layers. Allowed values: 0 <= x <= 0.6
bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the
layers.

```

Method embed(): Method for embedding documents. Please do not confuse this type of em-
beddings with the embeddings of texts created by an object of class [TextEmbeddingModel](#). These
embed documents according to their similarity to specific classes.

Usage:

```
TEClassifierProtoNet$embed(embeddings_q = NULL, batch_size = 32)
```

Arguments:

embeddings_q Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) containing
the text embeddings for all cases which should be embedded into the classification space.

batch_size int batch size.

Returns: Returns a list containing the following elements

- embeddings_q: embeddings for the cases (query sample).
- embeddings_prototypes: embeddings of the prototypes which were learned during train-
ing. They represents the center for the different classes.

Method plot_embeddings(): Method for creating a plot to visualize embeddings and their
corresponding centers (prototypes).

Usage:

```
TEClassifierProtoNet$plot_embeddings(
  embeddings_q,
  classes_q = NULL,
  batch_size = 12,
  alpha = 0.5,
  size_points = 3,
  size_points_prototypes = 8,
  inc_unlabeled = TRUE
)
```

Arguments:

`embeddings_q` Object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` containing the text embeddings for all cases which should be embedded into the classification space.

`classes_q` Named factor containg the true classes for every case. Please note that the names must match the names/ids in `embeddings_q`.

`batch_size` int batch size.

`alpha` float Value indicating how transparent the points should be (important if many points overlap). Does not apply to points representing prototypes.

`size_points` int Size of the points excluding the points for prototypes.

`size_points_prototypes` int Size of points representing prototypes.

`inc_unlabeled` bool If TRUE plot includes unlabeled cases as data points.

Returns: Returns a plot of class ggplotvisualizing embeddings.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
TEClassifierProtoNet$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

Note

This model requires `pad_value=0`. If this condition is not met the padding value is switched automatically.

References

Oreshkin, B. N., Rodriguez, P. & Lacoste, A. (2018). TADAM: Task dependent adaptive metric for improved few-shot learning. <https://doi.org/10.48550/arXiv.1805.10123>

Snell, J., Swersky, K. & Zemel, R. S. (2017). Prototypical Networks for Few-shot Learning. <https://doi.org/10.48550/arXiv.1703.05175>

Zhang, X., Nie, J., Zong, L., Yu, H. & Liang, W. (2019). One Shot Learning with Margin. In Q. Yang, Z.-H. Zhou, Z. Gong, M.-L. Zhang & S.-J. Huang (Eds.), Lecture Notes in Computer Science. Advances in Knowledge Discovery and Data Mining (Vol. 11440, pp. 305–317). Springer International Publishing. https://doi.org/10.1007/978-3-030-16145-3_24

See Also

Other Classification: `TEClassifierParallel`, `TEClassifierParallelPrototype`, `TEClassifierRegular`, `TEClassifierSequential`, `TEClassifierSequentialPrototype`

`TEClassifierRegular` *Text embedding classifier with a neural net*

Description

Abstract class for neural nets with 'pytorch'.

This class is **deprecated**. Please use an Object of class `TEClassifierSequential` instead.

Value

Objects of this class are used for assigning texts to classes/categories. For the creation and training of a classifier an object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` on the one hand and a `factor` on the other hand are necessary.

The object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` contains the numerical text representations (text embeddings) of the raw texts generated by an object of class `TextEmbeddingModel`. For supporting large data sets it is recommended to use `LargeDataSetForTextEmbeddings` instead of `EmbeddedText`.

The factor contains the classes/categories for every text. Missing values (unlabeled cases) are supported and can be used for pseudo labeling.

For predictions an object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` has to be used which was created with the same `TextEmbeddingModel` as for training.

Super classes

```
aifedducation::AIFEBaseModel -> aifedducation::ModelsBasedOnTextEmbeddings -> aifedducation::ClassifiersBasedOnTextEmbeddings -> aifedducation::TEClassifiersBasedOnRegular -> TEClassifierRegular
```

Methods

Public methods:

- `TEClassifierRegular$new()`
- `TEClassifierRegular$configure()`
- `TEClassifierRegular$clone()`

Method `new()`: Creating a new instance of this class.

Usage:

```
TEClassifierRegular$new()
```

Returns: Returns an object of class `TEClassifierRegular` which is ready for configuration.

Method `configure()`: Creating a new instance of this class.

Usage:

```
TEClassifierRegular$configure(
  name = NULL,
  label = NULL,
  text_embeddings = NULL,
  feature_extractor = NULL,
  target_levels = NULL,
  bias = TRUE,
  dense_size = 4,
  dense_layers = 0,
  rec_size = 4,
  rec_layers = 2,
  rec_type = "GRU",
  rec_bidirectional = FALSE,
  self_attention_heads = 0,
  intermediate_size = NULL,
  attention_type = "Fourier",
  add_pos_embedding = TRUE,
  act_fct = "ELU",
  parametrizations = "None",
  rec_dropout = 0.1,
  repeat_encoder = 1,
  dense_dropout = 0.4,
  encoder_dropout = 0.1
)
```

Arguments:

name string Name of the new model. Please refer to common name conventions. Free text can be used with parameter **label**. If set to NULL a unique ID is generated automatically.
Allowed values: any

label string Label for the new model. Here you can use free text. Allowed values: any

text_embeddings EmbeddedText, LargeDataSetForTextEmbeddings Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#).

feature_extractor TFEFeatureExtractor Object of class [TEFeatureExtractor](#) which should be used in order to reduce the number of dimensions of the text embeddings. If no feature extractor should be applied set NULL.

target_levels vector containing the levels (categories or classes) within the target data.
Please note that order matters. For ordinal data please ensure that the levels are sorted correctly with later levels indicating a higher category/class. For nominal data the order does not matter.

bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.

dense_size int Number of neurons for each dense layer. Allowed values: $1 \leq x$

dense_layers int Number of dense layers. Allowed values: $0 \leq x$

rec_size int Number of neurons for each recurrent layer. Allowed values: $1 \leq x$

rec_layers int Number of recurrent layers. Allowed values: $0 \leq x$

rec_type string Type of the recurrent layers. **rec_type='GRU'** for Gated Recurrent Unit and **rec_type='LSTM'** for Long Short-Term Memory. Allowed values: 'GRU', 'LSTM'

```

rec_bidirectional bool If TRUE a bidirectional version of the recurrent layers is used.
self_attention_heads int determining the number of attention heads for a self-attention
layer. Only relevant if attention_type='multihead' Allowed values: 0 <= x
intermediate_size int determining the size of the projection layer within a each transformer
encoder. Allowed values: 1 <= x
attention_type string Choose the attention type. Allowed values: 'Fourier', 'MultiHead'
add_pos_embedding bool TRUE if positional embedding should be used.
act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU',
'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'
parametrizations string Re-Parametrizations of the weights of layers. Allowed values:
'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
rec_dropout double determining the dropout between recurrent layers. Allowed values: 0 <= x <= 0.6
repeat_encoder int determining how many times the encoder should be added to the net-
work. Allowed values: 0 <= x
dense_dropout double determining the dropout between dense layers. Allowed values: 0 <= x <= 0.6
encoder_dropout double determining the dropout for the dense projection within the trans-
former encoder layers. Allowed values: 0 <= x <= 0.6

```

Returns: Returns an object of class [TEClassifierRegular](#) which is ready for training.

Method clone(): The objects of this class are cloneable with this method.

Usage:

```
TEClassifierRegular$clone(deep = FALSE)
```

Arguments:

deep Whether to make a deep clone.

Note

This model requires pad_value=0. If this condition is not met the padding value is switched auto-
matically.

See Also

Other Classification: [TEClassifierParallel](#), [TEClassifierParallelPrototype](#), [TEClassifierProtoNet](#),
[TEClassifierSequential](#), [TEClassifierSequentialPrototype](#)

Description

Base class for classifiers relying on [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) as input which use the architecture of Protonets and its corresponding training techniques.

Objects of this class containing fields and methods used in several other classes in 'AI for Education'.

This class is **not** designed for a direct application and should only be used by developers.

Value

A new object of this class.

Super classes

```
aifedducation::AIFEBaseModel -> aifedducation::ModelsBasedOnTextEmbeddings -> aifedducation::ClassifiersBasedOnProtoNet -> TEClassifiersBasedOnProtoNet
```

Methods

Public methods:

- [TEClassifiersBasedOnProtoNet\\$train\(\)](#)
- [TEClassifiersBasedOnProtoNet\\$predict_with_samples\(\)](#)
- [TEClassifiersBasedOnProtoNet\\$embed\(\)](#)
- [TEClassifiersBasedOnProtoNet\\$get_metric_scale_factor\(\)](#)
- [TEClassifiersBasedOnProtoNet\\$plot_embeddings\(\)](#)
- [TEClassifiersBasedOnProtoNet\\$clone\(\)](#)

Method train(): Method for training a neural net.

Training includes a routine for early stopping. In the case that loss<0.0001 and Accuracy=1.00 and Average Iota=1.00 training stops. The history uses the values of the last trained epoch for the remaining epochs.

After training the model with the best values for Average Iota, Accuracy, and Loss on the validation data set is used as the final model.

Usage:

```
TEClassifiersBasedOnProtoNet$train(
  data_embeddings = NULL,
  data_targets = NULL,
  data_folds = 5,
  data_val_size = 0.25,
  loss_pt_fct_name = "MultiWayContrastiveLoss",
  use_sc = FALSE,
  sc_method = "knnor",
  sc_min_k = 1,
  sc_max_k = 10,
  use_pl = FALSE,
  pl_max_steps = 3,
  pl_max = 1,
```

