

---

# **BigARTM Documentation**

***Release 1.0***

**Konstantin Vorontsov**

February 15, 2017



<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Downloads</b>	<b>5</b>
<b>3</b>	<b>Formats</b>	<b>7</b>
<b>4</b>	<b>Installation</b>	<b>9</b>
4.1	Installation for Windows users . . . . .	9
4.2	Installation for Linux and Mac OS-X users . . . . .	10
<b>5</b>	<b>Tutorial references</b>	<b>13</b>
5.1	BigARTM command line utility . . . . .	13
5.2	Running BigARTM from Python API . . . . .	17
5.3	Low-level API in C . . . . .	17
<b>6</b>	<b>Python Interface</b>	<b>25</b>
6.1	ARTM model . . . . .	25
6.2	LDA model . . . . .	30
6.3	hARTM . . . . .	33
6.4	Batches Utils . . . . .	36
6.5	Dictionary . . . . .	37
6.6	Regularizers . . . . .	39
6.7	Scores . . . . .	43
6.8	Score Tracker . . . . .	45
6.9	Master Component . . . . .	48
<b>7</b>	<b>Release Notes</b>	<b>55</b>
7.1	Changes in Python API . . . . .	55
7.2	Changes in Protobuf Messages . . . . .	58
7.3	Changes in BigARTM CLI . . . . .	61
7.4	Changes in c_interface . . . . .	61
7.5	BigARTM v0.7.X Release Notes . . . . .	62
<b>8</b>	<b>BigARTM Developer's Guide</b>	<b>77</b>
8.1	Downloads (Windows) . . . . .	77
8.2	Source code . . . . .	78
8.3	Build C++ code on Windows . . . . .	78
8.4	Python code on Windows . . . . .	79
8.5	Compiling .proto files on Windows . . . . .	80

8.6	Working with iPython notebooks remotely . . . . .	80
8.7	Build C++ code on Linux . . . . .	81
8.8	Code style . . . . .	81
<b>9</b>	<b>Legacy documentation pages</b>	<b>83</b>
9.1	Basic BigARTM tutorial for Linux and Mac OS-X users . . . . .	83
9.2	Basic BigARTM tutorial for Windows users . . . . .	85
9.3	Enabling Basic BigARTM Regularizers . . . . .	87
9.4	BigARTM as a Service . . . . .	89
9.5	BigARTM: The Algorithm Under The Hood . . . . .	90
9.6	Messages . . . . .	91
9.7	C++ interface . . . . .	120
9.8	Windows distribution . . . . .	124
	<b>Python Module Index</b>	<b>127</b>

### Getting help

- Learn more about BigARTM from [IPython Notebooks](#), [NLPub.ru](#), [MachineLearning.ru](#) and several [publications](#).
- Search for information in the archives of the [bigartm-users](#) mailing list, or [post a question](#).
- Report bugs with BigARTM in our [ticket tracker](#).
- Try the [Q&A](#) – it's got answers to many common questions.



---

## Introduction

---

**Warning:** Please note that this is a beta version of the BigARTM library which is still undergoing final testing before its official release. Should you encounter any bugs, lack of functionality or other problems with our library, please let us know immediately. Your help in this regard is greatly appreciated.

This is the documentation for the BigARTM library. BigARTM is a tool to infer [topic models](#), based on a novel technique called [Additive Regularization of Topic Models](#). This technique effectively builds multi-objective models by adding the weighted sums of regularizers to the optimization criterion. BigARTM is known to combine well very different objectives, including sparsing, smoothing, topics decorrelation and many others. Such combinations of regularizers significantly improves several quality measures at once almost without any loss of the perplexity.

**Online.** BigARTM never stores the entire text collection in the main memory. Instead the collection is split into small chunks called ‘batches’, and BigARTM always loads a limited number of batches into memory at any time.

**Parallel.** BigARTM can concurrently process several batches, and by doing so it substantially improves the throughput on multi-core machines. The library hosts all computation in several threads withing a single process, which enables efficient usage of shared memory across application threads.

**Extensible API.** BigARTM comes with an API in Python, but can be easily extended for all other languages that have an implementation of [Google Protocol Buffers](#).

**Cross-platform.** BigARTM is known to be compatible with gcc, clang and the Microsoft compiler (VS 2012). We have tested our library on Windows, Ubuntu and Fedora.

**Open source.** BigARTM is released under the [New BSD License](#). If you plan to use our library commercially, please beware that BigARTM depends on ZeroMQ. Please, make sure to review [ZeroMQ license](#).

---

**Acknowledgements.** BigARTM project is supported by Russian Foundation for Basic Research (grants 14-07-00847, 14-07-00908, 14-07-31176), Skolkovo Institute of Science and Technology (project 081-R), Moscow Institute of Physics and Technology.



**Partners**





---

## Downloads

---

- **Windows**

- Latest 64 bit release: [BigARTM\\_v0.8.2\\_win64](#)
- Latest build from master branch: [BigARTM\\_master\\_win64.7z](#) (warning, use this with caution)
- All previous releases are available at <https://github.com/bigartm/bigartm/releases>

Please refer to [Basic BigARTM tutorial for Windows users](#) for step by step installation procedure.

- **Linux, Mac OS-X**

To run BigARTM on Linux and Mac OS-X you need to clone BigARTM repository (<https://github.com/bigartm/bigartm>) and build it as described in [Basic BigARTM tutorial for Linux and Mac OS-X users](#).

- **Datasets**

Download one of the following datasets to start experimenting with BigARTM. See [Formats](#) page for the description of input data formats. Note that `docword.*` and `vocab.*` files indicate UCI BOW format, while `vw.*` file indicate Vowpal Wabbit format.

Task	Source	#Words	#Items	Files
kos	UCI	6906	3430	<ul style="list-style-type: none"> <li>- docword.kos.txt.gz (1 MB)</li> <li>- vocab.kos.txt (54 KB)</li> </ul>
nips	UCI	12419	1500	<ul style="list-style-type: none"> <li>- docword.nips.txt.gz (2.1 MB)</li> <li>- vocab.nips.txt (98 KB)</li> </ul>
enron	UCI	28102	39861	<ul style="list-style-type: none"> <li>- docword.enron.txt.gz (11.7 MB)</li> <li>- vocab.enron.txt (230 KB)</li> </ul>
nytimes	UCI	102660	300000	<ul style="list-style-type: none"> <li>- docword.nytimes.txt.gz (223 MB)</li> <li>- vocab.nytimes.txt (1.2 MB)</li> </ul>
pubmed	UCI	141043	8200000	<ul style="list-style-type: none"> <li>- docword.pubmed.txt.gz (1.7 GB)</li> <li>- vocab.pubmed.txt (1.3 MB)</li> </ul>
wiki	Gensim	100000	3665223	<ul style="list-style-type: none"> <li>- vw.wiki-en.txt.zip (1.8 GB)</li> </ul>
wiki_enru	Wiki	196749	216175	<ul style="list-style-type: none"> <li>- vw.wiki_enru.txt.zip (285 MB)</li> </ul>
eurlex	eurlex	19800	21000	<ul style="list-style-type: none"> <li>- vw.eurlex.txt.zip (13 MB)</li> <li>- vw.eurlex-test.txt.zip (13 MB)</li> </ul>
lastfm	lastfm		1k, 360k	<ul style="list-style-type: none"> <li>- vw.lastfm_1k.txt.zip (100 MB)</li> <li>- vw.lastfm_360k.txt.zip (330 MB)</li> </ul>
6				<b>Chapter 2. Downloads</b>
	mmro	7805	1061	<ul style="list-style-type: none"> <li>-</li> </ul>

## Formats

This page describes input data formats compatible with BigARTM. Currently all formats correspond to [Bag-of-words representation](#), meaning that all linguistic processing (lemmatization, tokenization, detection of n-grams, etc) needs to be done outside BigARTM.

1. [Vowpal Wabbit](#) is a single-format file, based on the following principles:

- each document is represented in a single line
- all tokens are represented as strings (no need to convert them into an integer identifier)
- token frequency defaults to 1.0, and can be optionally specified after a colon (:)
- namespaces (*Batch.class\_id*) can be identified by a pipe (|)

### Example 1

```
doc1 Alpha Bravo:10 Charlie:5 |author Ola_Nordmann
doc2 Bravo:5 Delta Echo:3 |author Ivan_Ivanov
```

### Example 2

```
user123 |track-like track2 track5 track7 |track-play track1:10 track2:25 track3:2 track7:8 |track-stop track1:10 track2:25 track3:2 track7:8 |track-stop
user345 |track-like track2 track5 track7 |track-play track1:10 track2:25 track3:2 track7:8 |track-stop
```

2. [UCI Bag-of-words](#) format consists of two files - `vocab.*.txt` and `docword.*.txt`. The format of the `docword.*.txt` file is 3 header lines, followed by NNZ triples:

```
D
W
NNZ
docID wordID count
docID wordID count
...
docID wordID count
```

The file must be sorted on docID. Values of wordID must be unity-based (not zero-based). The format of the `vocab.*.txt` file is line containing `wordID=n`. Note that words must not have spaces or tabs. In `vocab.*.txt` file it is also possible to specify the namespace (*Batch.class\_id*) for tokens, as it is shown in this example:

```
token1 @default_class
token2 custom_class
token3 @default_class
token4
```

Use space or tab to separate token from its class. Token that are not followed by class label automatically get “@default\_class” as a label (see “token4” in the example).

**Unicode support.** For non-ASCII characters save `vocab.*.txt` file in **UTF-8** format.

### 3. Batches (binary BigARTM-specific format).

This is compact and efficient format, based on several protobuf messages in public BigARTM interface (*Batch*, *Item* and *Field*).

- A batch is a collection of several items
- An item is a collection of several fields
- A field is a collection of pairs (`token_id`, `token_weight`).

The following example shows a Python code that generates a synthetic batch.

```
import artm.messages, random, uuid

num_tokens = 60
num_items = 100
batch = artm.messages.Batch()
batch.id = str(uuid.uuid4())
for token_id in range(0, num_tokens):
    batch.token.append('token' + str(token_id))

for item_id in range(0, num_items):
    item = batch.item.add()
    item.id = item_id
    field = item.field.add()
    for token_id in range(0, num_tokens):
        field.token_id.append(token_id)
        background_count = random.randint(1, 5) if (token_id >= 40) else 0
        topical_count = 10 if (token_id < 40) and ((token_id % 10) == (item_id % 10)) else 0
        field.token_weight.append(background_count + topical_count)
```

Note that the batch has its local dictionary, `batch.token`. This dictionary which maps `token_id` into the actual token. In order to create a batch from textual files involve one needs to find all distinct words, and map them into sequential indices.

`batch.id` must be set to a unique GUID in a format of 00000000-0000-0000-0000-000000000000.

---

## Installation

---

### Installation for Windows users

#### Download

Download latest binary distribution of BigARTM from <https://github.com/bigartm/bigartm/releases>. Explicit download links can be found at [Downloads](#) section (for 32 bit and 64 bit configurations).

The distribution will contain pre-build binaries, command-line interface and BigARTM API for Python. The distribution also contains a simple dataset. More datasets in BigARTM-compatible format are available in the [Downloads](#) section.

Refer to [Windows distribution](#) for details about other files, included in the binary distribution package.

#### Configure BigARTM Python API

1. Install Python, for example from the following links:

- Python 2.7.11, 64 bit – <https://www.python.org/ftp/python/2.7.11/python-2.7.11.amd64.msi>, or
- Python 2.7.11, 32 bit – <https://www.python.org/ftp/python/2.7.11/python-2.7.11.msi>

Remember that the version of BigARTM package must match your version Python installed on your machine. If you have 32 bit operating system then you must select 32 bit for Python and BigARTM package. If you have 64 bit operating system then you are free to select either version. However, please note that memory usage of 32 bit processes is limited by 2 GB. For this reason we recommend to select 64 bit configurations.

Please note that you must use Python 2.7, because Python 3 is not supported by BigARTM.

Also you need to have several Python libraries to be installed on your machine:

- numpy >= 1.9.2
- pandas >= 0.16.2

2. Add C:\BigARTM\bin folder to your PATH system variable, and add C:\BigARTM\python to your PYTHONPATH system variable:

```
set PATH=%PATH%;C:\BigARTM\bin
set PATH=%PATH%;C:\Python27;C:\Python27\Scripts
set PYTHONPATH=%PYTHONPATH%;C:\BigARTM\Python
```

Remember to change C:\BigARTM and C:\Python27 with your local folders.

3. Setup *Google Protocol Buffers* library, included in the BigARTM release package.

- Copy `C:\BigARTM\bin\protoc.exe` file into `C:\BigARTM\protobuf\src` folder
- Run the following commands from command prompt

```
cd C:\BigARTM\protobuf\Python
python setup.py build
python setup.py install
```

Avoid `python setup.py test` step, as it produces several confusing errors. Those errors are harmless. For further details about protobuf installation refer to [protobuf/python/README](#).

## Installation for Linux and Mac OS-X users

Currently there is no distribution package of BigARTM for Linux. BigARTM had been tested on several Linux distributions, and it is known to work well, but you have to get the source code and compile it locally on your machine.

### System dependencies

Building BigARTM requires the following components:

- [git](#) (any recent version) – for obtaining source code;
- [cmake](#) (at least of version 2.8), *make*, *g++* or *clang* compiler with *c++11* support, *boost* (at least of version 1.40) – for building library and binary executable;
- [python](#) (version 2.7) – for building Python API for BigARTM.

To simplify things, you may type:

- **On deb-based distributions:** `sudo apt-get install git make cmake build-essential libboost-all-dev`
- **On rpm-based distributions:** `sudo yum install git make cmake gcc-c++ glibc-static libstdc++-static boost boost-static python` (for Fedora 22 or higher use `dnf` instead of `yum`)
- **On Mac OS distributions:** `brew install git cmake boost`

### Download sources and build

Clone the latest BigARTM code from our github repository, and build it via CMake as in the following script.

```
cd ~
git clone --branch=stable https://github.com/bigartm/bigartm.git
cd bigartm
mkdir build && cd build
cmake ..
make
```

**Note for Linux users:** By default building binary executable `bigartm` requires static versions of Boost, C and C++ libraries. To alter it, run `cmake` command with option `-DBUILD_BIGARTM_CLI_STATIC=OFF`.

## System-wide installation

To install command-line utility, shared library module and Python interface for BigARTM, you can type:

```
sudo make install
```

Normally this will install:

- bigartm utility into folder `/usr/local/bin/`;
- shared library `libartm.so` (`artm.dylib` for Mac OS-X) into folder `/usr/local/lib/`;
- Python interface for BigARTM into Python-specific system directories, along with necessary dependencies.

If you want to alter target folders for binary and shared library objects, you may specify common prefix while running `cmake` command via option `-DCMAKE_INSTALL_PREFIX=path_to_folder`. By default `CMAKE_INSTALL_PREFIX=/usr/local/`.

## Configure BigARTM Python API

If you want to use only Python interface for BigARTM, you may run following commands:

```
# Step 1 - install Google Protobuf as dependency
cd ~/bigartm/3rdparty/protobuf/python
sudo python setup.py install

# Step 2 - install Python interface for BigARTM
cd ~/bigartm/python
sudo python setup.py install

# Step 3 - point ARTM_SHARED_LIBRARY variable to libartm.so (libartm.dylib) location
export ARTM_SHARED_LIBRARY=~/.bigartm/build/lib/libartm.so      # for linux
export ARTM_SHARED_LIBRARY=~/.bigartm/build/lib/libartm.dylib  # for Mac OS X
```

We strongly recommend system-wide installation as there is no need to keep BigARTM code after it, so you may safely remove folder `~/bigartm/`.

## Troubleshooting

If you build BigARTM in existing folder `build` (e.g. you built BigARTM before) and encounter any errors, it may be due to out-of-date file `CMakeCache.txt` in folder `build`. In that case we strongly recommend to delete this file and try to build again.

Using BigARTM Python API you can encounter this error:

```
Traceback (most recent call last):
File "<stdin>", line 1, in <module>
File "build/bdist.linux-x86_64/egg/artm/wrapper/api.py", line 19, in __init__
File "build/bdist.linux-x86_64/egg/artm/wrapper/api.py", line 53, in _load_cdll
OSError: libartm.so: cannot open shared object file: No such file or directory
Failed to load artm shared library. Try to add the location of `libartm.so` file into your LD_LIBRARY_PATH
```

This error indicates that BigARTM's python interface can not locate `libartm.so` (`libartm.dylib`) files. In such case type `export ARTM_SHARED_LIBRARY=path_to_artm_shared_library`.

## BigARTM on Travis-CI

To get a live usage example of BigARTM you may check BigARTM's [.travis.yml](#) script and the latest [continuous integration build](#).



---

## Tutorial references

---

### BigARTM command line utility

This document provides an overview of `bigartm` command-line utility shipped with BigARTM.

For a detailed description of `bigartm` command line interface refer to [bigartm.exe notebook](#) (in Russian).

In brief, you need to download some input data (a textual collection represented in bag-of-words format). We recommend to download sample collections in **vowpal wabbit** format by links provided in [Downloads](#) section of the tutorial. Then you can use `bigartm` as described by `bigartm --help`. You may also get more information about builtin regularizers by typing `bigartm --help --regularizer`.

```
BigARTM v0.8.2 - library for advanced topic modeling (http://bigartm.org):

Input data:
  -c [ --read-vw-corpus ] arg      Raw corpus in Vowpal Wabbit format
  -d [ --read-uci-docword ] arg    docword file in UCI format
  -v [ --read-uci-vocab ] arg      vocab file in UCI format
  --read-cooc arg                  read co-occurrences format
  --batch-size arg (=500)          number of items per batch
  --use-batches arg                folder with batches to use

Dictionary:
  --dictionary-min-df arg          filter out tokens present in less than
                                   N documents / less than P% of documents
  --dictionary-max-df arg          filter out tokens present in less than
                                   N documents / less than P% of documents
  --dictionary-size arg (=0)       limit dictionary size by filtering out
                                   tokens with high document frequency
  --use-dictionary arg             filename of binary dictionary file to
                                   use

Model:
  --load-model arg                 load model from file before processing
  -t [ --topics ] arg (=16)       number of topics
  --use-modality arg               modalities (class_ids) and their
                                   weights
  --predict-class arg              target modality to predict by theta
                                   matrix

Learning:
  -p [ --num-collection-passes ] arg (=0)
                                   number of outer iterations (passes)
```

<code>--num-document-passes arg (=10)</code>	through the collection) number of inner iterations (passes through the document)
<code>--update-every arg (=0)</code>	[online algorithm] requests an update of the model after <code>update_every</code> document
<code>--tau0 arg (=1024)</code>	[online algorithm] weight option from online update formula
<code>--kappa arg (=0.6999999988)</code>	[online algorithm] exponent option from online update formula
<code>--reuse-theta</code>	reuse theta between iterations
<code>--regularizer arg</code>	regularizers (SmoothPhi, SparsePhi, SmoothTheta, SparseTheta, Decorrelation)
<code>--threads arg (=1)</code>	number of concurrent processors (default: auto-detect)
<code>--async</code>	invoke asynchronous version of the online algorithm
Output:	
<code>--save-model arg</code>	save the model to binary file after processing
<code>--save-batches arg</code>	batch folder
<code>--save-dictionary arg</code>	filename of dictionary file
<code>--write-model-readable arg</code>	output the model in a human-readable format
<code>--write-dictionary-readable arg</code>	output the dictionary in a human-readable format
<code>--write-predictions arg</code>	write prediction in a human-readable format
<code>--write-class-predictions arg</code>	write class prediction in a human-readable format
<code>--write-scores arg</code>	write scores in a human-readable format
<code>--write-vw-corpus arg</code>	convert batches into plain text file in Vowpal Wabbit format
<code>--force</code>	force overwrite existing output files
<code>--csv-separator arg (=;)</code>	columns separator <b>for</b>
	<code>--write-model-readable</code> and <code>--write-predictions</code> . Use <code>\t</code> or TAB to indicate tab.
<code>--score-level arg (=2)</code>	score level (0, 1, 2, or 3
<code>--score arg</code>	scores (Perplexity, SparsityTheta, SparsityPhi, TopTokens, ThetaSnippet, or TopicKernel)
<code>--final-score arg</code>	final scores (same as scores)
Other options:	
<code>-h [ --help ]</code>	display this <b>help</b> message
<code>--rand-seed arg</code>	specify seed <b>for</b> random number generator, use system timer when not specified
<code>--guid-batch-name</code>	applies to <code>save-batches</code> and indicate that batch names should be <b>guids</b> (not sequential codes)
<code>--response-file arg</code>	response file
<code>--paused</code>	start paused and waits <b>for</b> a keystroke (allows to attach a debugger)
<code>--disk-cache-folder arg</code>	disk cache folder
<code>--disable-avx-opt</code>	disable AVX optimization (gives similar

```

--time-limit arg (=0)      behavior of the Processor component to
--log-dir arg              BigARTM v0.5.4)
                           limit execution time in milliseconds
                           target directory for logging
                           (GLOG_log_dir)
--log-level arg            min logging level (GLOG_minloglevel;
                           INFO=0, WARNING=1, ERROR=2, and
                           FATAL=3)

```

#### Examples:

```

* Download input data:
  wget https://s3-eu-west-1.amazonaws.com/artm/docword.kos.txt
  wget https://s3-eu-west-1.amazonaws.com/artm/vocab.kos.txt
  wget https://s3-eu-west-1.amazonaws.com/artm/vw.mmro.txt
  wget https://s3-eu-west-1.amazonaws.com/artm/vw.wiki-enru.txt.zip

* Parse docword and vocab files from UCI bag-of-words format; then fit topic model with 20 topics:
  bigartm -d docword.kos.txt -v vocab.kos.txt -t 20 --num_collection_passes 10

* Parse VW format; then save the resulting batches and dictionary:
  bigartm --read-vw-corpus vw.mmro.txt --save-batches mmro_batches --save-dictionary mmro.dict

* Parse VW format from standard input; note usage of single dash '-' after --read-vw-corpus:
  cat vw.mmro.txt | bigartm --read-vw-corpus - --save-batches mmro2_batches --save-dictionary mmro2.dict

* Re-save batches back into VW format:
  bigartm --use-batches mmro_batches --write-vw-corpus vw.mmro.txt

* Parse only specific modalities from VW file, and save them as a new VW file:
  bigartm --read-vw-corpus vw.wiki-enru.txt --use-modality @russian --write-vw-corpus vw.wiki-ru.txt

* Load and filter the dictionary on document frequency; save the result into a new file:
  bigartm --use-dictionary mmro.dict --dictionary-min-df 5 dictionary-max-df 40% --save-dictionary mmro.dict

* Load the dictionary and export it in a human-readable format:
  bigartm --use-dictionary mmro.dict --write-dictionary-readable mmro.dict.txt

* Use batches to fit a model with 20 topics; then save the model in a binary format:
  bigartm --use-batches mmro_batches --num_collection_passes 10 -t 20 --save-model mmro.model

* Load the model and export it in a human-readable format:
  bigartm --load-model mmro.model --write-model-readable mmro.model.txt

* Load the model and use it to generate predictions:
  bigartm --read-vw-corpus vw.mmro.txt --load-model mmro.model --write-predictions mmro.predict.txt

* Fit model with two modalities (@default_class and @target), and use it to predict @target label:
  bigartm --use-batches <batches> --use-modality @default_class,@target --topics 50 --num_collection_passes 10
  bigartm --use-batches <batches> --use-modality @default_class,@target --topics 50 --load-model mmro.model
  --write-predictions pred.txt --csv-separator=tab
  --predict-class @target --write-class-predictions pred_class.txt --score ClassPrecision

* Fit simple regularized model (increase sparsity up to 60-70%):
  bigartm -d docword.kos.txt -v vocab.kos.txt --dictionary-max-df 50% --dictionary-min-df 2
  --num_collection_passes 10 --batch-size 50 --topics 20 --write-model-readable model.txt
  --regularizer "0.05 SparsePhi" "0.05 SparseTheta"

```

```

* Fit more advanced regularize model, with 10 sparse objective topics, and 2 smooth background topics
  bigartm -d docword.kos.txt -v vocab.kos.txt --dictionary-max-df 50% --dictionary-min-df 2
    --num_collection_passes 10 --batch-size 50 --topics obj:10;background:2 --write-model-read
    --regularizer "0.05 SparsePhi #obj"
    --regularizer "0.05 SparseTheta #obj"
    --regularizer "0.25 SmoothPhi #background"
    --regularizer "0.25 SmoothTheta #background"

* Upgrade batches in the old format (from folder 'old_folder' into 'new_folder'):
  bigartm --use-batches old_folder --save-batches new_folder

* Configure logger to output into stderr:
  tset GLOG_logtostderr=1 & bigartm -d docword.kos.txt -v vocab.kos.txt -t 20 --num_collection_passes

```

#### Additional information about regularizers:

```

>bigartm.exe --regularizer --help
List of regularizers available in BigARTM CLI:

--regularizer "tau SmoothTheta #topics"
--regularizer "tau SparseTheta #topics"
--regularizer "tau SmoothPhi #topics @class_ids !dictionary"
--regularizer "tau SparsePhi #topics @class_ids !dictionary"
--regularizer "tau Decorrelation #topics @class_ids"
--regularizer "tau TopicSelection #topics"
--regularizer "tau LabelRegularization #topics @class_ids !dictionary"
--regularizer "tau ImproveCoherence #topics @class_ids !dictionary"
--regularizer "tau Biterms #topics @class_ids !dictionary"

```

List of regularizers available in BigARTM, but not exposed in CLI:

```

--regularizer "tau SpecifiedSparsePhi"
--regularizer "tau SmoothPtdw"
--regularizer "tau HierarchySparsingTheta"

```

If you are interested to see any of these regularizers in BigARTM CLI please send a message to [bigartm-users@googlegroups.com](mailto:bigartm-users@googlegroups.com).