```

    pl_anchor = 1,
    pl_min = 0,
    sustain_track = TRUE,
    sustain_iso_code = NULL,
    sustain_region = NULL,
    sustain_interval = 15,
    epochs = 40,
    batch_size = 35,
    Ns = 5,
    Nq = 3,
    loss_alpha = 0.5,
    loss_margin = 0.05,
    sampling_separate = FALSE,
    sampling_shuffle = TRUE,
    trace = TRUE,
    ml_trace = 1,
    log_dir = NULL,
    log_write_interval = 10,
    n_cores = auto_n_cores(),
    lr_rate = 0.001,
    lr_warm_up_ratio = 0.02,
    optimizer = "AdamW"
)

```

Arguments:

`data_embeddings` `EmbeddedText`, `LargeDataSetForTextEmbeddings` Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#).

`data_targets` factor containing the labels for cases stored in embeddings. Factor must be named and has to use the same names as used in the embeddings. .

`data_folds` int determining the number of cross-fold samples. Allowed values: $1 \leq x$

`data_val_size` double between 0 and 1, indicating the proportion of cases which should be used for the validation sample during the estimation of the model. The remaining cases are part of the training data. Allowed values: $0 < x < 1$

`loss_pt_fct_name` string Name of the loss function to use during training. Allowed values: '`MultiWayContrastiveLoss`'

`use_sc` bool TRUE if the estimation should integrate synthetic cases. FALSE if not.

`sc_method` string containing the method for generating synthetic cases. Allowed values: '`knor`'

`sc_min_k` int determining the minimal number of k which is used for creating synthetic units. Allowed values: $1 \leq x$

`sc_max_k` int determining the maximal number of k which is used for creating synthetic units. Allowed values: $1 \leq x$

`use_pl` bool TRUE if the estimation should integrate pseudo-labeling. FALSE if not.

`pl_max_steps` int determining the maximum number of steps during pseudo-labeling. Allowed values: $1 \leq x$

`pl_max` double setting the maximal level of confidence for considering a case for pseudo-labeling. Allowed values: $0 < x \leq 1$

pl_anchor double indicating the reference point for sorting the new cases of every label. Allowed values: $0 \leq x \leq 1$
 pl_min double setting the minimal level of confidence for considering a case for pseudo-labeling. Allowed values: $0 \leq x < 1$
 sustain_track bool If TRUE energy consumption is tracked during training via the python library 'codecarbon'.
 sustain_iso_code string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes. Allowed values: any
 sustain_region string Region within a country. Only available for USA and Canada See the documentation of codecarbon for more information. [https://mlco2.github.io/codecarbon/parameters.html](https://mlco2.github.io/codecarbon(parameters.html) Allowed values: any
 sustain_interval int Interval in seconds for measuring power usage. Allowed values: $1 \leq x$
 epochs int Number of training epochs. Allowed values: $1 \leq x$
 batch_size int Size of the batches for training. Allowed values: $1 \leq x$
 Ns int Number of cases for every class in the sample. Allowed values: $1 \leq x$
 Nq int Number of cases for every class in the query. Allowed values: $1 \leq x$
 loss_alpha double Value between 0 and 1 indicating how strong the loss should focus on pulling cases to its corresponding prototypes or pushing cases away from other prototypes. The higher the value the more the loss concentrates on pulling cases to its corresponding prototypes. Allowed values: $0 \leq x \leq 1$
 loss_margin double Value greater 0 indicating the minimal distance of every case from prototypes of other classes. Please note that in contrast to the original work by Zhang et al. (2019) this implementation reaches better performance if the margin is a magnitude lower (e.g. 0.05 instead of 0.5). Allowed values: $0 \leq x \leq 1$
 sampling_separate bool If TRUE the cases for every class are divided into a data set for sample and for query. These are never mixed. If TRUE sample and query cases are drawn from the same data pool. That is, a case can be part of sample in one epoch and in another epoch it can be part of query. It is ensured that a case is never part of sample and query at the same time. In addition, it is ensured that every cases exists only once during a training step.
 sampling_shuffle bool if TRUE cases are randomly drawn from the data during every step. If FALSE the cases are not shuffled.
 trace bool TRUE if information about the estimation phase should be printed to the console.
 ml_trace int ml_trace=0 does not print any information about the training process from pytorch on the console. Allowed values: $0 \leq x \leq 1$
 log_dir string Path to the directory where the log files should be saved. If no logging is desired set this argument to NULL. Allowed values: any
 log_write_interval int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if log_dir is not NULL. Allowed values: $1 \leq x$
 n_cores int Number of cores which should be used during the calculation of synthetic cases. Only relevant if use_sc=TRUE. Allowed values: $1 \leq x$
 lr_rate double Initial learning rate for the training. Allowed values: $0 < x \leq 1$
 lr_warm_up_ratio double Number of epochs used for warm up. Allowed values: $0 < x < 0.5$
 optimizer string determining the optimizer used for training. Allowed values: 'Adam', 'RMSprop', 'AdamW', 'SGD'

`loss_balance_class_weights` bool If TRUE class weights are generated based on the frequencies of the training data with the method Inverse Class Frequency. If FALSE each class has the weight 1.

`loss_balance_sequence_length` bool If TRUE sample weights are generated for the length of sequences based on the frequencies of the training data with the method Inverse Class Frequency. If FALSE each sequences length has the weight 1.

Details:

- `sc_max_k`: All values from `sc_min_k` up to `sc_max_k` are successively used. If the number of `sc_max_k` is too high, the value is reduced to a number that allows the calculating of synthetic units.
- `pl_anchor`: With the help of this value, the new cases are sorted. For this aim, the distance from the anchor is calculated and all cases are arranged into an ascending order.

Returns: Function does not return a value. It changes the object into a trained classifier.

Method `predict_with_samples()`: Method for predicting the class of given data (query) based on provided examples (sample).

Usage:

```
TEClassifiersBasedOnProtoNet$predict_with_samples(
  newdata,
  batch_size = 32,
  ml_trace = 1,
  embeddings_s = NULL,
  classes_s = NULL
)
```

Arguments:

`newdata` Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) containing the text embeddings for all cases which should be predicted. They form the query set.

`batch_size` int batch size.

`ml_trace` int `ml_trace=0` does not print any information about the training process from pytorch on the console. Allowed values: $0 \leq x \leq 1$

`embeddings_s` Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) containing the text embeddings for all reference examples. They form the sample set.

`classes_s` Named factor containing the classes for every case within `embeddings_s`.

Returns: Returns a `data.frame` containing the predictions and the probabilities of the different labels for each case.

Method `embed()`: Method for embedding documents. Please do not confuse this type of embeddings with the embeddings of texts created by an object of class [TextEmbeddingModel](#). These embeddings embed documents according to their similarity to specific classes.

Usage:

```
TEClassifiersBasedOnProtoNet$embed(
  embeddings_q = NULL,
  embeddings_s = NULL,
  classes_s = NULL,
  batch_size = 32,
  ml_trace = 1
)
```

Arguments:

`embeddings_q` Object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` containing the text embeddings for all cases which should be embedded into the classification space.

`embeddings_s` Object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` containing the text embeddings for all reference examples. They form the sample set. If set to NULL the trained prototypes are used.

`classes_s` Named factor containing the classes for every case within `embeddings_s`. If set to NULL the trained prototypes are used.

`batch_size` int batch size.

`ml_trace` int `ml_trace=0` does not print any information about the training process from pytorch on the console. Allowed values: $0 \leq x \leq 1$

Returns: Returns a list containing the following elements

- `embeddings_q`: embeddings for the cases (query sample).
- `distances_q`: matrix containing the distance of every query case to every prototype.
- `embeddings_prototypes`: embeddings of the prototypes which were learned during training. They represents the center for the different classes.

Method `get_metric_scale_factor()`: Method returns the scaling factor of the metric.

Usage:

```
TEClassifiersBasedOnProtoNet$get_metric_scale_factor()
```

Returns: Returns the scaling factor of the metric as float.

Method `plot_embeddings()`: Method for creating a plot to visualize embeddings and their corresponding centers (prototypes).

Usage:

```
TEClassifiersBasedOnProtoNet$plot_embeddings(
  embeddings_q,
  classes_q = NULL,
  embeddings_s = NULL,
  classes_s = NULL,
  batch_size = 12,
  alpha = 0.5,
  size_points = 3,
  size_points_prototypes = 8,
  inc_unlabeled = TRUE,
  inc_margin = TRUE
)
```

Arguments:

`embeddings_q` Object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` containing the text embeddings for all cases which should be embedded into the classification space.

`classes_q` Named factor containg the true classes for every case. Please note that the names must match the names/ids in `embeddings_q`.

`embeddings_s` Object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` containing the text embeddings for all reference examples. They form the sample set. If set to NULL the trained prototypes are used.

`classes_s` Named factor containing the classes for every case within `embeddings_s`. If set to NULL the trained prototypes are used.

`batch_size` int batch size.

`alpha` float Value indicating how transparent the points should be (important if many points overlap). Does not apply to points representing prototypes.

`size_points` int Size of the points excluding the points for prototypes.

`size_points_prototypes` int Size of points representing prototypes.

`inc_unlabeled` bool If TRUE plot includes unlabeled cases as data points.

`inc_margin` bool If TRUE plot includes the margin around every prototype. Adding margin requires a trained model. If the model is not trained this argument is treated as set to FALSE.

Returns: Returns a plot of class ggplotvisualizing embeddings.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

`TEClassifiersBasedOnProtoNet$clone(deep = FALSE)`

Arguments:

`deep` Whether to make a deep clone.

See Also

Other R6 Classes for Developers: [AIFEBaseModel](#), [ClassifiersBasedOnTextEmbeddings](#), [LargeDataSetBase](#), [ModelsBasedOnTextEmbeddings](#), [TEClassifiersBasedOnRegular](#)

TEClassifiersBasedOnRegular

Base class for regular classifiers relying on [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) as input

Description

Abstract class for all regular classifiers that use numerical representations of texts instead of words.

Objects of this class containing fields and methods used in several other classes in 'AI for Education'.

This class is **not** designed for a direct application and should only be used by developers.

Value

A new object of this class.

Super classes

```
aifedducation::AIFEBaseModel -> aifedducation::ModelsBasedOnTextEmbeddings -> aifedducation::ClassifiersB
-> TEClassifiersBasedOnRegular
```

Methods

Public methods:

- `TEClassifiersBasedOnRegular$train()`
- `TEClassifiersBasedOnRegular$clone()`

Method train(): Method for training a neural net.

Training includes a routine for early stopping. In the case that loss<0.0001 and Accuracy=1.00 and Average Iota=1.00 training stops. The history uses the values of the last trained epoch for the remaining epochs.

After training the model with the best values for Average Iota, Accuracy, and Loss on the validation data set is used as the final model.

Usage:

```
TEClassifiersBasedOnRegular$train(
  data_embeddings = NULL,
  data_targets = NULL,
  data_folds = 5,
  data_val_size = 0.25,
  loss_balance_class_weights = TRUE,
  loss_balance_sequence_length = TRUE,
  loss_cls_fct_name = "FocalLoss",
  use_sc = FALSE,
  sc_method = "knnor",
  sc_min_k = 1,
  sc_max_k = 10,
  use_pl = FALSE,
  pl_max_steps = 3,
  pl_max = 1,
  pl_anchor = 1,
  pl_min = 0,
  sustain_track = TRUE,
  sustain_iso_code = NULL,
  sustain_region = NULL,
  sustain_interval = 15,
  epochs = 40,
  batch_size = 32,
  trace = TRUE,
  ml_trace = 1,
  log_dir = NULL,
  log_write_interval = 10,
  n_cores = auto_n_cores(),
  lr_rate = 0.001,
  lr_warm_up_ratio = 0.02,
  optimizer = "AdamW"
)
```

Arguments:

`data_embeddings` `EmbeddedText`, `LargeDataSetForTextEmbeddings` Object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings`.

`data_targets` factor containing the labels for cases stored in embeddings. Factor must be named and has to use the same names as used in the embeddings. .

`data_folds` int determining the number of cross-fold samples. Allowed values: $1 \leq x$

`data_val_size` double between 0 and 1, indicating the proportion of cases which should be used for the validation sample during the estimation of the model. The remaining cases are part of the training data. Allowed values: $0 < x < 1$

`loss_balance_class_weights` bool If TRUE class weights are generated based on the frequencies of the training data with the method Inverse Class Frequency. If FALSE each class has the weight 1.

`loss_balance_sequence_length` bool If TRUE sample weights are generated for the length of sequences based on the frequencies of the training data with the method Inverse Class Frequency. If FALSE each sequences length has the weight 1.

`loss_cls_fct_name` string Name of the loss function to use during training. Allowed values: 'FocalLoss', 'CrossEntropyLoss'

`use_sc` bool TRUE if the estimation should integrate synthetic cases. FALSE if not.

`sc_method` string containing the method for generating synthetic cases. Allowed values: 'knor'

`sc_min_k` int determining the minimal number of k which is used for creating synthetic units. Allowed values: $1 \leq x$

`sc_max_k` int determining the maximal number of k which is used for creating synthetic units. Allowed values: $1 \leq x$

`use_pl` bool TRUE if the estimation should integrate pseudo-labeling. FALSE if not.

`pl_max_steps` int determining the maximum number of steps during pseudo-labeling. Allowed values: $1 \leq x$

`pl_max` double setting the maximal level of confidence for considering a case for pseudo-labeling. Allowed values: $0 < x \leq 1$

`pl_anchor` double indicating the reference point for sorting the new cases of every label. Allowed values: $0 \leq x \leq 1$

`pl_min` double setting the minimal level of confidence for considering a case for pseudo-labeling. Allowed values: $0 \leq x < 1$

`sustain_track` bool If TRUE energy consumption is tracked during training via the python library 'codecarbon'.

`sustain_iso_code` string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes. Allowed values: any

`sustain_region` string Region within a country. Only available for USA and Canada See the documentation of codecarbon for more information. [https://mlco2.github.io/codecarbon\(parameters.html\)](https://mlco2.github.io/codecarbon(parameters.html)) Allowed values: any

`sustain_interval` int Interval in seconds for measuring power usage. Allowed values: $1 \leq x$

`epochs` int Number of training epochs. Allowed values: $1 \leq x$

`batch_size` int Size of the batches for training. Allowed values: $1 \leq x$

`trace` bool TRUE if information about the estimation phase should be printed to the console.

`ml_trace` int `ml_trace=0` does not print any information about the training process from pytorch on the console. Allowed values: $0 \leq x \leq 1$