By default all regularizers act on the full **set** of topics and modalities. To limit action onto specific **set** of topics use **hash** sign (**#**), *followed by* list of topics (**for** example, **#topic1;topic2**) or *topic groups* (**#obj**). Similarly, to limit action onto specific **set** of class ids use **at** sign (**@**), by the list of class ids (**for** example, **@default\_class**). Some regularizers accept a dictionary. To specify the dictionary use exclamation mark (**!**), followed by the path to the dictionary(.dict file in your file system). Depending on regularizer the dictionary can be either optional or required. Some regularizers expect an dictionary with tokens and their frequencies; Other regularizers expect an dictionary with tokens co-occurencies; For more information about regularizers refer to wiki-page:

<https://github.com/bigartm/bigartm/wiki/Implemented-regularizers>

To get full **help** run ``bigartm --help`` without `--regularizer` switch.

## Running BigARTM from Python API

Refer to ARTM tutorial ([in Russian](#) or [in English](#)), which describes artm.ARTM model from high-level Python API of BigARTM.

Refer to LDA tutorial ([in Russian](#) or [in English](#)), which describes artm.LDA model from high-level Python API of BigARTM.

Refer to ARTM notebook with model experiment ([in Russian](#) or [in English](#)), which shows an example of usage of artm.ARTM model from high-level Python API of BigARTM.

If some of these link are not available, try to open the repository manually: <https://github.com/bigartm/bigartm-book>

## Low-level API in C

This document explains all public methods of the low level BigARTM interface, written in plain C language.

### Introduction

The goal of low level API is to expose all functionality of the library in a set of simple functions written in plain C language. This makes it easier to consume BigARTM from various programming environments. For example, the [Python Interface](#) of BigARTM uses `ctypes` module to call the low level BigARTM interface. Most programming environments also have similar functionality: [PInvoke](#) in C#, [loadlibrary](#) in Matlab, etc.

Typical methods of low-level API may look as follows:

```
int ArtmCreateMasterModel(int length, const char* master_model_config);
int ArtmFitOfflineMasterModel(int master_id, int length, const char* fit_offline_master_model_args);
int ArtmRequestTopicModel(int master_id, int length, const char* get_model_args);
```

This methods, similarly to most other methods in low level API, accept a serialized binary representation of some Google Protocol Buffer message. From BigARTM v0.8.2 it is also possible to pass JSON-serialized protobuf message. This might be useful if you are planing to use low-level C interface from environment where configuring protobuf libraries would be challenging. Please, refer to [Messages](#) for more details about each particular message, and [Protobuf documentation](#) regarding JSON mapping.

Note that this documentation is incomplete. For the actual list the methods of low-level C API please refer to [c\\_interface.h](#). Same is true about messages documentation. It is always recommended to review the [messages.proto](#) definition.

If you plan to implement a high-wrapper around low-level API we recoomend to review the source code of existing wrappers [cpp\\_interface.h](#), [cpp\\_interface.cc](#) (for C++ wrapper) and [spec.py](#), [api.py](#) (for python wrapper).

### List of all methods with corresponding protobuf types

	ArtmConfigureLogging	(	artm.ConfigureLoggingArgs);
const char*	= ArtmGetVersion();		
const char*	= ArtmGetLastErrorMessage();		
artm.CollectionParserInfo	= ArtmParseCollection	(	artm.CollectionParserConfig);
master_id	= ArtmCreateMasterModel	(	artm.MasterModelConfig);
	ArtmReconfigureMasterModel	(master_id, artm.MasterModelConfig);	
	ArtmReconfigureTopicName	(master_id, artm.MasterModelConfig);	

	ArtmDisposeMasterComponent	(master_id);
	ArtmImportBatches	(master_id, artm.ImportBatchesArgs);
	ArtmGatherDictionary	(master_id, artm.GatherDictionaryArgs);
	ArtmFilterDictionary	(master_id, artm.FilterDictionaryArgs);
	ArtmCreateDictionary	(master_id, artm.DictionaryData);
	ArtmImportDictionary	(master_id, artm.ImportDictionaryArgs);
	ArtmExportDictionary	(master_id, artm.ExportDictionaryArgs);
artm.DictionaryData	= ArtmRequestDictionary	(master_id, artm.GetDictionaryArgs);
	ArtmInitializeModel	(master_id, artm.InitializeModelArgs);
	ArtmExportModel	(master_id, artm.ExportModel);
	ArtmImportModel	(master_id, artm.ImportModel);
	ArtmOverwriteTopicModel	(master_id, artm.TopicModel);
	ArtmFitOfflineMasterModel	(master_id, artm.FitOfflineMasterModelArgs);
	ArtmFitOnlineMasterModel	(master_id, artm.FitOnlineMasterModelArgs);
artm.ThetaMatrix	= ArtmRequestTransformMasterModel	(master_id, artm.TransformMasterModelArgs);
artm.ThetaMatrix	= ArtmRequestTransformMasterModelExternal	(master_id, artm.TransformMasterModelArgs);
artm.MasterModelConfig	= ArtmRequestMasterModelConfig	(master_id);
artm.ThetaMatrix	= ArtmRequestThetaMatrix	(master_id, artm.GetThetaMatrix);
artm.ThetaMatrix	= ArtmRequestThetaMatrixExternal	(master_id, artm.GetThetaMatrix);
artm.TopicModel	= ArtmRequestTopicModel	(master_id, artm.GetTopicModel);
artm.TopicModel	= ArtmRequestTopicModelExternal	(master_id, artm.GetTopicModel);
artm.ScoreData	= ArtmRequestScore	(master_id, artm.GetScoreValueArgs);
artm.ScoreArray	= ArtmRequestScoreArray	(master_id, artm.GetScoreArrayArgs);
artm.MasterComponentInfo	= ArtmRequestMasterComponentInfo	(master_id, artm.GetMasterComponentInfoArgs);
	ArtmDisposeModel	(master_id, const char* model_name);
	ArtmDisposeDictionary	(master_id, const char* dictionary_name);
	ArtmDisposeBatch	(master_id, const char* batch_name);
	ArtmClearThetaCache	(master_id, artm.ClearThetaCacheArgs);
	ArtmClearScoreCache	(master_id, artm.ClearScoreCacheArgs);
	ArtmClearScoreArrayCache	(master_id, artm.ClearScoreArrayCacheArgs);
	ArtmCopyRequestedMessage	(int length, char* address);
	ArtmCopyRequestedObject	(int length, char* address);
	ArtmSetProtobufMessageFormatToJson();	
	ArtmSetProtobufMessageFormatToBinary();	
int	= ArtmProtobufMessageFormatIsJson();	

Below we give a short description of these methods.

- `ArtmConfigureLogging` allows to configure logging parameters; this method is optional, you may not use it
- `ArtmGetVersion` returns the version of BigARTM library
- `ArtmParseCollection` parse collection in VW or UCI-BOW formats, creates batches and stores them to disk
- `ArtmCreateMasterModel` / `ArtmReconfigureMasterModel` / `ArtmDisposeMasterComponent` create master model / updates its parameters / dispose given instance of master model.
- `ArtmImportBatches` loads batches from disk into memory for quicker processing. This is optional, most methods that require batches can work directly with files on disk.

- `ArtmGatherDictionary` / `ArtmFilterDictionary` / `ArtmImportDictionary` / `ArtmExportDictionary` Main methods to work with dictionaries. *Gather* initialized the dictionary based on a folder with batches, *Filter* eliminates tokens based on their frequency, *Import/Export* save and re-load dictionary to/from disk.
- You may also create the dictionary from `artm.DictionaryData` message, that contains the list of all tokens to be included in the dictionary. To do this use method `ArtmCreateDictionary` (to create a dictionary) and `ArtmRequestDictionary` (to retrieve `artm.DictionaryData` for an existing dictionary).
- `ArtmInitializeModel` / `ArtmExportModel` / `ArtmImportModel` handle *models* (e.g. matrices of size  $|T| \times |W|$  such as *pwt*, *nwt* or *rwt*). *Initialize\** fills the matrix with random 0..1 values. *Export* and *Import* saves the matrix to disk and re-loads it back.
- `ArtmOverwriteTopicModel` allows to overwrite values in topic model (for example to manually specify initial approximation).
- `ArtmFitOfflineMasterModel` — fit the model with *offline* algorithm
- `ArtmFitOnlineMasterModel` — fit the model with *online* algorithm
- `ArtmRequestTransformMasterModel` — apply the model to new data
- `ArtmRequestMasterModelConfig` — retrieve configuration of master model
- `ArtmRequestThetaMatrix` — retrieve cached theta matrix
- `ArtmRequestTopicModel` — retrieve a model (e.g. *pwt*, *nwt* or *rwt* matrix)
- `ArtmRequestScore` — retrieve score (such as perplexity, sparsity, etc)
- `ArtmRequestScoreArray` — retrieve historical information for a given score
- `ArtmRequestMasterComponentInfo` — retrieve diagnostics information and internal state of the master model
- `ArtmDisposeModel` / `ArtmDisposeDictionary` / `ArtmDisposeBatch` — dispose specific objects
- `ArtmClearThetaCache` / `ArtmClearScoreCache` / `ArtmClearScoreArrayCache` — clear specific caches
- `ArtmSetProtobufMessageFormatToJson` / `ArtmSetProtobufMessageFormatToBinary` / `ArtmProtobufMessageFormatIsJson` — configure the low-level API to work with JSON-serialized protobuf messages instead of binary-serialized protobuf messages

The following operations are less important part of low-level BigARTM CLI. In most cases you won't need them, unless you have a very specific needs.

<code>master_id</code>	<code>= ArtmDuplicateMasterComponent</code>	<code>(master_id, artm.DuplicateMasterComponent</code>
	<code>ArtmCreateRegularizer</code>	<code>(master_id, artm.RegularizerConfig);</code>
	<code>ArtmReconfigureRegularizer</code>	<code>(master_id, artm.RegularizerConfig);</code>
	<code>ArtmDisposeRegularizer</code>	<code>(master_id, const char* regularizer_name);</code>
	<code>ArtmOverwriteTopicModelNamed</code>	<code>(master_id, artm.TopicModel, const char* name);</code>
	<code>ArtmCreateDictionaryNamed</code>	<code>(master_id, artm.DictionaryData, const char* name);</code>
	<code>ArtmAttachModel</code>	<code>(master_id, artm.AttachModelArgs, int argc, const char* argv[]);</code>
<code>artm.ProcessBatchesResult</code>	<code>= ArtmRequestProcessBatches</code>	<code>(master_id, artm.ProcessBatchesArgs);</code>
<code>artm.ProcessBatchesResult</code>	<code>= ArtmRequestProcessBatchesExternal</code>	<code>(master_id, artm.ProcessBatchesArgs);</code>
	<code>ArtmAsyncProcessBatches</code>	<code>(master_id, artm.ProcessBatchesArgs);</code>
	<code>ArtmMergeModel</code>	<code>(master_id, artm.MergeModelArgs);</code>
	<code>ArtmRegularizeModel</code>	<code>(master_id, artm.RegularizeModelArgs);</code>
	<code>ArtmNormalizeModel</code>	<code>(master_id, artm.NormalizeModelArgs);</code>
<code>artm.Batch</code>	<code>= ArtmRequestLoadBatch</code>	<code>(const char* filename);</code>
	<code>ArtmAwaitOperation</code>	<code>(int operation_id, artm.AwaitOperationArgs);</code>
	<code>ArtmSaveBatch</code>	<code>(const char* disk_path, artm.Batch);</code>

## Protocol for retrieving results

The methods in low-level API can be split into two groups — those that *execute* certain action, and those that *request* certain data. For example `ArtmCreateMasterModel` and `ArtmFitOfflineMasterModel` just execute an action, while `ArtmRequestTopicModel` is a request for data. Naming convention is that such requests always start with `ArtmRequest` prefix.

1. To call execute-action method is fairly straightforward — first you create a protobuf message that describe the arguments of the operation. For example, `ArtmCreateMasterModel` expects `artm.MasterModelConfig` message, as defined in the documentation of `ArtmCreateMasterModel`. Then you serialize protobuf message, and pass it to the method along with the length of the serialized message. In some cases you also pass the *master\_id*, returned by `ArtmCreateMasterModel`, as described in details further below on this page. The execute-action method will typically return an error code, with zero value (or `ARTM_SUCCESS`) indicating successful execution.
2. To call request-data method is more tricky. First you follow the same procedure as when calling an execute-action method, e.g. create and serialize protobuf message and pass it to your `ArtmRequestXxx` operation. For example, `ArtmRequestTopicModel` expects `artm.GetTopicModelArgs` message. Then the method like `ArtmRequestTopicModel` will return the size (in bytes) of the memory buffer that needs to be allocated by caller. To fill this buffer with actual data you need to call method

```
int ArtmCopyRequestedMessage(int length, char* address)
```

where `address` give a pointer to the memory buffer, and `length` must give the length of the buffer (e.g. must match the value returned by `ArtmRequestXxx` call). After `ArtmCopyRequestedMessage` the buffer will contain protobuf-serialized message. To deserialize this message you need to know its protobuf type, which will be defined by the documentation of the `ArtmRequestXxx` method that you are calling. For `ArtmRequestTopicModel` it will be a `artm.TopicModel` message.

3. Note that few `ArtmRequestXxx` methods has a more complex protocol that require two subsequent calls — first, to `ArtmCopyRequestedMessage`, and then to `ArtmCopyRequestedObject`. If that's the case the name of the method will be `ArtmRequestXxxExternal` (for example `ArtmRequestThetaMatrixExternal` or `ArtmRequestTopicModelExternal`). Typically this is used to copy out large objects, such as theta or phi matrices, and store them directly as dense matrices, bypassing protobuf serialization. For more information see [cpp\\_interface.cc](#).

A side-note on thread safety: in between calls to `ArtmRequestXxx` and `ArtmCopyRequestedMessage` the result is stored in a thread local storage. This allows you to call multiple `ArtmRequestXxx` methods from different threads.

## Error handling

All methods in this API return an integer value. Negative return values represent an error code. See [error codes](#) for the list of all error codes. To get corresponding error message as string use `ArtmGetLastErrorMessage()`. Non-negative return values represent a success, and for some API methods might also incorporate some useful information. For example, `ArtmCreateMasterModel()` returns the ID of newly created master component, and `ArtmRequestTopicModel()` returns the length of the buffer that should be allocated before calling `ArtmCopyRequestedMessage()`.

## MasterId and MasterModel

The concept of *Master Model* is central in low-level API. Almost any interaction with the low-level API starts by calling method `ArtmCreateMasterModel`, which creates an instance of so-called *Master Model* (or *Master Component*), and returns its *master\_id* – an integer identifier that refers to that instance. You need *master\_id* in the remaining methods of the low-level API, such as `ArtmFitOfflineMasterModel`. *master\_id* creates a context, or scope, that isolate different models from each other. An operation applied to a specific *master\_id* will



not affect other master components. Each master model occupy some memory — potentially a very large amount, depending on the number of topics and tokens in the model. Once you are done with a specific instance of master component you need to dispose its resources by calling `ArtmDisposeMasterComponent(master_id)`. After that `master_id` is no longer valid, and it must not be used as argument to other methods.

You may use method `ArtmRequestMasterComponentInfo` to retrieve internal diagnostics information about master component. It will reveal its internal state and tell the config of the master component, the list of scores and regularizers, the list of phi matrices, the list of dictionaries, cache entries, and other informatino that will help to understand how master component is functioning.

Note there might be confusion between terms *MasterComponent* and *MasterModel*, throughout this page as well as in the actual naming of the methods. This is due to historical reasons, and for all practical purposes you may think that this terms refer to the same thing.

## ArtmConfigureLogging

You may use `ArtmConfigureLogging` call to set logging parameters, such as verbosity level or directory to output logs. You are not require to call `ArtmConfigureLogging`, in which case logging is automatically initialized to INFO level, and logs are placed in the active working folder.

Note that you can set log directory just one time. Once it is set you can not change it afterwards. Method `ArtmConfigureLogging` will return error code `INVALID_OPERATION` if it detects an attempt to change logging folder after logging had been initialized. In order to set log directory the call to `ArtmConfigureLogging` must happen prior to calling any other methods in low-level C API. (with exception to `ArtmSetProtobufMessageFormatToJson`, `ArtmSetProtobufMessageFormatToBinary` and `ArtmProtobufMessageFormatIsJson`). This is because methods in `c_interface` may automatically initialize logging into current working directory, which later can not be changed.

Setting log directory require that target folder already exist on disk.

The following parameters can be customized with `ArtmConfigureLogging`.

```
message ConfigureLoggingArgs {
  // If specified, logfiles are written into this directory
  // instead of the default logging directory.
  optional string log_dir = 1;

  // Messages logged at a lower level than this
  // do not actually get logged anywhere
  optional int32 minloglevel = 2;

  // Log messages at a level >= this flag are automatically
  // sent to stderr in addition to log files.
  optional int32 stderrthreshold = 3;

  // Log messages go to stderr instead of logfiles
  optional bool logtostderr = 4;

  // color messages logged to stderr (if supported by terminal)
  optional bool colorlogtostderr = 5;

  // log messages go to stderr in addition to logfiles
  optional bool alsologtostderr = 6;

  // Buffer log messages for at most this many seconds
  optional int32 logbufsecs = 7;

  // Buffer log messages logged at this level or lower
```

```
// (-1 means do not buffer; 0 means buffer INFO only; ...)
optional int32 logbuflevel = 8;

// approx. maximum log file size (in MB). A value of 0 will be silently overridden to 1.
optional int32 max_log_size = 9;

// Stop attempting to log to disk if the disk is full.
optional bool stop_logging_if_full_disk = 10;
}
```

We recommend to set `logbuflevel = -1` to not buffer log messages. However by default BigARTM does not set this parameter, using the same default as provided by glog.

## ArtemReconfigureTopicName

To explain `ArtemReconfigureTopicName` we need to first start with `ArtemReconfigureMasterModel`. `ArtemReconfigureMasterModel` allow user to rename topics, but the number of topics must stay the same, and the order and the content of all existing phi matrices remains unchanged. On contrary, `ArtemReconfigureTopicName` may change the number of topics by adding or removing topics, as well as re-order columns of existing phi matrices. In `ArtemReconfigureMasterModel` the list of topic names is treated as new identifiers that should be set for existing columns. In `ArtemReconfigureTopicName` the list of topic names is matched against previous topic names. New topic names are added to phi matrices, topic names removed from the list are excluded from phi matrices, and topic names present in both old and new lists are re-ordered accordingly to match new topic name list.

Examples for `ArtemReconfigureMasterModel`:

- `t1, t2, t3 -> t4, t5, t6` sets new topic names for existing columns in phi matrices.

Examples for `ArtemReconfigureTopicName`:

- `t1, t2, t3 -> t1, t2, t3, t4` adds a new column to phi matrices, initialized with zeros
- `t1, t2, t3 -> t1, t2` removes last column from phi matrices
- `t1, t2, t3 -> t2, t3` removes the first column from phi matrices
- `t1, t2, t3 -> t3, t2` removes the first column from phi matrices and swaps the remaining two columns
- `t1, t2, t3 -> t4, t5, t6` removes all columns from phi matrices and creates three new columns, initialized with zeros

Note that both `ArtemReconfigureTopicName` and `ArtemReconfigureMasterModel` only affect phi matrices where set of topic names match the configuration of the master model. User-created matrices with custom set of topic names, for example created via `ArtemMergeModel`, will stay unchanged.

If you change topic names you should also consider changing your configuration of scores and regularizers. Also take into account that `ArtemReconfigureTopicName` and `ArtemReconfigureMasterModel` do not update theta cache. It is a good idea to call `ArtemClearThetaCache` after changing topic names.

## ArtemGetLastErrorMessage

```
const char* ArtemGetLastErrorMessage ()
```

Retrieves the textual error message, occurred during the last failing request.

## Error codes

```
#define ARTM_SUCCESS 0
#define ARTM_STILL_WORKING -1
#define ARTM_INTERNAL_ERROR -2
#define ARTM_ARGUMENT_OUT_OF_RANGE -3
#define ARTM_INVALID_MASTER_ID -4
#define ARTM_CORRUPTED_MESSAGE -5
#define ARTM_INVALID_OPERATION -6
#define ARTM_DISK_READ_ERROR -7
#define ARTM_DISK_WRITE_ERROR -8
```

### **ARTM\_SUCCESS**

The API call succeeded.

### **ARTM\_STILL\_WORKING**

This error code is applicable only to `ArtmAwaitOperation()`. It indicates that library is still processing the collection. Try to retrieve results later.

### **ARTM\_INTERNAL\_ERROR**

The API call failed due to internal error in BigARTM library. Please, collect steps to reproduce this issue and report it with BigARTM issue tracker.

### **ARTM\_ARGUMENT\_OUT\_OF\_RANGE**

The API call failed because one or more values of an argument are outside the allowable range of values as defined by the invoked method.

### **ARTM\_INVALID\_MASTER\_ID**

An API call that require *master\_id* parameter failed because MasterComponent with given ID does not exist.

### **ARTM\_CORRUPTED\_MESSAGE**

Unable to deserialize protocol buffer message.

### **ARTM\_INVALID\_OPERATION**

The API call is invalid in current state or due to provided parameters.

### **ARTM\_DISK\_READ\_ERROR**

The required files could not be read from disk.

### **ARTM\_DISK\_WRITE\_ERROR**

The required files could not be writtent to disk.



---

## Python Interface

---

This document describes all classes and functions in python interface of BigARTM library.

### ARTM model

This page describes ARTM class.

```
class artm.ARTM(num_topics=None, topic_names=None, num_processors=None, class_ids=None,
                 scores=None, regularizers=None, num_document_passes=10, reuse_theta=False,
                 dictionary=None, cache_theta=False, theta_columns_naming='id', seed=-1)

__init__(num_topics=None, topic_names=None, num_processors=None, class_ids=None,
         scores=None, regularizers=None, num_document_passes=10, reuse_theta=False, dic-
         tionary=None, cache_theta=False, theta_columns_naming='id', seed=-1)
```

#### Parameters

- **num\_topics** (*int*) – the number of topics in model, will be overwritten if topic\_names is set
- **num\_processors** (*int*) – how many threads will be used for model training, if not specified then number of threads will be detected by the lib
- **topic\_names** (*list of str*) – names of topics in model
- **class\_ids** (*dict*) – list of class\_ids and their weights to be used in model, key — class\_id, value — weight, if not specified then all class\_ids will be used
- **cache\_theta** (*bool*) – save or not the Theta matrix in model. Necessary if ARTM.get\_theta() usage expects
- **scores** (*list*) – list of scores (objects of artm.\*Score classes)
- **regularizers** (*list*) – list with regularizers (objects of artm.\*Regularizer classes)
- **num\_document\_passes** (*int*) – number of inner iterations over each document
- **dictionary** (*str or reference to Dictionary object*) – dictionary to be used for initialization, if None nothing will be done
- **reuse\_theta** (*bool*) – reuse Theta from previous iteration or not
- **theta\_columns\_naming** (*str*) – either 'id' or 'title', determines how to name columns (documents) in theta dataframe
- **seed** (*unsigned int or -1*) – seed for random initialization, -1 means no seed

**Important public fields**

- **regularizers**: contains dict of regularizers, included into model
- **scores**: contains dict of scores, included into model
- **score\_tracker**: contains dict of scoring results: key — score name, value — ScoreTracker object, which contains info about values of score on each synchronization (e.g. collection pass) in list

**Note**

- Here and anywhere in BigARTM empty `topic_names` or `class_ids` means that model (or regularizer, or score) should use all topics or `class_ids`.
- If some fields of regularizers or scores are not defined by user — internal lib defaults would be used.
- If field ‘`topic_names`’ is `None`, it will be generated by BigARTM and will be available using `ARTM.topic_names()`.