```

log_dir string Path to the directory where the log files should be saved. If no logging is
desired set this argument to NULL. Allowed values: any
log_write_interval int Time in seconds determining the interval in which the logger should
try to update the log files. Only relevant if log_dir is not NULL. Allowed values: 1 <= x
n_cores int Number of cores which should be used during the calculation of synthetic cases.
Only relevant if use_sc=TRUE. Allowed values: 1 <= x
lr_rate double Initial learning rate for the training. Allowed values: 0 < x <= 1
lr_warm_up_ratio double Number of epochs used for warm up. Allowed values: 0 < x < 0.5
optimizer string determining the optimizer used for training. Allowed values: 'Adam',
'RMSprop', 'AdamW', 'SGD'
```

Details:

- sc_max_k: All values from sc_min_k up to sc_max_k are successively used. If the number of sc_max_k is too high, the value is reduced to a number that allows the calculating of synthetic units.
- pl_anchor: With the help of this value, the new cases are sorted. For this aim, the distance from the anchor is calculated and all cases are arranged into an ascending order.

Returns: Function does not return a value. It changes the object into a trained classifier.

Method clone(): The objects of this class are cloneable with this method.

Usage:

```
TEClassifiersBasedOnRegular$clone(deep = FALSE)
```

Arguments:

deep Whether to make a deep clone.

See Also

Other R6 Classes for Developers: [AIFEBaseModel](#), [ClassifiersBasedOnTextEmbeddings](#), [LargeDataSetBase](#), [ModelsBasedOnTextEmbeddings](#), [TEClassifiersBasedOnProtoNet](#)

TEClassifierSequential

Text embedding classifier with a neural net

Description

Classification Type

This is a probability classifier that predicts a probability distribution for different classes/categories. This is the standard case most common in literature.

Sequential Core Architecture

This model is based on a sequential architecture. The input is passed to a specific number of layers step by step. All layers are grouped by their kind into stacks.

Transformer Encoder Layers

Description

The transformer encoder layers follow the structure of the encoder layers used in transformer models. A single layer is designed as described by Chollet, Kalinowski, and Allaire (2022, p. 373) with the exception that single components of the layers (such as the activation function, the kind of residual connection, the kind of normalization or the kind of attention) can be customized. All parameters with the prefix *tf_* can be used to configure this layer.

Feature Layer

Description

The feature layer is a dense layer that can be used to increase or decrease the number of features of the input data before passing the data into your model. The aim of this layer is to increase or reduce the complexity of the data for your model. The output size of this layer determines the number of features for all following layers. In the special case that the requested number of features equals the number of features of the text embeddings this layer is reduced to a dropout layer with masking capabilities. All parameters with the prefix *feat_* can be used to configure this layer.

Dense Layers

Description

A fully connected layer. The layer is applied to every step of a sequence. All parameters with the prefix *dense_* can be used to configure this layer.

Multiple N-Gram Layers

Description

This type of layer focuses on sub-sequence and performs an 1d convolutional operation. On a word and token level these sub-sequences can be interpreted as n-grams (Jacovi, Shalom & Goldberg 2018). The convolution is done across all features. The number of filters equals the number of features of the input tensor. Thus, the shape of the tensor is retained (Pham, Kruszewski & Boleda 2016).

The layer is able to consider multiple n-grams at the same time. In this case the convolution of the n-grams is done separately and the resulting tensors are concatenated along the feature dimension. The number of filters for every n-gram is set to *num_features/num_n-grams*. Thus, the resulting tensor has the same shape as the input tensor.

Sub-sequences that are masked in the input are also masked in the output.

The output of this layer can be understand as the results of the n-gram filters. Stacking this layer allows the model to perform n-gram detection of n-grams (meta perspective). All parameters with the prefix *ng_conv_* can be used to configure this layer.

Recurrent Layers

Description

A regular recurrent layer either as Gated Recurrent Unit (GRU) or Long Short-Term Memory (LSTM) layer. Uses PyTorchs implementation. All parameters with the prefix *rec_* can be used to configure this layer.

Classification Pooling Layer

Description

Layer transforms sequences into a lower dimensional space that can be passed to dense layers. It performs two types of pooling. First, it extractes features across the time dimension selecting the maximal and/or minimal features. Second, it performs pooling over the remaining features selecting a specific number of the highest and/or lowest features.

In the case of selecting the minimal *and* maximal features at the same time the minimal features are concatenated to the tensor of the maximal features resulting in the shape \$(Batch, Times, 2*Features)\$ at the end of the first step. In the second step the number of requested features is halved. The first half is used for the maximal features and the second for the minimal features. All parameters with the prefix `cls_pooling_` can be used to configure this layer.

Training and Prediction

For the creation and training of a classifier an object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` on the one hand and a `factor` on the other hand are necessary.

The object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` contains the numerical text representations (text embeddings) of the raw texts generated by an object of class `TextEmbeddingModel`. For supporting large data sets it is recommended to use `LargeDataSetForTextEmbeddings` instead of `EmbeddedText`.

The factor contains the classes/categories for every text. Missing values (unlabeled cases) are supported and can be used for pseudo labeling.

For predictions an object of class `EmbeddedText` or `LargeDataSetForTextEmbeddings` has to be used which was created with the same `TextEmbeddingModel` as for training.

Value

Returns a new object of this class ready for configuration or for loading a saved classifier.

Super classes

```
aifedducation::AIFE BaseModel -> aifedducation::ModelsBasedOnTextEmbeddings -> aifedducation::ClassifiersB
-> aifedducation::TEClassifiersBasedOnRegular -> TEClassifierSequential
```

Methods

Public methods:

- `TEClassifierSequential$configure()`
- `TEClassifierSequential$clone()`

Method `configure()`: Creating a new instance of this class.

Usage:

```
TEClassifierSequential$configure(
  name = NULL,
  label = NULL,
  text_embeddings = NULL,
  feature_extractor = NULL,
  target_levels = NULL,
  skip_connection_type = "ResidualGate",
  cls_pooling_features = NULL,
  cls_pooling_type = "MinMax",
  feat_act_fct = "ELU",
  feat_size = 50,
  feat_bias = TRUE,
  feat_dropout = 0,
```