**dispose()**

**Description** free all native memory, allocated for this model

**Note**

- This method does not free memory occupied by dictionaries, because dictionaries are shared across all models
- ARTM class implements `__exit__` and `__del__` methods, which automatically call `dispose`.

**fit\_offline** (*batch\_vectorizer=None, num\_collection\_passes=1*)

**Description** proceeds the learning of topic model in offline mode

**Parameters**

- **batch\_vectorizer** (*object\_reference*) – an instance of BatchVectorizer class
- **num\_collection\_passes** (*int*) – number of iterations over whole given collection

**fit\_online** (*batch\_vectorizer=None, tau0=1024.0, kappa=0.7, update\_every=1, apply\_weight=None, decay\_weight=None, update\_after=None, async=False*)

**Description** proceeds the learning of topic model in online mode

**Parameters**

- **batch\_vectorizer** (*object\_reference*) – an instance of BatchVectorizer class
- **update\_every** (*int*) – the number of batches; model will be updated once per it
- **tau0** (*float*) – coefficient (see ‘Update formulas’ paragraph)
- **kappa** (*float*) (*float*) – power for tau0, (see ‘Update formulas’ paragraph)
- **update\_after** (*list of int*) – number of batches to be passed for Phi synchronizations
- **apply\_weight** (*list of float*) – weight of applying new counters
- **decay\_weight** (*list of float*) – weight of applying old counters
- **async** (*bool*) – use or not the async implementation of the EM-algorithm

**Note** `async=True` leads to impossibility of score extraction via `score_tracker`. Use `get_score()` instead.

### Update formulas

- The formulas for `decay_weight` and `apply_weight`:
- `update_count = current_processed_docs / (batch_size * update_every)`;
- `rho = pow(tau0 + update_count, -kappa)`;
- `decay_weight = 1-rho`;
- `apply_weight = rho`;
- if `apply_weight`, `decay_weight` and `update_after` are set, they will be used, otherwise the code below will be used (with `update_every`, `tau0` and `kappa`)

**get\_phi** (*topic\_names=None, class\_ids=None, model\_name=None*)

**Description** get custom Phi matrix of model. The extraction of the whole Phi matrix expects **ARTM.phi\_** call.

### Parameters

- **topic\_names** (*list of str*) – list with topics to extract, None value means all topics
- **class\_ids** (*list of str*) – list with class ids to extract, None means all class ids
- **model\_name** (*str*) – `self.model_pwt` by default, `self.model_nwt` is also reasonable to extract unnormalized counters

### Returns

- `pandas.DataFrame`: (data, columns, rows), where:
- columns — the names of topics in topic model;
- rows — the tokens of topic model;
- data — content of Phi matrix.

**get\_phi\_sparse** (*topic\_names=None, class\_ids=None, model\_name=None, eps=None*)

**Description** get phi matrix in sparse format

### Parameters

- **topic\_names** (*list of str*) – list with topics to extract, None value means all topics
- **class\_ids** (*list of str*) – list with class ids to extract, None means all class ids
- **model\_name** (*str*) – `self.model_pwt` by default, `self.model_nwt` is also reasonable to extract unnormalized counters
- **eps** (*float*) – threshold to consider values as zero

### Returns

- a 3-tuple of (data, rows, columns), where
- data — `scipy.sparse.csr_matrix` with values
- columns — the names of topics in topic model;
- rows — the tokens of topic model;

**get\_score** (*score\_name*)

**Description** get score after `fit_offline`, `fit_online` or `transform`

**Parameters** **score\_name** (*str*) – the name of the score to return

**get\_theta** (*topic\_names=None*)

**Description** get Theta matrix for training set of documents (or cached after transform)

**Parameters** **topic\_names** (*list of str*) – list with topics to extract, None means all topics

**Returns**

- pandas.DataFrame: (data, columns, rows), where:
- columns — the ids of documents, for which the Theta matrix was requested;
- rows — the names of topics in topic model, that was used to create Theta;
- data — content of Theta matrix.

**get\_theta\_sparse** (*topic\_names=None, eps=None*)

**Description** get Theta matrix in sparse format

**Parameters**

- **topic\_names** (*list of str*) – list with topics to extract, None means all topics
- **eps** (*float*) – threshold to consider values as zero

**Returns**

- a 3-tuple of (data, rows, columns), where
- data — scipy.sparse.csr\_matrix with values
- columns — the ids of documents;
- rows — the names of topics in topic model;

**info**

**Description** returns internal diagnostics information about the model

**initialize** (*dictionary=None*)

**Description** initialize topic model before learning

**Parameters** **dictionary** (*str or reference to Dictionary object*) – loaded BigARTM collection dictionary

**library\_version**

**Description** the version of BigARTM library in a MAJOR.MINOR.PATCH format

**load** (*filename, model\_name='p\_wt'*)

**Description** loads from disk the topic model saved by ARTM.save()

**Parameters**

- **filename** (*str*) – the name of file containing model
- **model\_name** (*str*) – the name of matrix to be saved, 'p\_wt' or 'n\_wt'

**Note**

- Loaded model will overwrite ARTM.topic\_names and class\_ids fields.
- All class\_ids weights will be set to 1.0, you need to specify them by hand if it's necessary.
- The method call will empty ARTM.score\_tracker.



- All regularizers and scores will be forgotten.
- etc.
- We strongly recommend you to reset all important parameters of the ARTM model, used earlier.

**remove\_theta()**

**Description** removes cached theta matrix

**reshape\_topics(topic\_names)**

**Description** update topic names of the model.

Adds, removes, and reorders columns of phi matrices according to the new set of topic names. New topics are initialized with zeros.

**save(filename, model\_name='p\_wt')**

**Description** saves one Phi-like matrix to disk

**Parameters**

- **filename** (*str*) – the name of file to store model
- **model\_name** (*str*) – the name of matrix to be saved, 'p\_wt' or 'n\_wt'

**topic\_names**

**Description** Gets or sets the list of topic names of the model.

**Note**

- Setting topic name allows you to put new labels on the existing topics. To add, remove or reorder topics use ARTM.reshape\_topics() method.
- In ARTM topic names are used just as string identifiers, which give a unique name to each column of the phi matrix. Typically you want to set topic names as something like "topic0", "topic1", etc. Later operations like get\_phi() allow you to specify which topics you need to retrieve. Most regularizers allow you to limit the set of topics they act upon. If you configure a rich set of regularizers it is important design your topic names according to how they are regularized. For example, you may use names obj0, obj1, ..., objN for *objective* topics (those where you enable sparsity regularizers), and back0, back1, ..., backM for *background* topics (those where you enable smoothing regularizers).

**transform(batch\_vectorizer=None, theta\_matrix\_type='dense\_theta', predict\_class\_id=None)**

**Description** find Theta matrix for new documents

**Parameters**

- **batch\_vectorizer** (*object\_reference*) – an instance of BatchVectorizer class
- **theta\_matrix\_type** (*str*) – type of matrix to be returned, possible values: 'dense\_theta', 'dense\_ptdw', 'cache', None, default='dense\_theta'
- **predict\_class\_id** (*str*) – class\_id of a target modality to predict. When this option is enabled the resulting columns of theta matrix will correspond to unique labels of a target modality. The values will represent p(c|d), which give the probability of class label c for document d.

**Returns**

- pandas.DataFrame: (data, columns, rows), where:
- columns — the ids of documents, for which the Theta matrix was requested;

- `rows` — the names of topics in topic model, that was used to create Theta;
- `data` — content of Theta matrix.

**Note**

- ‘dense\_ptdw’ mode provides simple access to values of  $p(t|w,d)$ . The resulting pandas.DataFrame object will contain a flat theta matrix (no 3D) where each item has multiple columns - as many as the number of tokens in that document. These columns will have the same item\_id. The order of columns with equal item\_id is the same as the order of tokens in the input data (batch.item.token\_id).

**transform\_sparse** (*batch\_vectorizer*, *eps=None*)

**Description** find Theta matrix for new documents as sparse scipy matrix

**Parameters**

- **batch\_vectorizer** (*object\_reference*) – an instance of BatchVectorizer class
- **eps** (*float*) – threshold to consider values as zero

**Returns**

- a 3-tuple of (data, rows, columns), where
- data — scipy.sparse.csr\_matrix with values
- columns — the ids of documents;
- rows — the names of topics in topic model;

## LDA model

This page describes LDA class.

```
class artm.LDA(num_topics=None, num_processors=None, cache_theta=False, dictionary=None, num_document_passes=10, seed=-1, alpha=0.01, beta=0.01, theta_columns_naming='id')
```

```
__init__(num_topics=None, num_processors=None, cache_theta=False, dictionary=None, num_document_passes=10, seed=-1, alpha=0.01, beta=0.01, theta_columns_naming='id')
```

**Parameters**

- **num\_topics** (*int*) – the number of topics in model, will be overwritten if topic\_names is set
- **num\_processors** (*int*) – how many threads will be used for model training, if not specified then number of threads will be detected by the lib
- **cache\_theta** (*bool*) – save or not the Theta matrix in model. Necessary if ARTM.get\_theta() usage expects
- **num\_document\_passes** (*int*) – number of inner iterations over each document
- **dictionary** (*str or reference to Dictionary object*) – dictionary to be used for initialization, if None nothing will be done
- **reuse\_theta** (*bool*) – reuse Theta from previous iteration or not
- **seed** (*unsigned int or -1*) – seed for random initialization, -1 means no seed

- **alpha** (*float*) – hyperparameter of Theta smoothing regularizer
- **beta** (*float or list of floats with len == num\_topics*) – hyperparameter of Phi smoothing regularizer
- **theta\_columns\_naming** (*str*) – either ‘id’ or ‘title’, determines how to name columns (documents) in theta dataframe

**Note**

- the type (not value!) of beta should not change after initialization: if it was scalar - it should stay scalar, if it was list - it should stay list.

**fit\_offline** (*batch\_vectorizer, num\_collection\_passes=1*)

**Description** proceeds the learning of topic model in offline mode

**Parameters**

- **batch\_vectorizer** (*object\_reference*) – an instance of BatchVectorizer class
- **num\_collection\_passes** (*int*) – number of iterations over whole given collection

**fit\_online** (*batch\_vectorizer, tau0=1024.0, kappa=0.7, update\_every=1*)

**Description** proceeds the learning of topic model in online mode

**Parameters**

- **batch\_vectorizer** (*object\_reference*) – an instance of BatchVectorizer class
- **update\_every** (*int*) – the number of batches; model will be updated once per it
- **tau0** (*float*) – coefficient (see ‘Update formulas’ paragraph)
- **kappa** (*float*) (*float*) – power for tau0, (see ‘Update formulas’ paragraph)
- **update\_after** (*list of int*) – number of batches to be passed for Phi synchronizations

**Update formulas**

- The formulas for decay\_weight and apply\_weight:
- $\text{update\_count} = \text{current\_processed\_docs} / (\text{batch\_size} * \text{update\_every})$ ;
- $\text{rho} = \text{pow}(\text{tau0} + \text{update\_count}, -\text{kappa})$ ;
- $\text{decay\_weight} = 1 - \text{rho}$ ;
- $\text{apply\_weight} = \text{rho}$ ;

**get\_theta** ()

**Description** get Theta matrix for training set of documents

**Returns**

- pandas.DataFrame: (data, columns, rows), where:
- columns — the ids of documents, for which the Theta matrix was requested;
- rows — the names of topics in topic model, that was used to create Theta;
- data — content of Theta matrix.