```

feat_parametrizations = "None",
feat_normalization_type = "LayerNorm",
ng_conv_act_fct = "ELU",
ng_conv_n_layers = 1,
ng_conv_ks_min = 2,
ng_conv_ks_max = 4,
ng_conv_bias = FALSE,
ng_conv_dropout = 0.1,
ng_conv_parametrizations = "None",
ng_conv_normalization_type = "LayerNorm",
ng_conv_residual_type = "ResidualGate",
dense_act_fct = "ELU",
dense_n_layers = 1,
dense_dropout = 0.5,
dense_bias = FALSE,
dense_parametrizations = "None",
dense_normalization_type = "LayerNorm",
dense_residual_type = "ResidualGate",
rec_act_fct = "Tanh",
rec_n_layers = 1,
rec_type = "GRU",
rec_bidirectional = FALSE,
rec_dropout = 0.2,
rec_bias = FALSE,
rec_parametrizations = "None",
rec_normalization_type = "LayerNorm",
rec_residual_type = "ResidualGate",
tf_act_fct = "ELU",
tf_dense_dim = 50,
tf_n_layers = 1,
tf_dropout_rate_1 = 0.1,
tf_dropout_rate_2 = 0.5,
tf_attention_type = "MultiHead",
tf_positional_type = "absolute",
tf_num_heads = 1,
tf_bias = FALSE,
tf_parametrizations = "None",
tf_normalization_type = "LayerNorm",
tf_residual_type = "ResidualGate"
)

```

Arguments:

name string Name of the new model. Please refer to common name conventions. Free text can be used with parameter **label**. If set to NULL a unique ID is generated automatically.
 Allowed values: any

label string Label for the new model. Here you can use free text. Allowed values: any
text_embeddings EmbeddedText, LargeDataSetForTextEmbeddings Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#).

feature_extractor TFeatureExtractor Object of class [TFeatureExtractor](#) which should

be used in order to reduce the number of dimensions of the text embeddings. If no feature extractor should be applied set NULL.

target_levels vector containing the levels (categories or classes) within the target data. Please note that order matters. For ordinal data please ensure that the levels are sorted correctly with later levels indicating a higher category/class. For nominal data the order does not matter.

skip_connection_type string Type of residual connection for all layers and stack of layers. Allowed values: 'ResidualGate', 'Addition', 'None'

cls_pooling_features int Number of features to be extracted at the end of the model. Allowed values: $1 \leq x$

cls_pooling_type string Type of extracting intermediate features. Allowed values: 'Max', 'Min', 'MinMax'

feat_act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU', 'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'

feat_size int Number of neurons for each dense layer. Allowed values: $2 \leq x$

feat_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.

feat_dropout double determining the dropout for the dense projection of the feature layer. Allowed values: $0 \leq x \leq 0.6$

feat_parametrizations string Re-Parametrizations of the weights of layers. Allowed values: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'

feat_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'

ng_conv_act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU', 'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'

ng_conv_n_layers int determining how many times the n-gram layers should be added to the network. Allowed values: $0 \leq x$

ng_conv_ks_min int determining the minimal window size for n-grams. Allowed values: $2 \leq x$

ng_conv_ks_max int determining the maximal window size for n-grams. Allowed values: $2 \leq x$

ng_conv_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.

ng_conv_dropout double determining the dropout for n-gram convolution layers. Allowed values: $0 \leq x \leq 0.6$

ng_conv_parametrizations string Re-Parametrizations of the weights of layers. Allowed values: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'

ng_conv_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'

ng_conv_residual_type string Type of residual connection for all layers and stack of layers. Allowed values: 'ResidualGate', 'Addition', 'None'

dense_act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU', 'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'

dense_n_layers int Number of dense layers. Allowed values: $0 \leq x$

dense_dropout double determining the dropout between dense layers. Allowed values: $0 \leq x \leq 0.6$

dense_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.
dense_parametrizations string Re-Parametrizations of the weights of layers. Allowed values: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
dense_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'
dense_residual_type string Type of residual connection for all layers and stack of layers. Allowed values: 'ResidualGate', 'Addition', 'None'
rec_act_fct string Activation function for all layers. Allowed values: 'Tanh'
rec_n_layers int Number of recurrent layers. Allowed values: $0 \leq x$
rec_type string Type of the recurrent layers. `rec_type='GRU'` for Gated Recurrent Unit and `rec_type='LSTM'` for Long Short-Term Memory. Allowed values: 'GRU', 'LSTM'
rec_bidirectional bool If TRUE a bidirectional version of the recurrent layers is used.
rec_dropout double determining the dropout between recurrent layers. Allowed values: $0 \leq x \leq 0.6$
rec_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.
rec_parametrizations string Re-Parametrizations of the weights of layers. Allowed values: 'None'
rec_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'
rec_residual_type string Type of residual connection for all layers and stack of layers. Allowed values: 'ResidualGate', 'Addition', 'None'
tf_act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU', 'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'
tf_dense_dim int determining the size of the projection layer within a each transformer encoder. Allowed values: $1 \leq x$
tf_n_layers int determining how many times the encoder should be added to the network. Allowed values: $0 \leq x$
tf_dropout_rate_1 double determining the dropout after the attention mechanism within the transformer encoder layers. Allowed values: $0 \leq x \leq 0.6$
tf_dropout_rate_2 double determining the dropout for the dense projection within the transformer encoder layers. Allowed values: $0 \leq x \leq 0.6$
tf_attention_type string Choose the attention type. Allowed values: 'Fourier', 'Multi-Head'
tf_positional_type string Type of processing positional information. Allowed values: 'absolute'
tf_num_heads int determining the number of attention heads for a self-attention layer. Only relevant if `attention_type='multihead'` Allowed values: $0 \leq x$
tf_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.
tf_parametrizations string Re-Parametrizations of the weights of layers. Allowed values: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
tf_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'

`tf_residual_type` string Type of residual connection for all layers and stack of layers.
 Allowed values: 'ResidualGate', 'Addition', 'None'

Returns: Function does nothing return. It modifies the current object.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

`TEClassifierSequential$clone(deep = FALSE)`

Arguments:

`deep` Whether to make a deep clone.

See Also

Other Classification: [TEClassifierParallel](#), [TEClassifierParallelPrototype](#), [TEClassifierProtoNet](#),
[TEClassifierRegular](#), [TEClassifierSequentialPrototype](#)

TEClassifierSequentialPrototype

Text embedding classifier with a ProtoNet

Description

Classification Type

This object is a metric based classifier and represents in implementation of a prototypical network for few-shot learning as described by Snell, Swersky, and Zemel (2017). The network uses a multi way contrastive loss described by Zhang et al. (2019). The network learns to scale the metric as described by Oreshkin, Rodriguez, and Lacoste (2018).

Sequential Core Architecture

This model is based on a sequential architecture. The input is passed to a specific number of layers step by step. All layers are grouped by their kind into stacks.

Transformer Encoder Layers

Description

The transformer encoder layers follow the structure of the encoder layers used in transformer models. A single layer is designed as described by Chollet, Kalinowski, and Allaire (2022, p. 373) with the exception that single components of the layers (such as the activation function, the kind of residual connection, the kind of normalization or the kind of attention) can be customized. All parameters with the prefix `tf_` can be used to configure this layer.

Feature Layer

Description

The feature layer is a dense layer that can be used to increase or decrease the number of features of the input data before passing the data into your model. The aim of this layer is to increase or reduce the complexity of the data for your model. The output size of this layer determines the number of features for all following layers. In the special case that the requested number of features equals

the number of features of the text embeddings this layer is reduced to a dropout layer with masking capabilities. All parameters with the prefix *feat_* can be used to configure this layer.

Dense Layers

Description

A fully connected layer. The layer is applied to every step of a sequence. All parameters with the prefix *dense_* can be used to configure this layer.

Multiple N-Gram Layers

Description

This type of layer focuses on sub-sequence and performs an 1d convolutional operation. On a word and token level these sub-sequences can be interpreted as n-grams (Jacovi, Shalom & Goldberg 2018). The convolution is done across all features. The number of filters equals the number of features of the input tensor. Thus, the shape of the tensor is retained (Pham, Kruszewski & Boleda 2016).

The layer is able to consider multiple n-grams at the same time. In this case the convolution of the n-grams is done separately and the resulting tensors are concatenated along the feature dimension. The number of filters for every n-gram is set to num_features/num_n-grams. Thus, the resulting tensor has the same shape as the input tensor.

Sub-sequences that are masked in the input are also masked in the output.

The output of this layer can be understand as the results of the n-gram filters. Stacking this layer allows the model to perform n-gram detection of n-grams (meta perspective). All parameters with the prefix *ng_conv_* can be used to configure this layer.

Recurrent Layers

Description

A regular recurrent layer either as Gated Recurrent Unit (GRU) or Long Short-Term Memory (LSTM) layer. Uses PyTorchs implementation. All parameters with the prefix *rec_* can be used to configure this layer.

Classification Pooling Layer

Description

Layer transforms sequences into a lower dimensional space that can be passed to dense layers. It performs two types of pooling. First, it extractes features across the time dimension selecting the maximal and/or minimal features. Second, it performs pooling over the remaining features selecting a specific number of the highest and/or lowest features.

In the case of selecting the minimal *and* maximal features at the same time the minimal features are concatenated to the tensor of the maximal features resulting the in the shape \$(Batch, Times, 2*Features)\$ at the end of the first step. In the second step the number of requested features is halved. The first half is used for the maximal features and the second for the minimal features. All parameters with the prefix *cls_pooling_* can be used to configure this layer.

Training and Prediction

For the creation and training of a classifier an object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) on the one hand and a [factor](#) on the other hand are necessary.

The object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) contains the numerical text representations (text embeddings) of the raw texts generated by an object of class [TextEmbeddingModel](#). For supporting large data sets it is recommended to use [LargeDataSetForTextEmbeddings](#) instead of [EmbeddedText](#).

The factor contains the classes/categories for every text. Missing values (unlabeled cases) are supported and can be used for pseudo labeling.

For predictions an object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) has to be used which was created with the same [TextEmbeddingModel](#) as for training..

Value

Returns a new object of this class ready for configuration or for loading a saved classifier.

Super classes

```
aifedducation::AIFEBaseModel -> aifedducation::ModelsBasedOnTextEmbeddings -> aifedducation::ClassifiersB  
-> aifedducation::TEClassifiersBasedOnProtoNet -> TEClassifierSequentialPrototype
```

Methods

Public methods:

- [TEClassifierSequentialPrototype\\$configure\(\)](#)
- [TEClassifierSequentialPrototype\\$clone\(\)](#)

Method [configure\(\)](#): Creating a new instance of this class.

Usage:

```
TEClassifierSequentialPrototype$configure(  
  name = NULL,  
  label = NULL,  
  text_embeddings = NULL,  
  feature_extractor = NULL,  
  target_levels = NULL,  
  skip_connection_type = "ResidualGate",  
  cls_pooling_features = 50,  
  cls_pooling_type = "MinMax",  
  metric_type = "Euclidean",  
  feat_act_fct = "ELU",  
  feat_size = 50,  
  feat_bias = TRUE,  
  feat_dropout = 0,  
  feat_parametrizations = "None",  
  feat_normalization_type = "LayerNorm",  
  ng_conv_act_fct = "ELU",  
  ng_conv_n_layers = 1,  
  ng_conv_ks_min = 2,  
  ng_conv_ks_max = 4,  
  ng_conv_bias = FALSE,  
  ng_conv_dropout = 0.1,  
  ng_conv_parametrizations = "None",  
  ng_conv_normalization_type = "LayerNorm",  
  ng_conv_residual_type = "ResidualGate",  
  dense_act_fct = "ELU",
```

```

dense_n_layers = 1,
dense_dropout = 0.5,
dense_bias = FALSE,
dense_parametrizations = "None",
dense_normalization_type = "LayerNorm",
dense_residual_type = "ResidualGate",
rec_act_fct = "Tanh",
rec_n_layers = 1,
rec_type = "GRU",
rec_bidirectional = FALSE,
rec_dropout = 0.2,
rec_bias = FALSE,
rec_parametrizations = "None",
rec_normalization_type = "LayerNorm",
rec_residual_type = "ResidualGate",
tf_act_fct = "ELU",
tf_dense_dim = 50,
tf_n_layers = 1,
tf_dropout_rate_1 = 0.1,
tf_dropout_rate_2 = 0.5,
tf_attention_type = "MultiHead",
tf_positional_type = "absolute",
tf_num_heads = 1,
tf_bias = FALSE,
tf_parametrizations = "None",
tf_normalization_type = "LayerNorm",
tf_residual_type = "ResidualGate",
embedding_dim = 2
)

```

Arguments:

name string Name of the new model. Please refer to common name conventions. Free text can be used with parameter **label**. If set to NULL a unique ID is generated automatically.
 Allowed values: any

label string Label for the new model. Here you can use free text. Allowed values: any

text_embeddings EmbeddedText, LargeDataSetForTextEmbeddings Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#).

feature_extractor TFEFeatureExtractor Object of class [TFEFeatureExtractor](#) which should be used in order to reduce the number of dimensions of the text embeddings. If no feature extractor should be applied set NULL.

target_levels vector containing the levels (categories or classes) within the target data.
 Please note that order matters. For ordinal data please ensure that the levels are sorted correctly with later levels indicating a higher category/class. For nominal data the order does not matter.

skip_connection_type string Type of residual connection for all layers and stack of layers.
 Allowed values: 'ResidualGate', 'Addition', 'None'

cls_pooling_features int Number of features to be extracted at the end of the model. Allowed values: $1 \leq x$

cls_pooling_type string Type of extracting intermediate features. Allowed values: 'Max', 'Min', 'MinMax'
 metric_type string Type of metric used for calculating the distance. Allowed values: 'Euclidean'
 feat_act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU', 'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'
 feat_size int Number of neurons for each dense layer. Allowed values: $2 \leq x$
 feat_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.
 feat_dropout double determining the dropout for the dense projection of the feature layer. Allowed values: $0 \leq x \leq 0.6$
 feat_parametrizations string Re-Parametrizations of the weights of layers. Allowed values: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
 feat_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'
 ng_conv_act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU', 'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'
 ng_conv_n_layers int determining how many times the n-gram layers should be added to the network. Allowed values: $0 \leq x$
 ng_conv_ks_min int determining the minimal window size for n-grams. Allowed values: $2 \leq x$
 ng_conv_ks_max int determining the maximal window size for n-grams. Allowed values: $2 \leq x$
 ng_conv_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.
 ng_conv_dropout double determining the dropout for n-gram convolution layers. Allowed values: $0 \leq x \leq 0.6$
 ng_conv_parametrizations string Re-Parametrizations of the weights of layers. Allowed values: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
 ng_conv_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'
 ng_conv_residual_type string Type of residual connection for all layers and stack of layers. Allowed values: 'ResidualGate', 'Addition', 'None'
 dense_act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU', 'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'
 dense_n_layers int Number of dense layers. Allowed values: $0 \leq x$
 dense_dropout double determining the dropout between dense layers. Allowed values: $0 \leq x \leq 0.6$
 dense_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the layers.
 dense_parametrizations string Re-Parametrizations of the weights of layers. Allowed values: 'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
 dense_normalization_type string Type of normalization applied to all layers and stack layers. Allowed values: 'LayerNorm', 'None'
 dense_residual_type string Type of residual connection for all layers and stack of layers. Allowed values: 'ResidualGate', 'Addition', 'None'

```

rec_act_fct string Activation function for all layers. Allowed values: 'Tanh'
rec_n_layers int Number of recurrent layers. Allowed values: 0 <= x
rec_type string Type of the recurrent layers. rec_type='GRU' for Gated Recurrent Unit and
    rec_type='LSTM' for Long Short-Term Memory. Allowed values: 'GRU', 'LSTM'
rec_bidirectional bool If TRUE a bidirectional version of the recurrent layers is used.
rec_dropout double determining the dropout between recurrent layers. Allowed values: 0 <= x <= 0.6
rec_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the
    layers.
rec_parametrizations string Re-Parametrizations of the weights of layers. Allowed val-
        ues: 'None'
rec_normalization_type string Type of normalization applied to all layers and stack layers.
    Allowed values: 'LayerNorm', 'None'
rec_residual_type string Type of residual connection for all layers and stack of layers.
    Allowed values: 'ResidualGate', 'Addition', 'None'
tf_act_fct string Activation function for all layers. Allowed values: 'ELU', 'LeakyReLU',
    'ReLU', 'GELU', 'Sigmoid', 'Tanh', 'PReLU'
tf_dense_dim int determining the size of the projection layer within a each transformer en-
    coder. Allowed values: 1 <= x
tf_n_layers int determining how many times the encoder should be added to the network.
    Allowed values: 0 <= x
tf_dropout_rate_1 double determining the dropout after the attention mechanism within the
    transformer encoder layers. Allowed values: 0 <= x <= 0.6
tf_dropout_rate_2 double determining the dropout for the dense projection within the trans-
    former encoder layers. Allowed values: 0 <= x <= 0.6
tf_attention_type string Choose the attention type. Allowed values: 'Fourier', 'Multi-
    Head'
tf_positional_type string Type of processing positional information. Allowed values: 'ab-
        solute'
tf_num_heads int determining the number of attention heads for a self-attention layer. Only
    relevant if attention_type='multihead' Allowed values: 0 <= x
tf_bias bool If TRUE a bias term is added to all layers. If FALSE no bias term is added to the
    layers.
tf_parametrizations string Re-Parametrizations of the weights of layers. Allowed values:
    'None', 'OrthogonalWeights', 'WeightNorm', 'SpectralNorm'
tf_normalization_type string Type of normalization applied to all layers and stack layers.
    Allowed values: 'LayerNorm', 'None'
tf_residual_type string Type of residual connection for all layers and stack of layers.
    Allowed values: 'ResidualGate', 'Addition', 'None'
embedding_dim int determining the number of dimensions for the embedding. Allowed val-
    ues: 2 <= x

```

Returns: Function does nothing return. It modifies the current object.

Method clone(): The objects of this class are cloneable with this method.

Usage:

```
TEClassifierSequentialPrototype$clone(deep = FALSE)
```

Arguments:

deep Whether to make a deep clone.

References

- Oreshkin, B. N., Rodriguez, P. & Lacoste, A. (2018). TADAM: Task dependent adaptive metric for improved few-shot learning. <https://doi.org/10.48550/arXiv.1805.10123>
- Snell, J., Swersky, K. & Zemel, R. S. (2017). Prototypical Networks for Few-shot Learning. <https://doi.org/10.48550/arXiv.1703.05175>
- Zhang, X., Nie, J., Zong, L., Yu, H. & Liang, W. (2019). One Shot Learning with Margin. In Q. Yang, Z.-H. Zhou, Z. Gong, M.-L. Zhang & S.-J. Huang (Eds.), Lecture Notes in Computer Science. Advances in Knowledge Discovery and Data Mining (Vol. 11440, pp. 305–317). Springer International Publishing. https://doi.org/10.1007/978-3-030-16145-3_24

See Also

Other Classification: [TEClassifierParallel](#), [TEClassifierParallelPrototype](#), [TEClassifierProtoNet](#), [TEClassifierRegular](#), [TEClassifierSequential](#)

TEFeatureExtractor	<i>Feature extractor for reducing the number for dimensions of text embeddings.</i>
--------------------	---

Description

Abstract class for auto encoders with 'pytorch'.

Objects of this class are used for reducing the number of dimensions of text embeddings created by an object of class [TextEmbeddingModel](#).

For training an object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) generated by an object of class [TextEmbeddingModel](#) is necessary. Passing raw texts is not supported.

For prediction an ob object class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) is necessary that was generated with the same [TextEmbeddingModel](#) as during training. Prediction outputs a new object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) which contains a text embedding with a lower number of dimensions.

All models use tied weights for the encoder and decoder layers (except method="LSTM") and apply the estimation of orthogonal weights. In addition, training tries to train the model to achieve uncorrelated features.

Objects of class [TEFeatureExtractor](#) are designed to be used with classifiers such as [TEClassifierRegular](#) and [TEClassifierProtoNet](#).

Value

A new instances of this class.

Super classes

[aifedducation::AIFEBaseModel](#) -> [aifedducation::ModelsBasedOnTextEmbeddings](#) -> TEFeatureExtractor

Methods

Public methods:

- [TEFeatureExtractor\\$configure\(\)](#)
- [TEFeatureExtractor\\$train\(\)](#)
- [TEFeatureExtractor\\$extract_features\(\)](#)
- [TEFeatureExtractor\\$extract_features_large\(\)](#)
- [TEFeatureExtractor\\$plot_training_history\(\)](#)
- [TEFeatureExtractor\\$clone\(\)](#)

Method [configure\(\)](#): Creating a new instance of this class.

Usage:

```
TEFeatureExtractor$configure(
  name = NULL,
  label = NULL,
  text_embeddings = NULL,
  features = 128,
  method = "dense",
  orthogonal_method = "matrix_exp",
  noise_factor = 0.2
)
```

Arguments:

name string Name of the new model. Please refer to common name conventions. Free text can be used with parameter **label**. If set to NULL a unique ID is generated automatically.

Allowed values: any

label string Label for the new model. Here you can use free text. Allowed values: any

text_embeddings EmbeddedText, LargeDataSetForTextEmbeddings Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#).

features int Number of features the model should use. Allowed values: $1 \leq x$

method string Method to use for the feature extraction. 'lstm' for an extractor based on LSTM-layers or 'Dense' for dense layers. Allowed values: 'Dense', 'LSTM'

orthogonal_method string Method to use for the feature extraction. 'lstm' for an extractor based on LSTM-layers or 'Dense' for dense layers. Allowed values: 'Dense', 'LSTM'

noise_factor double Value between 0 and a value lower 1 indicating how much noise should be added to the input during training. Allowed values: $0 \leq x \leq 1$

Returns: Returns an object of class [TEFeatureExtractor](#) which is ready for training.

Method [train\(\)](#): Method for training a neural net.

Usage:

```
TEFeatureExtractor$train(
  data_embeddings = NULL,
  data_val_size = 0.25,
  sustain_track = TRUE,
  sustain_iso_code = NULL,
  sustain_region = NULL,
  sustain_interval = 15,
  epochs = 40,
  batch_size = 32,
  trace = TRUE,
  ml_trace = 1,
  log_dir = NULL,
  log_write_interval = 10,
  lr_rate = 0.001,
  lr_warm_up_ratio = 0.02,
  optimizer = "AdamW"
)
```

Arguments:

`data_embeddings` Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#).

`data_val_size` double between 0 and 1, indicating the proportion of cases which should be used for the validation sample.

`sustain_track` bool If TRUE energy consumption is tracked during training via the python library 'codecarbon'.

`sustain_iso_code` string ISO code (Alpha-3-Code) for the country. This variable must be set if sustainability should be tracked. A list can be found on Wikipedia: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes.

`sustain_region` Region within a country. Only available for USA and Canada See the documentation of 'codecarbon' for more information. <https://mlco2.github.io/codecarbon/parameters.html>

`sustain_interval` int Interval in seconds for measuring power usage.

`epochs` int Number of training epochs.

`batch_size` int Size of batches.

`trace` bool TRUE, if information about the estimation phase should be printed to the console.

`ml_trace` int `ml_trace=0` does not print any information about the training process from pytorch on the console. `ml_trace=1` prints a progress bar.

`log_dir` string Path to the directory where the log files should be saved. If no logging is desired set this argument to NULL.

`log_write_interval` int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if `log_dir` is not NULL.

`lr_rate` double Initial learning rate for the training.

`lr_warm_up_ratio` double Number of epochs used for warm up.

`optimizer` string determining the optimizer used for training. Allowed values: 'Adam', 'RMSprop', 'AdamW', 'SGD'

Returns: Function does not return a value. It changes the object into a trained classifier.

Method extract_features(): Method for extracting features. Applying this method reduces the number of dimensions of the text embeddings. Please note that this method should only be used if a small number of cases should be compressed since the data is loaded completely into memory. For a high number of cases please use the method `extract_features_large`.

Usage:

```
TEFeatureExtractor$extract_features(data_embeddings, batch_size)
```

Arguments:

`data_embeddings` Object of class [EmbeddedText](#), [LargeDataSetForTextEmbeddings](#), `datasets.arrow_dataset.Dataset` or array containing the text embeddings which should be reduced in their dimensions.

`batch_size` int batch size.

Returns: Returns an object of class [EmbeddedText](#) containing the compressed embeddings.

Method extract_features_large(): Method for extracting features from a large number of cases. Applying this method reduces the number of dimensions of the text embeddings.

Usage:

```
TEFeatureExtractor$extract_features_large(
  data_embeddings,
  batch_size,
  trace = FALSE
)
```

Arguments:

`data_embeddings` Object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#) containing the text embeddings which should be reduced in their dimensions.

`batch_size` int batch size.

`trace` bool If TRUE information about the progress is printed to the console.

Returns: Returns an object of class [LargeDataSetForTextEmbeddings](#) containing the compressed embeddings.

Method plot_training_history(): Method for requesting a plot of the training history. This method requires the R package 'ggplot2' to work.

Usage:

```
TEFeatureExtractor$plot_training_history(
  y_min = NULL,
  y_max = NULL,
  text_size = 10
)
```

Arguments:

`y_min` Minimal value for the y-axis. Set to NULL for an automatic adjustment.

`y_max` Maximal value for the y-axis. Set to NULL for an automatic adjustment.

`text_size` Size of the text.

Returns: Returns a plot of class ggplot visualizing the training process.

Method clone(): The objects of this class are cloneable with this method.

Usage:

```
TEFeatureExtractor$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

Note

`features` refers to the number of features for the compressed text embeddings.

This model requires `pad_value=0`. If this condition is not met the padding value is switched automatically.

This model requires that the underlying [TextEmbeddingModel](#) uses `pad_value=0`. If this condition is not met the pad value is switched before training.

See Also

Other Text Embedding: [TextEmbeddingModel](#)

`tensor_list_to_numpy` *Convert list of tensors into numpy arrays*

Description

Function converts tensors within a list into numpy arrays in order to allow further operations in *R*.

Usage

```
tensor_list_to_numpy(tensor_list)
```

Arguments

`tensor_list` list of objects.

Value

Returns the same list with the exception that objects of class `torch.Tensor` are transformed into numpy arrays. If the tensor requires a gradient and/or is on gpu it is detached and converted. If the object in a list is not of this class the original object is returned.

See Also

Other Utils Python Data Management Developers: [class_vector_to_py_dataset\(\)](#), [data.frame_to_py_dataset\(\)](#), [get_batches_index\(\)](#), [prepare_r_array_for_dataset\(\)](#), [py_dataset_to_embeddings\(\)](#), [reduce_to_unique\(\)](#), [tensor_to_numpy\(\)](#)

`tensor_to_matrix_c` *Transform tensor to matrix*

Description

Function written in C++ for transformation the tensor (with size batch x times x features) to the matrix (with size batch x times*features)

Usage

```
tensor_to_matrix_c(tensor, times, features)
```

Arguments

<code>tensor</code>	3-D array (cube) data as tensor (with size batch x times x features)
<code>times</code>	unsigned integer times number
<code>features</code>	unsigned integer features number

Value

Returns matrix (with size batch x times*features)

See Also

Other Utils Developers: [auto_n_cores\(\)](#), [create_object\(\)](#), [create_synthetic_units_from_matrix\(\)](#), [generate_id\(\)](#), [get_n_chunks\(\)](#), [get_synthetic_cases_from_matrix\(\)](#), [matrix_to_array_c\(\)](#), [to_categorical_c\(\)](#)

`tensor_to_numpy` *Tensor_to_numpy*

Description

Function converts a tensor into a numpy array in order to allow further operations in *R*.