**get\_top\_tokens** (*num\_tokens=10, with\_weights=False*)

**Description** returns most probable tokens for each topic

**Parameters**

- **num\_tokens** (*int*) – number of top tokens to be returned
- **with\_weights** (*bool*) – return only tokens, or tuples (token, its p\_wt)

**Returns**

- list of lists of str, each internal list corresponds one topic in natural order, if with\_weights == False, or list, or list of lists of tules, each tuple is (str, float)

**initialize** (*dictionary*)

**Description** initialize topic model before learning

**Parameters** **dictionary** (*str or reference to Dictionary object*) – loaded BigARTM collection dictionary

**load** (*filename, model\_name='p\_wt'*)

**Description** loads from disk the topic model saved by LDA.save()

**Parameters**

- **filename** (*str*) – the name of file containing model
- **model\_name** (*str*) – the name of matrix to be saved, 'p\_wt' or 'n\_wt'

**Note**

- We strongly recommend you to reset all important parameters of the LDA model, used earlier.

**remove\_theta** ()

**Description** removes cached theta matrix

**save** (*filename, model\_name='p\_wt'*)

**Description** saves one Phi-like matrix to disk

**Parameters**

- **filename** (*str*) – the name of file to store model
- **model\_name** (*str*) – the name of matrix to be saved, 'p\_wt' or 'n\_wt'

**transform** (*batch\_vectorizer, theta\_matrix\_type='dense\_theta'*)

**Description** find Theta matrix for new documents

**Parameters**

- **batch\_vectorizer** (*object\_reference*) – an instance of BatchVectorizer class
- **theta\_matrix\_type** (*str*) – type of matrix to be returned, possible values: 'dense\_theta', None, default='dense\_theta'

**Returns**

- pandas.DataFrame: (data, columns, rows), where:
- columns — the ids of documents, for which the Theta matrix was requested;
- rows — the names of topics in topic model, that was used to create Theta;
- data — content of Theta matrix.

## hARTM

This page describes hARTM class.

```
class artm.hARTM(num_processors=None,    class_ids=None,    scores=None,    regularizers=None,
                  num_document_passes=10, reuse_theta=False, dictionary=None, cache_theta=False,
                  theta_columns_naming='id', seed=-1, tmp_files_path='')
```

```
    __init__(num_processors=None,    class_ids=None,    scores=None,    regularizers=None,
            num_document_passes=10, reuse_theta=False, dictionary=None, cache_theta=False,
            theta_columns_naming='id', seed=-1, tmp_files_path='')
```

**Description** a class for constructing topic hierarchy that is a sequence of tied artm.ARTM() models (levels)

### Parameters

- **num\_processors** (*int*) – how many threads will be used for model training, if not specified then number of threads will be detected by the lib
- **class\_ids** (*dict*) – list of class\_ids and their weights to be used in model, key — class\_id, value — weight, if not specified then all class\_ids will be used
- **cache\_theta** (*bool*) – save or not the Theta matrix in model. Necessary if ARTM.get\_theta() usage expects
- **scores** (*list*) – list of scores (objects of artm.\*Score classes)
- **regularizers** (*list*) – list with regularizers (objects of artm.\*Regularizer classes)
- **num\_document\_passes** (*int*) – number of inner iterations over each document
- **dictionary** (*str or reference to Dictionary object*) – dictionary to be used for initialization, if None nothing will be done
- **reuse\_theta** (*bool*) – reuse Theta from previous iteration or not
- **theta\_columns\_naming** (*str*) – either 'id' or 'title', determines how to name columns (documents) in theta dataframe
- **seed** (*unsigned int or -1*) – seed for random initialization, -1 means no seed
- **tmp\_files\_path** (*str*) – a path where to save temporary files (temporary solution), default value: current directory

### Usage

- **to construct hierarchy you have to learn several ARTM models:** `hier = artm.hARTM()` `level0 = hier.add_level(num_topics=5)` # returns artm.ARTM() instance # work with level0 as with usual model `level1 = hier.add_level(num_topics=25, parent_level_weight=1)` # work with level1 as with usual model # ...
- **to get the i-th level's model, use**  
`level = hier[i]`  
 or `level = hier.get_level(i)`
- **to iterate through levels use**  
`for level in hier:` # some work with level
- method `hier.del_level(...)` removes i-th level and all levels after it

- other hARTM methods correspond to those in ARTM class and call them sequentially for all levels of hierarchy from 0 to the last one. For example, to fit levels offline you may call `fit_offline` method of hARTM instance or of each level individually.

**add\_level** (*num\_topics=None, topic\_names=None, parent\_level\_weight=1*)

**Description** adds new level to the hierarchy

**Parameters**

- **num\_topics** (*int*) – the number of topics in level model, will be overwritten if parameter `topic_names` is set
- **topic\_names** (*list of str*) – names of topics in model
- **parent\_level\_weight** (*float*) – the coefficient of smoothing `n_wt` by `n_wa`, `a` enumerates parent topics

**Returns** ARTM or derived ARTM\_Level instance

**Notes**

- hierarchy structure assumes the number of topics on each following level is greater than on previous one
- work with returned value as with usual ARTM model
- to access any level, use `[]` or `get_level` method
- Important! You cannot add next level before previous one is initialized and fit.

**del\_level** (*level\_idx*)

**Description** removes *i*-th level and all following levels.

**Parameters** **level\_idx** (*int*) – the number of level from what to start removing if -1, the last level is removed

**dispose** ()

**Description** free all native memory, allocated for this hierarchy

**Note**

- This method does not free memory occupied by models' dictionaries, because dictionaries are shared across all models
- hARTM class implements `__exit__` and `__del__` methods, which automatically call `dispose`.

**fit\_offline** (*batch\_vectorizer, num\_collection\_passes=1*)

**Description** proceeds the learning of all hierarchy levels from 0 to the last one

**Parameters**

- **batch\_vectorizer** (*object\_referenece*) – an instance of BatchVectorizer class
- **num\_collection\_passes** (*int*) – number of iterations over whole given collection for each level

**Note**

- You cannot add add next level before previous one is fit. So use this method only when all levels are added, initialized and fit, for example, when you added one more regularizer or loaded hierarchy from disk.

**get\_level** (*level\_idx*)

**Description** access level

**Parameters** **level\_idx** (*int*) – the number of level to return

**Returns** specified level that is ARTM or derived ARTM\_Level instance

**get\_phi** (*class\_ids=None, model\_name=None*)

**Description** get level-wise horizontally stacked Phi matrices

**Parameters**

- **class\_ids** (*list of str*) – list with class ids to extract, None means all class ids
- **model\_name** (*str*) – self.model\_pwt by default, self.model\_nwt is also reasonable to extract unnormalized counters

**Returns**

- pandas.DataFrame: (data, columns, rows), where:
- **columns** — the names of topics in format **level\_X\_Y** where X is level index and Y is topic name;
- rows — the tokens of topic model;
- data — content of Phi matrix.

**Note**

- if you need to extract specified topics, use get\_phi() method of individual level model

**get\_theta** (*topic\_names=None*)

**Description** get level-wise vertically stacked Theta matrices for training set of documents

**Parameters** **topic\_names** (*list of str*) – list with topics to extract, None means all topics

**Returns**

- pandas.DataFrame: (data, columns, rows), where:
- columns — the ids of documents, for which the Theta matrix was requested;
- **rows** — the names of topics in format **level\_X\_Y** where X is level index and Y is topic name;
- data — content of Theta matrix.

**library\_version**

**Description** the version of BigARTM library in a MAJOR.MINOR.PATCH format

**load** (*path*)

**Description** loads models of already constructed hierarchy

**Parameters** **path** (*str*) – a path where hierarchy was saved by hARTM.save method

**Notes**

- Loaded models will overwrite ARTM.topic\_names and class\_ids fields of each level.
- All class\_ids weights will be set to 1.0, you need to specify them by hand if it's necessary.

- The method call will empty ARTM.score\_tracker of each level.
- All regularizers and scores will be forgotten.
- etc.
- We strongly recommend you to reset all important parameters of the ARTM models and hARTM, used earlier.

**save** (*path*)

**Description** saves all levels

**Parameters** **path** (*str*) – a path where to save hierarchy files This must be existing empty path, otherwise exception is raised

**transform** (*batch\_vectorizer*)

**Description** get level-wise vertically stacked Theta matrices for new documents

**Parameters** **batch\_vectorizer** (*object\_reference*) – an instance of BatchVectorizer class

**Returns**

- pandas.DataFrame: (data, columns, rows), where:
- columns — the ids of documents, for which the Theta matrix was requested;
- rows — the names of topics in format level\_X\_Y where X is level index and Y is topic name;
- data — content of Theta matrix.

**Note**

- to access p(tld, w) matrix or to predict class use transform method of hierarchy level individually

## Batches Utils

This page describes BatchVectorizer class.

```
class artm.BatchVectorizer (batches=None, collection_name=None, data_path='',
                           data_format='batches', target_folder=None, batch_size=1000,
                           batch_name_type='code', data_weight=1.0, n_wd=None, vocabulary=None, gather_dictionary=True, class_ids=None)
```

```
__init__ (batches=None, collection_name=None, data_path='', data_format='batches', target_folder=None, batch_size=1000, batch_name_type='code', data_weight=1.0, n_wd=None, vocabulary=None, gather_dictionary=True, class_ids=None)
```

**Parameters**

- **collection\_name** (*str*) – the name of text collection (required if data\_format == 'bow\_uci')
- **data\_path** (*str*) –
  1. if data\_format == 'bow\_uci' => folder containing 'docword.collection\_name.txt' and vocab.collection\_name.txt files; 2) if data\_format == 'vowpal\_wabbit' => file in Vowpal Wabbit format; 3) if data\_format == 'bow\_n\_wd' => useless parameter
  - 4) if data\_format == 'batches' => folder containing batches



- **data\_format** (*str*) – the type of input data: 1) ‘bow\_uci’ — Bag-Of-Words in UCI format; 2) ‘vowpal\_wabbit’ — Vowpal Wabbit format; 3 ‘bow\_n\_wd’ — result of CountVectorizer or similar tool; 4) ‘batches’ — the BigARTM data format
- **batch\_size** (*int*) – number of documents to be stored in each batch
- **target\_folder** (*str*) – full path to folder for future batches storing; if not set, no batches will be produced for further work
- **batches** (*list of str*) – list with non-full file names of batches (necessary parameters are batches + data\_path + data\_format==‘batches’ in this case)
- **batch\_name\_type** (*str*) – name batches in natural order (‘code’) or using random guids (guid)
- **data\_weight** (*float*) – weight for a group of batches from data\_path; it can be a list of floats, then data\_path (and target\_folder if not data\_format == ‘batches’) should also be lists; one weight corresponds to one path from the data\_path list;
- **n\_wd** (*array*) – matrix with n\_wd counters
- **vocabulary** (*dict*) – dict with vocabulary, key - index of n\_wd, value - token
- **gather\_dictionary** (*bool*) – create or not the default dictionary in vectorizer; if data\_format == ‘bow\_n\_wd’ - automatically set to True; and if data\_weight is list - automatically set to False
- **class\_ids** (*list of str*) – list of class ids to parse and include in batches

**batch\_size**

**Returns** the user-defined size of the batches

**batches\_list**

**Returns** list of batches names

**data\_path**

**Returns** the disk path of batches

**dictionary**

**Returns** Dictionary object, if parameter gather\_dictionary was True, else None

**num\_batches**

**Returns** the number of batches

**weights**

**Returns** list of batches weights

## Dictionary

This page describes Dictionary class.