Usage

```
tensor_to_numpy(object)
```

Arguments

<code>object</code>	Object of any class.
---------------------	----------------------

Value

In the case the object is of class `torch.Tensor` it returns a numpy error. If the tensor requires a gradient and/or is on gpu it is detached and converted. If the object is not of class `torch.Tensor` the original object is returned.

See Also

Other Utils Python Data Management Developers: `class_vector_to_py_dataset()`, `data.frame_to_py_dataset()`, `get_batches_index()`, `prepare_r_array_for_dataset()`, `py_dataset_to_embeddings()`, `reduce_to_unique()`, `tensor_list_to_numpy()`

TextEmbeddingModel *Text embedding model*

Description

This R6 class stores a text embedding model which can be used to tokenize, encode, decode, and embed raw texts. The object provides a unique interface for different text processing methods.

Value

Objects of class `TextEmbeddingModel` transform raw texts into numerical representations which can be used for downstream tasks. For this aim objects of this class allow to tokenize raw texts, to encode tokens to sequences of integers, and to decode sequences of integers back to tokens.

Public fields

`last_training ('list()')`

List for storing the history and the results of the last training. This information will be overwritten if a new training is started.

`tokenizer_statistics ('matrix()')`

Matrix containing the tokenizer statistics for the creation of the tokenizer and all training runs according to Kaya & Tantug (2024).

Kaya, Y. B., & Tantug, A. C. (2024). Effect of tokenization granularity for Turkish large language models. Intelligent Systems with Applications, 21, 200335. <https://doi.org/10.1016/j.iswa.2024.200335>

Methods

Public methods:

- `TextEmbeddingModel$configure()`
- `TextEmbeddingModel$load_from_disk()`
- `TextEmbeddingModel$load()`
- `TextEmbeddingModel$save()`
- `TextEmbeddingModel$encode()`
- `TextEmbeddingModel$decode()`

- `TextEmbeddingModel$get_special_tokens()`
- `TextEmbeddingModel$embed()`
- `TextEmbeddingModel$embed_large()`
- `TextEmbeddingModel$fill_mask()`
- `TextEmbeddingModel$set_publication_info()`
- `TextEmbeddingModel$get_publication_info()`
- `TextEmbeddingModel$set_model_license()`
- `TextEmbeddingModel$get_model_license()`
- `TextEmbeddingModel$set_documentation_license()`
- `TextEmbeddingModel$get_documentation_license()`
- `TextEmbeddingModel$set_model_description()`
- `TextEmbeddingModel$get_model_description()`
- `TextEmbeddingModel$get_model_info()`
- `TextEmbeddingModel$get_package_versions()`
- `TextEmbeddingModel$get_basic_components()`
- `TextEmbeddingModel$get_n_features()`
- `TextEmbeddingModel$get_transformer_components()`
- `TextEmbeddingModel$get_sustainability_data()`
- `TextEmbeddingModel$get_ml_framework()`
- `TextEmbeddingModel$get_pad_value()`
- `TextEmbeddingModel$count_parameter()`
- `TextEmbeddingModel$is_configured()`
- `TextEmbeddingModel$get_private()`
- `TextEmbeddingModel$get_all_fields()`
- `TextEmbeddingModel$plot_training_history()`
- `TextEmbeddingModel$clone()`

Method `configure()`: Method for creating a new text embedding model

Usage:

```
TextEmbeddingModel$configure(
  model_name = NULL,
  model_label = NULL,
  model_language = NULL,
  max_length = 0,
  chunks = 2,
  overlap = 0,
  emb_layer_min = "Middle",
  emb_layer_max = "2_3_layer",
  emb_pool_type = "Average",
  pad_value = -100,
  model_dir = NULL,
  trace = FALSE
)
```

Arguments:

`model_name` string containing the name of the new model.
`model_label` string containing the label/title of the new model.
`model_language` string containing the language which the model represents (e.g., English).
`max_length` int determining the maximum length of token sequences used in transformer models. Not relevant for the other methods.
`chunks` int Maximum number of chunks. Must be at least 2.
`overlap` int determining the number of tokens which should be added at the beginning of the next chunk. Only relevant for transformer models.
`emb_layer_min` int or string determining the first layer to be included in the creation of embeddings. An integer corresponds to the layer number. The first layer has the number 1. Instead of an integer the following strings are possible: "start" for the first layer, "Middle" for the middle layer, "2_3_layer" for the layer two-third layer, and "Last" for the last layer.
`emb_layer_max` int or string determining the last layer to be included in the creation of embeddings. An integer corresponds to the layer number. The first layer has the number 1. Instead of an integer the following strings are possible: "start" for the first layer, "Middle" for the middle layer, "2_3_layer" for the layer two-third layer, and "Last" for the last layer.
`emb_pool_type` string determining the method for pooling the token embeddings within each layer. If "CLS" only the embedding of the CLS token is used. If "Average" the token embedding of all tokens are averaged (excluding padding tokens). "cls" is not supported for `method="funnel"`.
`pad_value` int Value indicating padding. This value should not be in the range of regular values for computations. Thus it is not recommended to change this value. Default is -100. Allowed values: $x \leq -100$
`model_dir` string path to the directory where the BERT model is stored.
`trace` bool TRUE prints information about the progress. FALSE does not.

Returns: Returns an object of class [TextEmbeddingModel](#).

Method `load_from_disk()`: loads an object from disk and updates the object to the current version of the package.

Usage:

```
TextEmbeddingModel$load_from_disk(dir_path)
```

Arguments:

`dir_path` Path where the object set is stored.

Returns: Method does not return anything. It loads an object from disk.

Method `load()`: Method for loading a transformers model into R.

Usage:

```
TextEmbeddingModel$load(dir_path)
```

Arguments:

`dir_path` string containing the path to the relevant model directory.

Returns: Function does not return a value. It is used for loading a saved transformer model into the R interface.

Method `save()`: Method for saving a transformer model on disk. Relevant only for transformer models.

Usage:

```
TextEmbeddingModel$save(dir_path, folder_name)
```

Arguments:

`dir_path` string containing the path to the relevant model directory.

`folder_name` string Name for the folder created within the directory. This folder contains all model files.

Returns: Function does not return a value. It is used for saving a transformer model to disk.

Method `encode()`: Method for encoding words of raw texts into integers.

Usage:

```
TextEmbeddingModel$encode(
  raw_text,
  token_encodings_only = FALSE,
  to_int = TRUE,
  trace = FALSE
)
```

Arguments:

`raw_text` vector containing the raw texts.

`token_encodings_only` bool If TRUE, only the token encodings are returned. If FALSE, the complete encoding is returned which is important for some transformer models.

`to_int` bool If TRUE the integer ids of the tokens are returned. If FALSE the tokens are returned. Argument only applies for transformer models and if `token_encodings_only`=TRUE.

`trace` bool If TRUE, information of the progress is printed. FALSE if not requested.

Returns: list containing the integer or token sequences of the raw texts with special tokens.

Method `decode()`: Method for decoding a sequence of integers into tokens

Usage:

```
TextEmbeddingModel$decode(int_sequence, to_token = FALSE)
```

Arguments:

`int_sequence` list containing the integer sequences which should be transformed to tokens or plain text.

`to_token` bool If FALSE plain text is returned. If TRUE a sequence of tokens is returned. Argument only relevant if the model is based on a transformer.

Returns: list of token sequences

Method `get_special_tokens()`: Method for receiving the special tokens of the model

Usage:

```
TextEmbeddingModel$get_special_tokens()
```

Returns: Returns a matrix containing the special tokens in the rows and their type, token, and id in the columns.

Method `embed()`: Method for creating text embeddings from raw texts. This method should only be used if a small number of texts should be transformed into text embeddings. For a large number of texts please use the method `embed_large`.

Usage:

```
TextEmbeddingModel$embed(
  raw_text = NULL,
  doc_id = NULL,
  batch_size = 8,
  trace = FALSE,
  return_large_dataset = FALSE
)
```

Arguments:

`raw_text` vector containing the raw texts.

`doc_id` vector containing the corresponding IDs for every text.

`batch_size` int determining the maximal size of every batch.

`trace` bool TRUE, if information about the progression should be printed on console.

`return_large_dataset` 'bool' If TRUE the retuned object is of class [LargeDataSetForTextEmbeddings](#). If FALSE it is of class [EmbeddedText](#)

Returns: Method returns an object of class [EmbeddedText](#) or [LargeDataSetForTextEmbeddings](#). This object contains the embeddings as a `data.frame` and information about the model creating the embeddings.

Method `embed_large()`: Method for creating text embeddings from raw texts.

Usage:

```
TextEmbeddingModel$embed_large(
  large_datas_set,
  batch_size = 32,
  trace = FALSE,
  log_file = NULL,
  log_write_interval = 2
)
```

Arguments:

`large_datas_set` Object of class [LargeDataSetForText](#) containing the raw texts.

`batch_size` int determining the maximal size of every batch.

`trace` bool TRUE, if information about the progression should be printed on console.

`log_file` string Path to the file where the log should be saved. If no logging is desired set this argument to NULL.

`log_write_interval` int Time in seconds determining the interval in which the logger should try to update the log files. Only relevant if `log_file` is not NULL.

Returns: Method returns an object of class [LargeDataSetForTextEmbeddings](#).

Method `fill_mask()`: Method for calculating tokens behind mask tokens.

Usage:

```
TextEmbeddingModel$fill_mask(text, n_solutions = 5)
```

Arguments:

`text` string Text containing mask tokens.
`n_solutions` int Number estimated tokens for every mask.

Returns: Returns a list containing a `data.frame` for every mask. The `data.frame` contains the solutions in the rows and reports the score, token id, and token string in the columns.

Method `set_publication_info()`: Method for setting the bibliographic information of the model.

Usage:

```
TextEmbeddingModel$set_publication_info(type, authors, citation, url = NULL)
```

Arguments:

`type` string Type of information which should be changed/added. `developer`, and `modifier` are possible.

`authors` List of people.

`citation` string Citation in free text.

`url` string Corresponding URL if applicable.

Returns: Function does not return a value. It is used to set the private members for publication information of the model.

Method `get_publication_info()`: Method for getting the bibliographic information of the model.

Usage:

```
TextEmbeddingModel$get_publication_info()
```

Returns: list of bibliographic information.

Method `set_model_license()`: Method for setting the license of the model

Usage:

```
TextEmbeddingModel$set_model_license(license = "CC BY")
```

Arguments:

`license` string containing the abbreviation of the license or the license text.

Returns: Function does not return a value. It is used for setting the private member for the software license of the model.

Method `get_model_license()`: Method for requesting the license of the model

Usage:

```
TextEmbeddingModel$get_model_license()
```

Returns: string License of the model

Method `set_documentation_license()`: Method for setting the license of models' documentation.

Usage:

```
TextEmbeddingModel$set_documentation_license(license = "CC BY")
```

Arguments:

license string containing the abbreviation of the license or the license text.

Returns: Function does not return a value. It is used to set the private member for the documentation license of the model.

Method get_documentation_license(): Method for getting the license of the models' documentation.

Usage:

```
TextEmbeddingModel$get_documentation_license()
```

Arguments:

license string containing the abbreviation of the license or the license text.

Method set_model_description(): Method for setting a description of the model

Usage:

```
TextEmbeddingModel$set_model_description(  
  eng = NULL,  
  native = NULL,  
  abstract_eng = NULL,  
  abstract_native = NULL,  
  keywords_eng = NULL,  
  keywords_native = NULL  
)
```

Arguments:

eng string A text describing the training of the classifier, its theoretical and empirical background, and the different output labels in English.

native string A text describing the training of the classifier, its theoretical and empirical background, and the different output labels in the native language of the model.

abstract_eng string A text providing a summary of the description in English.

abstract_native string A text providing a summary of the description in the native language of the classifier.

keywords_eng vectorof keywords in English.

keywords_native vectorof keywords in the native language of the classifier.

Returns: Function does not return a value. It is used to set the private members for the description of the model.

Method get_model_description(): Method for requesting the model description.

Usage:

```
TextEmbeddingModel$get_model_description()
```

Returns: list with the description of the model in English and the native language.

Method get_model_info(): Method for requesting the model information

Usage:

```
TextEmbeddingModel$get_model_info()
```

Returns: list of all relevant model information

Method `get_package_versions()`: Method for requesting a summary of the R and python packages' versions used for creating the model.

Usage:

```
TextEmbeddingModel$get_package_versions()
```

Returns: Returns a list containing the versions of the relevant R and python packages.

Method `get_basic_components()`: Method for requesting the part of interface's configuration that is necessary for all models.

Usage:

```
TextEmbeddingModel$get_basic_components()
```

Returns: Returns a list.

Method `get_n_features()`: Method for requesting the number of features.

Usage:

```
TextEmbeddingModel$get_n_features()
```

Returns: Returns a double which represents the number of features. This number represents the hidden size of the embeddings for every chunk or time.

Method `get_transformer_components()`: Method for requesting the part of interface's configuration that is necessary for transformer models.

Usage:

```
TextEmbeddingModel$get_transformer_components()
```

Returns: Returns a list.

Method `get_sustainability_data()`: Method for requesting a log of tracked energy consumption during training and an estimate of the resulting CO2 equivalents in kg.

Usage:

```
TextEmbeddingModel$get_sustainability_data()
```

Returns: Returns a matrix containing the tracked energy consumption, CO2 equivalents in kg, information on the tracker used, and technical information on the training infrastructure for every training run.

Method `get_ml_framework()`: Method for requesting the machine learning framework used for the classifier.

Usage:

```
TextEmbeddingModel$get_ml_framework()
```

Returns: Returns a string describing the machine learning framework used for the classifier.

Method `get_pad_value()`: Value for indicating padding.

Usage:

```
TextEmbeddingModel$get_pad_value()
```

Returns: Returns an int describing the value used for padding.

Method `count_parameter()`: Method for counting the trainable parameters of a model.

Usage:

```
TextEmbeddingModel$count_parameter(with_head = FALSE)
```

Arguments:

`with_head` bool If TRUE the number of parameters is returned including the language modeling head of the model. If FALSE only the number of parameters of the core model is returned.

Returns: Returns the number of trainable parameters of the model.

Method `is_configured()`: Method for checking if the model was successfully configured. An object can only be used if this value is TRUE.

Usage:

```
TextEmbeddingModel$is_configured()
```

Returns: bool TRUE if the model is fully configured. FALSE if not.

Method `get_private()`: Method for requesting all private fields and methods. Used for loading and updating an object.

Usage:

```
TextEmbeddingModel$get_private()
```

Returns: Returns a list with all private fields and methods.

Method `get_all_fields()`: Return all fields.

Usage:

```
TextEmbeddingModel$get_all_fields()
```

Returns: Method returns a list containing all public and private fields of the object.

Method `plot_training_history()`: Method for requesting a plot of the training history. This method requires the R package 'ggplot2' to work.

Usage:

```
TextEmbeddingModel$plot_training_history(y_min = NULL, y_max = NULL)
```

Arguments:

`y_min` Minimal value for the y-axis. Set to NULL for an automatic adjustment.

`y_max` Maximal value for the y-axis. Set to NULL for an automatic adjustment.

Returns: Returns a plot of class ggplot visualizing the training process.

Method `clone()`: The objects of this class are cloneable with this method.

Usage:

```
TextEmbeddingModel$clone(deep = FALSE)
```

Arguments:

`deep` Whether to make a deep clone.

See Also

Other Text Embedding: [TEFeatureExtractor](#)

`to_categorical_c` *Transforming classes to one-hot encoding*

Description

Function transforming a vector of classes (int) into a binary class matrix.

Usage

```
to_categorical_c(class_vector, n_classes)
```

Arguments

<code>class_vector</code>	vector containing integers for every class. The integers must range from 0 to <code>n_classes</code> -1.
<code>n_classes</code>	int Total number of classes.

Value

Returns a `matrix` containing the binary representation for every class.

See Also

Other Utils Developers: [auto_n_cores\(\)](#), [create_object\(\)](#), [create_synthetic_units_from_matrix\(\)](#), [generate_id\(\)](#), [get_n_chunks\(\)](#), [get_synthetic_cases_from_matrix\(\)](#), [matrix_to_array_c\(\)](#), [tensor_to_matrix_c\(\)](#)

`update_aifedducation` *Updates an existing installation of 'aifedducation' on a machine*

Description

Function for updating 'aifedducation' on a machine.

The function tries to find an existing environment on the machine, removes the environment and installs the environment with the new python modules.

In the case `env_type = "auto"` the function tries to update an existing virtual environment. If no virtual environment exists it tries to update a conda environment.

Usage

```
update_aifedducation(
    update_aifedducation_studio = TRUE,
    env_type = "auto",
    cuda_version = "12.4",
    envname = "aifedducation"
)
```

Arguments

update_aifedducation_studio	bool If TRUE all necessary R packages are installed for using AI for Education Studio.
env_type	string If set to "venv" virtual environment is requested. If set to "conda" a 'conda' environment is requested. If set to "auto" the function tries to use a virtual environment with the given name. If this environment does not exist it tries to activate a conda environment with the given name. If this fails the default virtual environment is used.'
cuda_version	string determining the requested version of cuda.
envname	string Name of the environment where the packages should be installed.

Value

Function does nothing return. It installs python, optional R packages, and necessary 'python' packages on a machine.

Note

On MAC OS torch will be installed without support for cuda.

See Also

Other Installation and Configuration: [check_aif_py_modules\(\)](#), [install_aifedducation\(\)](#), [install_aifedducation_studio\(\)](#), [install_py_modules\(\)](#), [prepare_session\(\)](#), [set_transformers_logger\(\)](#)

`write_log`

Write log

Description

Function for writing a log file from R containing three rows and three columns. The log file can report the current status of maximal three processes. The first row describes the top process. The second row describes the status of the process within the top process. The third row can be used to describe the status of a process within the middle process.

The log can be read with [read_log](#).

Usage

```
write_log(
  log_file,
  value_top = 0,
  total_top = 1,
  message_top = NA,
  value_middle = 0,
  total_middle = 1,
```

```

message_middle = NA,
value_bottom = 0,
total_bottom = 1,
message_bottom = NA,
last_log = NULL,
write_interval = 2
)

```

Arguments

<code>log_file</code>	string Path to the file where the log should be saved and updated.
<code>value_top</code>	double Current value for the top process.
<code>total_top</code>	double Maximal value for the top process.
<code>message_top</code>	string Message describing the current state of the top process.
<code>value_middle</code>	double Current value for the middle process.
<code>total_middle</code>	double Maximal value for the middle process.
<code>message_middle</code>	string Message describing the current state of the middle process.
<code>value_bottom</code>	double Current value for the bottom process.
<code>total_bottom</code>	double Maximal value for the bottom process.
<code>message_bottom</code>	string Message describing the current state of the bottom process.
<code>last_log</code>	POSIXct Time when the last log was created. If there is no log file set this value to NULL.
<code>write_interval</code>	int Time in seconds. This time must be past before a new log is created.

Value

This function writes a log file to the given location. If `log_file` is `NULL` the function will not try to write a log file.

If `log_file` is a valid path to a file the function will write a log if the time specified by `write_interval` has passed. In addition the function will return an object of class `POSIXct` describing the time when the log file was successfully updated. If the initial attempt for writing log fails the function returns the value of `last_log` which is `NULL` by default.

See Also

Other Utils Log Developers: [cat_message\(\)](#), [clean_pytorch_log_transformers\(\)](#), [output_message\(\)](#), [print_message\(\)](#), [read_log\(\)](#), [read_loss_log\(\)](#), [reset_log\(\)](#), [reset_loss_log\(\)](#)

Index

- * **Classification**
 - TEClassifierParallel, 131
 - TEClassifierParallelPrototype, 137
 - TEClassifierProtoNet, 144
 - TEClassifierRegular, 148
 - TEClassifierSequential, 159
 - TEClassifierSequentialPrototype, 165
- * **Data Management Developers**
 - DataManagerClassifier, 67
- * **Data Management**
 - EmbeddedText, 73
 - LargeDataSetForText, 106
 - LargeDataSetForTextEmbeddings, 110
- * **Data Sets**
 - imdb_movie_reviews, 97
- * **Graphical User Interface**
 - start_aifeducation_studio, 130
- * **Installation and Configuration**
 - check_aif_py_modules, 56
 - install_aifeducation, 97
 - install_aifeducation_studio, 98
 - install_py_modules, 99
 - prepare_session, 122
 - set_transformers_logger, 129
 - update_aifeducation, 186
- * **Parameter Dictionary**
 - get_called_args, 83
 - get_depr_obj_names, 85
 - get_magnitude_values, 88
 - get_param_def, 90
 - get_param_dict, 91
 - get_param_doc_desc, 92
 - get_TEClassifiers_class_names, 95
- * **R6 Classes for Developers**
 - AIFEBaseModel, 44
 - ClassifiersBasedOnTextEmbeddings, 58
 - LargeDataSetBase, 103
- * **R6 classes for transformers**
 - .AIFEBaseTransformer, 5
 - .AIFEBertTransformer, 14
 - .AIFEFunnelTransformer, 19
 - .AIFELongformerTransformer, 24
 - .AIFEMpnetTransformer, 33
 - .AIFERobertaTransformer, 38
- * **Saving and Loading**
 - load_from_disk, 116
 - save_to_disk, 128
- * **Text Embedding**
 - TEFeatureExtractor, 171
 - TextEmbeddingModel, 177
- * **Transformers for developers**
 - .AIFEModernBertTransformer, 29
 - aife_transformer.load_model_mlm, 49
 - aife_transformer.load_tokenizer, 50
- * **Transformer**
 - aife_transformer.make, 51
 - AIFETrType, 49
- * **Utils Checks Developers**
 - check_all_args, 57
 - check_class_and_type, 57
- * **Utils Developers**
 - auto_n_cores, 51
 - create_object, 65
 - create_synthetic_units_from_matrix, 66
 - generate_id, 80
 - get_n_chunks, 89
 - get_synthetic_cases_from_matrix, 94
 - matrix_to_array_c, 118
 - tensor_to_matrix_c, 176

to_categorical_c, 186

* **Utils Documentation**

- build_documentation_for_model, 52
- build_layer_stack_documentation_for_vignette, 53
- get_desc_for_core_model_architecture, 86
- get_layer_documentation, 87
- get_parameter_documentation, 90

* **Utils File Management Developers**

- create_dir, 65
- get_file_extension, 86

* **Utils Log Developers**

- cat_message, 55
- clean_pytorch_log_transformers, 63
- output_message, 121
- print_message, 123
- read_log, 125
- read_loss_log, 125
- reset_log, 127
- reset_loss_log, 127
- write_log, 187