**class** artm.Dictionary (*name=None, dictionary\_path=None, data\_path=None*)

**\_\_init\_\_** (*name=None, dictionary\_path=None, data\_path=None*)

**Parameters**

- **name** (*str*) – name of the dictionary

- **dictionary\_path** (*str*) – can be used for default call of load() method in constructor
- **data\_path** (*str*) – can be used for default call of gather() method in constructor

Note: all parameters are optional

**copy** ()

**Description** returns a copy the dictionary loaded in lib with another name.

**create** (*dictionary\_data*)

**Description** creates dictionary using DictionaryData object

**Parameters** **dictionary\_data** (*DictionaryData instance*) – configuration of dictionary

**filter** (*class\_id=None, min\_df=None, max\_df=None, min\_df\_rate=None, max\_df\_rate=None, min\_tf=None, max\_tf=None, max\_dictionary\_size=None*)

**Description** filters the BigARTM dictionary of the collection, which was already loaded into the lib

**Parameters**

- **dictionary\_name** (*str*) – name of the dictionary in the lib to filter
- **dictionary\_target\_name** (*str*) – name for the new filtered dictionary in the lib
- **class\_id** (*str*) – class\_id to filter
- **min\_df** (*float*) – min df value to pass the filter
- **max\_df** (*float*) – max df value to pass the filter
- **min\_df\_rate** (*float*) – min df rate to pass the filter
- **max\_df\_rate** (*float*) – max df rate to pass the filter
- **min\_tf** (*float*) – min tf value to pass the filter
- **max\_tf** (*float*) – max tf value to pass the filter
- **max\_dictionary\_size** (*float*) – give an easy option to limit dictionary size; rare tokens will be excluded until dictionary reaches given size.

**Note** the current dictionary will be replaced with filtered

**gather** (*data\_path, cooc\_file\_path=None, vocab\_file\_path=None, symmetric\_cooc\_values=False*)

**Description** creates the BigARTM dictionary of the collection, represented as batches and load it in the lib

**Parameters**

- **data\_path** (*str*) – full path to batches folder
- **cooc\_file\_path** (*str*) – full path to the file with cooc info. Cooc info is a file with three columns, first two are the zero-based indices of tokens in vocab file, and third one is a value of their cooccurrence in collection (or another) pairwise statistic.
- **vocab\_file\_path** (*str*) – full path to the file with vocabulary. If given, the dictionary token will have the same order, as in this file, otherwise the order will be random. If given, the tokens from batches, that are not presented in vocab, will be skipped.

- **symmetric\_cooc\_values** (*bool*) – if the cooc matrix should considered to be symmetric or not

**load** (*dictionary\_path*)

**Description** loads the BigARTM dictionary of the collection into the lib

**Parameters** **dictionary\_path** (*str*) – full filename of the dictionary

**load\_text** (*dictionary\_path*, *encoding*='utf-8')

**Description** loads the BigARTM dictionary of the collection from the disk in the human readable text format

**Parameters**

- **dictionary\_path** (*str*) – full file name of the text dictionary file
- **encoding** (*str*) – an encoding of text in dictionary

**save** (*dictionary\_path*)

**Description** saves the BigARTM dictionary of the collection on the disk

**Parameters** **dictionary\_path** (*str*) – full file name for the dictionary

**save\_text** (*dictionary\_path*, *encoding*='utf-8')

**Description** saves the BigARTM dictionary of the collection on the disk in the human readable text format

**Parameters**

- **dictionary\_path** (*str*) – full file name for the text dictionary file
- **encoding** (*str*) – an encoding of text in dictionary

## Regularizers

This page describes *KlFunctionInfo* and *Regularizer* classes.

See [detailed description of regularizers](#) for understanding their sense.

**class** `artm.KlFunctionInfo` (*function\_type*='log', *power\_value*=2.0)

**\_\_init\_\_** (*function\_type*='log', *power\_value*=2.0)

**Parameters**

- **function\_type** (*str*) – the type of function, 'log' (logarithm) or 'pol' (polynomial)
- **power\_value** (*float*) – the double power of polynomial, ignored if type = 'log'

**class** `artm.SmoothSparsePhiRegularizer` (*name*=None, *tau*=1.0, *gamma*=None, *class\_ids*=None, *topic\_names*=None, *dictionary*=None, *kl\_function\_info*=None, *config*=None)

**\_\_init\_\_** (*name*=None, *tau*=1.0, *gamma*=None, *class\_ids*=None, *topic\_names*=None, *dictionary*=None, *kl\_function\_info*=None, *config*=None)

**Parameters**

- **name** (*str*) – the identifier of regularizer, will be auto-generated if not specified

- **tau** (*float*) – the coefficient of regularization for this regularizer
- **gamma** (*float*) – the coefficient of relative regularization for this regularizer
- **class\_ids** (*list of str*) – list of class\_ids to regularize, will regularize all classes if not specified
- **topic\_names** (*list of str*) – list of names of topics to regularize, will regularize all topics if not specified
- **dictionary** (*str or reference to Dictionary object*) – BigARTM collection dictionary, won't use dictionary if not specified
- **kl\_function\_info** (*KlFunctionInfo object*) – class with additional info about function under KL-div in regularizer
- **config** (*protobuf object*) – the low-level config of this regularizer

```
class artm.SmoothSparseThetaRegularizer(name=None, tau=1.0, topic_names=None,
                                         alpha_iter=None, kl_function_info=None,
                                         doc_titles=None, doc_topic_coef=None, config=None)
```

```
__init__(name=None, tau=1.0, topic_names=None, alpha_iter=None, kl_function_info=None,
          doc_titles=None, doc_topic_coef=None, config=None)
```

#### Parameters

- **name** (*str*) – the identifier of regularizer, will be auto-generated if not specified
- **tau** (*float*) – the coefficient of regularization for this regularizer
- **alpha\_iter** (*list of str*) – list of additional coefficients of regularization on each iteration over document. Should have length equal to `model.num_document_passes`
- **topic\_names** (*list of str*) – list of names of topics to regularize, will regularize all topics if not specified
- **kl\_function\_info** (*KlFunctionInfo object*) – class with additional info about function under KL-div in regularizer
- **doc\_titles** (*list of strings*) – list of titles of documents to be processed by this regularizer. Default empty value means processing of all documents. User should guarantee the existence and correctness of document titles in batches (e.g. in src files with data, like WV).
- **doc\_topic\_coef** (*list of doubles or list of lists of doubles*) – Two cases: 1) list of doubles with length equal to num of topics. Means additional multiplier in M-step formula besides alpha and tau, unique for each topic, but general for all processing documents. 2) list of lists of doubles with outer list length equal to length of doc\_titles, and each inner list length equal to num of topics. Means case 1 with unique list of additional multipliers for each document from doc\_titles. Other documents will not be regularized according to description of doc\_titles parameter. Note, that doc\_topic\_coef and topic\_names are both using.
- **config** (*protobuf object*) – the low-level config of this regularizer

```
class artm.DecorrelatorPhiRegularizer(name=None, tau=1.0, gamma=None, class_ids=None,
                                       topic_names=None, config=None)
```

```
__init__(name=None, tau=1.0, gamma=None, class_ids=None, topic_names=None, config=None)
```

#### Parameters

- **name** (*str*) – the identifier of regularizer, will be auto-generated if not specified
- **tau** (*float*) – the coefficient of regularization for this regularizer
- **gamma** (*float*) – the coefficient of relative regularization for this regularizer
- **class\_ids** (*list of str*) – list of class\_ids to regularize, will regularize all classes if not specified
- **topic\_names** (*list of str*) – list of names of topics to regularize, will regularize all topics if not specified
- **config** (*protobuf object*) – the low-level config of this regularizer

```
class artm.LabelRegularizationPhiRegularizer (name=None, tau=1.0, gamma=None,
                                             class_ids=None, topic_names=None, dictionary=None, config=None)
```

```
__init__(name=None, tau=1.0, gamma=None, class_ids=None, topic_names=None, dictionary=None, config=None)
```

#### Parameters

- **name** (*str*) – the identifier of regularizer, will be auto-generated if not specified
- **tau** (*float*) – the coefficient of regularization for this regularizer
- **gamma** (*float*) – the coefficient of relative regularization for this regularizer
- **class\_ids** (*list of str*) – list of class\_ids to regularize, will regularize all classes if not specified
- **topic\_names** (*list of str*) – list of names of topics to regularize, will regularize all topics if not specified
- **dictionary** (*str or reference to Dictionary object*) – BigARTM collection dictionary, won't use dictionary if not specified
- **config** (*protobuf object*) – the low-level config of this regularizer

```
class artm.SpecifiedSparsePhiRegularizer (name=None, tau=1.0, gamma=None,
                                          topic_names=None, class_id=None,
                                          num_max_elements=None, probability_threshold=None,
                                          sparse_by_columns=True, config=None)
```

```
__init__(name=None, tau=1.0, gamma=None, topic_names=None, class_id=None,
         num_max_elements=None, probability_threshold=None, sparse_by_columns=True,
         config=None)
```

#### Parameters

- **name** (*str*) – the identifier of regularizer, will be auto-generated if not specified
- **tau** (*float*) – the coefficient of regularization for this regularizer
- **gamma** (*float*) – the coefficient of relative regularization for this regularizer
- **class\_id** – class\_id to regularize
- **topic\_names** (*list of str*) – list of names of topics to regularize, will regularize all topics if not specified
- **num\_max\_elements** (*int*) – number of elements to save in row/column

- **probability\_threshold** (*float*) – if  $m$  elements in row/column sum into value  $\geq$  probability\_threshold,  $m < n \Rightarrow$  only these elements would be saved. Value should be in (0, 1), default=None
- **sparse\_by\_columns** (*bool*) – find max elements in column or in row
- **config** (*protobuf object*) – the low-level config of this regularizer

```
class artm.ImproveCoherencePhiRegularizer(name=None, tau=1.0, gamma=None,
                                          class_ids=None, topic_names=None, dictionary=None, config=None)
```

```
__init__(name=None, tau=1.0, gamma=None, class_ids=None, topic_names=None, dictionary=None, config=None)
```

#### Parameters

- **name** (*str*) – the identifier of regularizer, will be auto-generated if not specified
- **tau** (*float*) – the coefficient of regularization for this regularizer
- **gamma** (*float*) – the coefficient of relative regularization for this regularizer
- **class\_ids** (*list of str*) – list of class\_ids to regularize, will regularize all classes if not specified, dictionary should contain pairwise tokens cooccurancy info
- **topic\_names** (*list of str*) – list of names of topics to regularize, will regularize all topics if not specified
- **dictionary** (*str or reference to Dictionary object*) – BigARTM collection dictionary, won't use dictionary if not specified, in this case regularizer is useless
- **config** (*protobuf object*) – the low-level config of this regularizer

```
class artm.SmoothPtdwRegularizer(name=None, tau=1.0, config=None)
```

```
__init__(name=None, tau=1.0, config=None)
```

#### Parameters

- **name** (*str*) – the identifier of regularizer, will be auto-generated if not specified
- **tau** (*float*) – the coefficient of regularization for this regularizer
- **config** (*protobuf object*) – the low-level config of this regularizer

```
class artm.TopicSelectionThetaRegularizer(name=None, tau=1.0, topic_names=None, alpha_iter=None, config=None)
```

```
__init__(name=None, tau=1.0, topic_names=None, alpha_iter=None, config=None)
```

#### Parameters

- **name** (*str*) – the identifier of regularizer, will be auto-generated if not specified
- **tau** (*float*) – the coefficient of regularization for this regularizer
- **alpha\_iter** (*list of str*) – list of additional coefficients of regularization on each iteration over document. Should have length equal to model.num\_document\_passes
- **topic\_names** (*list of str*) – list of names of topics to regularize, will regularize all topics if not specified
- **config** (*protobuf object*) – the low-level config of this regularizer

```
class artm.TopicSegmentationPtdwRegularizer (name=None, window=None, threshold=None,
                                             background_topic_names=None, config=None)
```

```
__init__ (name=None, window=None, threshold=None, background_topic_names=None, config=None)
```

#### Parameters

- **name** (*str*) – the identifier of regularizer, will be auto-generated if not specified
- **window** (*int*) – a number of words to the one side over which smoothing will be performed
- **threshold** (*float*) – probability threshold for a word to be a topic-changing word
- **background\_topic\_names** (*list of str*) – list of names of topics to be considered background, will not consider background topics if not specified
- **config** (*protobuf object*) – the low-level config of this regularizer

## Scores

This page describes \*Scores classes.