* **Utils Python Data Management Developers**

- class_vector_to_py_dataset, 63
- data.frame_to_py_dataset, 67
- get_batches_index, 82
- prepare_r_array_for_dataset, 122
- py_dataset_to_embeddings, 123
- reduce_to_unique, 126
- tensor_list_to_numpy, 175
- tensor_to_numpy, 176

* **Utils Python Developers**

- get_py_package_version, 93
- get_py_package_versions, 93
- load_all_py_scripts, 116
- load_py_scripts, 117
- run_py_file, 128

* **Utils Studio Developers**

- add_missing_args, 43
- long_load_target_data, 117
- summarize_args_for_long_task, 130

* **Utils Sustainability Developers**

- get_alpha_3_codes, 82

* **Utils TestThat Developers**

- check_adjust_n_samples_on_CI, 55
- generate_args_for_tests, 79
- generate_embeddings, 80

generate_tensors, 81

get_current_args_for_print, 85

get_fixed_test_tensor, 87

get_test_data_for_classifiers, 95

random_bool_on_CI, 124

* **Utils Transformers Developers**

- calc_tokenizer_statistics, 54

* **datasets**

- AIFETrType, 49
- imdb_movie_reviews, 97

* **oversampling_approaches Developers**

- knnor_is_same_class, 102

* **oversampling_approaches**

- knnor, 101

* **performance measures**

- calc_standard_classification_measures, 53
- cohens_kappa, 64
- fleiss_kappa, 78
- get_coder_metrics, 83
- gwet_ac, 96
- kendalls_w, 100
- kripp_alpha, 102

- .AIFEBaseTransformer, 5, 19, 24, 28, 38, 43, 51
- .AIFEBertTransformer, 14, 14, 24, 28, 38, 43
- .AIFEFunnelTransformer, 14, 19, 19, 28, 38, 43
- .AIFELongformerTransformer, 14, 19, 24, 24, 38, 43
- .AIFEModernBertTransformer, 29, 50
- .AIFEMpnetTransformer, 14, 19, 24, 28, 33, 43
- .AIFERobertaTransformer, 14, 19, 24, 28, 38, 38
- .AIFETrConfig, 33, 50
- .AIFETrModel, 33, 50
- .AIFETrModelMLM, 33, 50
- .AIFETrObj, 33, 50
- .AIFETrTokenizer, 33, 50
- .aife_transformer.check_type, 33, 50

- add_missing_args, 43, 118, 131
- aife_transformer.load_model, 33, 50
- aife_transformer.load_model_config, 33, 50

- aife_transformer.load_model_mlm, 33, 49, 50

aife_transformer.load_tokenizer, 33, 50, 50
aife_transformer.make, 49, 51
aife_transformer.make(), 49
AIFEBaseModel, 44, 62, 106, 121, 156, 159
aifeducation::AIFEBaseTransformer, 14, 19, 24, 29, 34, 39
aifeducation::AIFEBaseModel, 58, 119, 133, 139, 144, 148, 151, 156, 161, 167, 172
aifeducation::ClassifiersBasedOnTextEmbeddings, 133, 139, 144, 148, 151, 156, 161, 167
aifeducation::LargeDataSetBase, 106, 111
aifeducation::ModelsBasedOnTextEmbeddings, 58, 133, 139, 144, 148, 151, 156, 161, 167, 172
aifeducation::TEClassifiersBasedOnProtoNet, 139, 144, 167
aifeducation::TEClassifiersBasedOnRegular, 133, 148, 161
AIFETrType, 49, 50, 51
auto_n_cores, 51, 66, 67, 81, 89, 95, 118, 176, 186

build_documentation_for_model, 52, 53, 86, 88, 90
build_layer_stack_documentation_for_vignette, 52, 53, 86, 88, 90

calc_standard_classification_measures, 53, 64, 78, 84, 97, 101, 103
calc_tokenizer_statistics, 54
cat_message, 55, 64, 121, 123, 125–128, 188
check_adjust_n_samples_on_CI, 55, 79–81, 85, 87, 96, 124
check_aif_py_modules, 56, 98–100, 123, 129, 187
check_all_args, 57, 58
check_class_and_type, 57, 57
class_vector_to_py_dataset, 63, 67, 82, 122, 124, 126, 175, 177
ClassifiersBasedOnTextEmbeddings, 48, 58, 106, 110, 121, 156, 159
clean_pytorch_log_transformers, 55, 63, 121, 123, 125–128, 188
cohens_kappa, 54, 64, 78, 84, 97, 101, 103

create_data_embeddings_description, 43, 118, 131
create_dir, 65, 87
create_object, 52, 65, 67, 81, 89, 95, 118, 176, 186
create_synthetic_units_from_matrix, 52, 66, 68, 81, 89, 95, 118, 176, 186

data.frame, 97, 181
data.frame_to_py_dataset, 63, 67, 82, 122, 124, 126, 175, 177
DataManagerClassifier, 67, 68–71

EmbeddedText, 58, 60, 61, 66, 69, 73, 73, 75, 110, 115, 116, 119, 120, 128, 133, 135, 139, 141, 144–149, 151, 152, 154–157, 161, 162, 166–168, 171–174, 181

factor, 133, 139, 148, 161, 166
feature extractor, 76, 77
fleiss_kappa, 54, 64, 78, 84, 97, 101, 103

generate_args_for_tests, 56, 79, 80, 81, 85, 87, 96, 124
generate_embeddings, 56, 79, 80, 81, 85, 87, 96, 124
generate_id, 52, 66, 67, 80, 89, 95, 118, 176, 186
generate_tensors, 56, 79, 80, 81, 85, 87, 96, 124
get_alpha_3_codes, 82
get_batches_index, 63, 67, 82, 122, 124, 126, 175, 177
get_called_args, 83, 85, 89, 91, 92, 95
get_coder_metrics, 54, 64, 78, 83, 97, 101, 103
get_current_args_for_print, 56, 79–81, 85, 87, 96, 124
get_depr_obj_names, 83, 85, 89, 91, 92, 95
get_desc_for_core_model_architecture, 52, 53, 86, 88, 90
get_dict_cls_type, 52, 53, 86, 88, 90
get_dict_core_models, 52, 53, 86, 88, 90
get_dict_input_types, 52, 53, 86, 88, 90
get_file_extension, 65, 86
get_fixed_test_tensor, 56, 79–81, 85, 87, 96, 124
get_layer_dict, 52, 53, 86, 88, 90

get_layer_documentation, 52, 53, 86, 87,
 90
 get_magnitude_values, 83, 85, 88, 91, 92, 95
 get_n_chunks, 52, 66, 67, 81, 89, 95, 118,
 176, 186
 get_param_def, 83, 85, 89, 90, 92, 95
 get_param_dict, 57, 83, 85, 89, 91, 91, 92, 95
 get_param_doc_desc, 83, 85, 89, 91, 92, 92,
 95
 get_parameter_documentation, 52, 53, 86,
 88, 90
 get_py_package_version, 93, 93, 116, 117,
 128
 get_py_package_versions, 93, 93, 116, 117,
 128
 get_synthetic_cases_from_matrix, 52, 66,
 67, 81, 89, 94, 118, 176, 186
 get_TEClassifiers_class_names, 83, 85,
 89, 91, 92, 95
 get_test_data_for_classifiers, 56,
 79–81, 85, 87, 95, 124
 gwet_ac, 54, 64, 78, 84, 96, 101, 103

 imdb_movie_reviews, 97
 install_aifedducation, 56, 97, 99, 100, 123,
 129, 187
 install_aifedducation_studio, 56, 98, 98,
 100, 123, 129, 187
 install_py_modules, 56, 98, 99, 99, 123,
 129, 187

 kendalls_w, 54, 64, 78, 84, 97, 100, 103
 knnor, 101
 knnor_is_same_class, 102
 kripp_alpha, 54, 64, 78, 84, 97, 101, 102

 LargeDataSetBase, 48, 62, 103, 105, 121,
 156, 159
 LargeDataSetForText, 11, 12, 16, 17, 20, 22,
 25, 27, 30, 32, 35, 37, 40, 41, 78,
 106, 106, 115, 116, 128, 181
 LargeDataSetForTextEmbeddings, 58, 60,
 61, 69, 73, 77, 78, 110, 110, 113,
 115, 116, 119, 120, 128, 133, 135,
 139, 141, 144–149, 151, 152,
 154–157, 161, 162, 166–168,
 171–174, 181
 load_all_py_scripts, 93, 116, 117, 128
 load_from_disk, 116, 129

 load_py_scripts, 93, 116, 117, 128
 long_load_target_data, 43, 117, 131

 matrix_to_array_c, 52, 66, 67, 81, 89, 95,
 118, 176, 186
 ModelsBasedOnTextEmbeddings, 48, 62, 106,
 119, 156, 159

 output_message, 55, 64, 121, 123, 125–128,
 188

 prepare_r_array_for_dataset, 63, 67, 82,
 122, 124, 126, 175, 177
 prepare_session, 56, 98–100, 122, 129, 187
 print_message, 55, 64, 121, 123, 125–128,
 188
 py_dataset_to_embeddings, 63, 67, 82, 122,
 123, 126, 175, 177

 random_bool_on_CI, 56, 79–81, 85, 87, 96,
 124
 read_log, 55, 64, 121, 123, 125, 126–128,
 187, 188
 read_loss_log, 55, 64, 121, 123, 125, 125,
 127, 128, 188
 reduce_to_unique, 63, 67, 82, 122, 124, 126,
 175, 177
 reset_log, 55, 64, 121, 123, 125, 126, 127,
 128, 188
 reset_loss_log, 55, 64, 121, 123, 125–127,
 127, 188
 run_py_file, 93, 116, 117, 128

 save_to_disk, 116, 128
 set_transformers_logger, 56, 98–100, 123,
 129, 187
 start_aifedducation_studio, 130
 summarize_args_for_long_task, 43, 118,
 130
 summarize_tracked_sustainability, 82

 TEClassifierParallel, 131, 143, 147, 150,
 165, 171
 TEClassifierParallelPrototype, 137, 137,
 147, 150, 165, 171
 TEClassifierProtoNet, 66, 73, 80, 116, 128,
 137, 143, 144, 144, 150, 165, 171
 TEClassifierRegular, 66, 73, 80, 116, 128,
 137, 143, 147, 148, 148, 150, 165,
 171

TEClassifiersBasedOnProtoNet, 48, 62, 106, 121, 150, 159
TEClassifiersBasedOnRegular, 48, 62, 106, 121, 156, 156
TEClassifierSequential, 137, 143, 147, 148, 150, 159, 171
TEClassifierSequentialPrototype, 137, 143, 144, 147, 150, 165, 165
TEFeatureExtractor, 48, 60, 73, 77, 110, 113, 114, 116, 128, 135, 141, 145, 149, 162, 168, 171, 171, 172, 185
tensor_list_to_numpy, 63, 67, 82, 122, 124, 126, 175, 177
tensor_to_matrix_c, 52, 66, 67, 81, 89, 95, 118, 176, 186
tensor_to_numpy, 63, 67, 82, 122, 124, 126, 175, 176
TextEmbeddingModel, 58, 60, 73, 75, 77, 80, 110, 113, 116, 120, 128, 129, 133, 139, 144, 146, 148, 154, 161, 166, 167, 171, 175, 177, 177, 179
to_categorical_c, 52, 66, 67, 81, 89, 95, 118, 176, 186

update_aifeducation, 56, 98–100, 123, 129, 186

write_log, 55, 64, 121, 123, 125–128, 187