See [detailed description of scores](#) for understanding their sense.

```
class artm.SparsityPhiScore (name=None, class_id=None, topic_names=None, model_name=None,
                             eps=None)
```

```
__init__ (name=None, class_id=None, topic_names=None, model_name=None, eps=None)
```

#### Parameters

- **name** (*str*) – the identifier of score, will be auto-generated if not specified
- **class\_id** (*str*) – class\_id to score
- **topic\_names** (*list of str*) – list of names of topics to regularize, will score all topics if not specified
- **model\_name** – phi-like matrix to be scored (typically ‘pwt’ or ‘nwt’), ‘pwt’ if not specified
- **eps** (*float*) – the tolerance const, everything < eps considered to be zero

```
class artm.ItemsProcessedScore (name=None)
```

```
__init__ (name=None)
```

**Parameters** **name** (*str*) – the identifier of score, will be auto-generated if not specified

```
class artm.PerplexityScore (name=None, class_ids=None, topic_names=None, dictionary=None)
```

```
__init__ (name=None, class_ids=None, topic_names=None, dictionary=None)
```

#### Parameters

- **name** (*str*) – the identifier of score, will be auto-generated if not specified
- **class\_ids** (*list of str*) – class\_id to score, means that tokens of all class\_ids will be used

- **dictionary** (*str or reference to Dictionary object*) – BigARTM collection dictionary, is strongly recommended to be used for correct replacing of zero counters.

```
class artm.SparsityThetaScore(name=None, topic_names=None, eps=None)
```

```
__init__(name=None, topic_names=None, eps=None)
```

#### Parameters

- **name** (*str*) – the identifier of score, will be auto-generated if not specified
- **topic\_names** (*list of str*) – list of names of topics to regularize, will score all topics if not specified
- **eps** (*float*) – the tolerance const, everything < eps considered to be zero

```
class artm.ThetaSnippetScore(name=None, item_ids=None, num_items=None)
```

```
__init__(name=None, item_ids=None, num_items=None)
```

#### Parameters

- **name** (*str*) – the identifier of score, will be auto-generated if not specified
- **item\_ids** (*list of int*) – list of names of items to show, default=None
- **num\_items** (*int*) – number of theta vectors to show from the beginning (no sense if item\_ids was given)

```
class artm.TopicKernelScore(name=None, class_id=None, topic_names=None, eps=None, dictionary=None, probability_mass_threshold=None)
```

```
__init__(name=None, class_id=None, topic_names=None, eps=None, dictionary=None, probability_mass_threshold=None)
```

#### Parameters

- **name** (*str*) – the identifier of score, will be auto-generated if not specified
- **class\_id** (*str*) – class\_id to score
- **topic\_names** (*list of str*) – list of names of topics to regularize, will score all topics if not specified
- **probability\_mass\_threshold** (*float*) – the threshold for p(tlw) values to get token into topic kernel. Should be in (0, 1)
- **dictionary** (*str or reference to Dictionary object*) – BigARTM collection dictionary, won't use dictionary if not specified
- **eps** (*float*) – the tolerance const, everything < eps considered to be zero

```
class artm.TopTokensScore(name=None, class_id=None, topic_names=None, num_tokens=None, dictionary=None)
```

```
__init__(name=None, class_id=None, topic_names=None, num_tokens=None, dictionary=None)
```

#### Parameters

- **name** (*str*) – the identifier of score, will be auto-generated if not specified
- **class\_id** (*str*) – class\_id to score



- **topic\_names** (*list of str*) – list of names of topics to regularize, will score all topics if not specified
- **num\_tokens** (*int*) – number of tokens with max probability in each topic
- **dictionary** (*str or reference to Dictionary object*) – Bi-gARTM collection dictionary, won't use dictionary if not specified

```
class artm.TopicMassPhiScore (name=None, class_id=None, topic_names=None, model_name=None,
                             eps=None)
```

```
__init__ (name=None, class_id=None, topic_names=None, model_name=None, eps=None)
```

#### Parameters

- **name** (*str*) – the identifier of score, will be auto-generated if not specified
- **class\_id** (*str*) – class\_id to score
- **topic\_names** (*list of str*) – list of names of topics to regularize, will score all topics if not specified
- **model\_name** – phi-like matrix to be scored (typically 'pwt' or 'nwt'), 'pwt' if not specified
- **eps** (*float*) – the tolerance const, everything < eps considered to be zero

```
class artm.ClassPrecisionScore (name=None)
```

```
__init__ (name=None)
```

**Parameters** **name** (*str*) – the identifier of score, will be auto-generated if not specified

```
class artm.BackgroundTokensRatioScore (name=None, class_id=None, delta_threshold=None,
                                       save_tokens=None, direct_kl=None)
```

```
__init__ (name=None, class_id=None, delta_threshold=None, save_tokens=None, direct_kl=None)
```

#### Parameters

- **name** (*str*) – the identifier of score, will be auto-generated if not specified
- **class\_id** (*str*) – class\_id to score
- **delta\_threshold** (*float*) – the threshold for KL-div between p(tlw) and p(t) to get token into background. Should be non-negative
- **save\_tokens** (*bool*) – save background tokens or not, save if field not specified
- **direct\_kl** (*bool*) – use  $KL(p(t) \parallel p(tlw))$  or via versa, true if field not specified

## Score Tracker

This page describes \*ScoreTracker classes.

```
class artm.score_tracker.SparsityPhiScoreTracker (score)
```

```
__init__ (score)
```

#### Properties

- Note: every field is a list of info about score on all synchronizations.

- value - values of Phi sparsity.
- zero\_tokens - number of zero rows in Phi.
- total\_tokens - number of all rows in Phi.
- Note: every field has a version with prefix '**last\_**', means retrieving only info about the last synchronization.

```
class artm.score_tracker.SparsityThetaScoreTracker (score)
```

```
    __init__ (score)
```

#### Properties

- Note: every field is a list of info about score on all synchronizations.
- value - values of Theta sparsity.
- zero\_topics - number of zero rows in Theta.
- total\_topics - number of all rows in Theta.
- Note: every field has a version with prefix '**last\_**', means retrieving only info about the last synchronization.

```
class artm.score_tracker.PerplexityScoreTracker (score)
```

```
    __init__ (score)
```

#### Properties

- Note: every field is a list of info about score on all synchronizations.
- value - values of perplexity.
- raw - raw values in formula for perplexity.
- normalizer - normalizer values in formula for perplexity.
- zero\_tokens - number of zero  $p(w|d) = \sum_t p(w|t) p(t|d)$ .
- Note: every field has a version with prefix '**last\_**', means retrieving only info about the last synchronization.

```
class artm.score_tracker.TopTokensScoreTracker (score)
```

```
    __init__ (score)
```

#### Properties

- Note: every field is a list of info about score on all synchronizations.
- num\_tokens - number of requested top tokens.
- coherence - each element of list is a dict, key - topic name, value - topic coherence counted using top-tokens
- average\_coherence - average coherencies of all scored topics.
- tokens - each element of list is a dict, key - topic name, value - list of top-tokens

- weights - each element of list is a dict, key - topic name, value - list of weights of corresponding top-tokens (weight of token ==  $p(w|t)$ )
- Note: every field has a version with prefix '**last\_**', means retrieving only info about the last synchronization.

```
class artm.score_tracker.TopicKernelScoreTracker(score)
```

```
    __init__(score)
```

#### Properties

- Note: every field is a list of info about score on all synchronizations.
- tokens - each element of list is a dict, key - topic name, value - list of kernel tokens
- size - each element of list is a dict, key - topic name, value - kernel size
- contrast - each element of list is a dict, key - topic name, value - kernel contrast
- purity - each element of list is a dict, key - topic name, value - kernel purity
- coherence - each element of list is a dict, key - topic name, value - topic coherence counted using kernel tokens
- average\_size - average kernel size of all scored topics.
- average\_contrast - average kernel contrast of all scored topics.
- average\_purity - average kernel purity of all scored topics.
- average\_coherence - average coherencies of all scored topics.
- Note: every field has a version with prefix '**last\_**', means retrieving only info about the last synchronization.

```
class artm.score_tracker.ItemsProcessedScoreTracker(score)
```

```
    __init__(score)
```

#### Properties

- Note: every field is a list of info about score on all synchronizations.
- value - numbers of processed documents.
- Note: every field has a version with prefix '**last\_**', means retrieving only info about the last synchronization.

```
class artm.score_tracker.ThetaSnippetScoreTracker(score)
```

```
    __init__(score)
```

#### Properties

- Note: every field is a list of info about score on all synchronizations.
- document\_ids - each element of list is a list of ids of returned documents.
- snippet - each element of list is a dict, key - doc id, value - list with corresponding  $p(t|d)$  values.

- Note: every field has a version with prefix '**last\_**', means retrieving only info about the last synchronization.

```
class artm.score_tracker.TopicMassPhiScoreTracker(score)
```

```
    __init__(score)
```

#### Properties

- Note: every field is a list of info about score on all synchronizations.
- value - values of ratio of sum\_t n\_t of scored topics and all topics
- topic\_mass - each value is a dict, key - topic name, value - topic mass n\_t
- topic\_ratio - each value is a dict, key - topic name, value - topic ratio
- Note: every field has a version with prefix '**last\_**', means retrieving only info about the last synchronization.

```
class artm.score_tracker.ClassPrecisionScoreTracker(score)
```

```
    __init__(score)
```

#### Properties

- Note: every field is a list of info about score on all synchronizations.
- value - values of ratio of correct predictions.
- error - numbers of error predictions.
- total - numbers of all predictions.
- Note: every field has a version with prefix '**last\_**', means retrieving only info about the last synchronization.

```
class artm.score_tracker.BackgroundTokensRatioScoreTracker(score)
```

```
    __init__(score)
```

#### Properties

- Note: every field is a list of info about score on all synchronizations.
- value - values of part of background tokens.
- tokens - each element of list is a list of background tokens (can be accessed if 'save\_tokens' was True)
- Note: every field has a version with prefix '**last\_**', means retrieving only info about the last synchronization.

## Master Component

This page describes MasterComponent class.

```
class artm.MasterComponent (library, topic_names=None, class_ids=None, scores=None, regulariz-  
ers=None, num_processors=None, pwt_name=None, nwt_name=None,  
num_document_passes=None, reuse_theta=None, cache_theta=False)
```

```
__init__ (library, topic_names=None, class_ids=None, scores=None, regulariz-  
ers=None, num_processors=None, pwt_name=None, nwt_name=None,  
num_document_passes=None, reuse_theta=None, cache_theta=False)
```

#### Parameters

- **library** – an instance of LibArtm
- **topic\_names** (*list of str*) – list of topic names to use in model
- **class\_ids** (*dict*) – key - class\_id, value - class\_weight
- **scores** (*dict*) – key - score name, value - config
- **regularizers** (*dict*) – key - regularizer name, value - tuple (config, tau) or triple (config, tau, gamma)
- **num\_processors** (*int*) – number of worker threads to use for processing the collection
- **pwt\_name** (*str*) – name of pwt matrix
- **nwt\_name** (*str*) – name of nwt matrix
- **num\_document\_passes** (*in*) – num passes through each document
- **reuse\_theta** (*bool*) – reuse Theta from previous iteration or not
- **cache\_theta** (*bool*) – save or not the Theta matrix

```
attach_model (model)
```

**Parameters** **model** (*str*) – name of matrix in BigARTM

#### Returns

- messahes.TopicModel() object with info about Phi matrix
- numpy.ndarray with Phi data (i.e., p(wlt) values)

```
clear_score_array_cache ()  
Clears all entries from score array cache
```

```
clear_score_cache ()  
Clears all entries from score cache
```

```
clear_theta_cache ()  
Clears all entries from theta matrix cache
```

```
create_dictionary (dictionary_data, dictionary_name=None)
```

#### Parameters

- **dictionary\_data** – an instance of DictionaryData with info about dictionary
- **dictionary\_name** (*str*) – name of exported dictionary

```
create_regularizer (name, config, tau, gamma=None)
```

#### Parameters

- **name** (*str*) – the name of the future regularizer
- **config** – the config of the future regularizer

- **tau** (*float*) – the coefficient of the regularization

**create\_score** (*name, config, model\_name=None*)

**Parameters**

- **name** (*str*) – the name of the future score
- **config** – an instance of `ScoreConfig`

**export\_dictionary** (*filename, dictionary\_name*)

**Parameters**

- **filename** (*str*) – full name of dictionary file
- **dictionary\_name** (*str*) – name of exported dictionary

**export\_model** (*model, filename*)

**filter\_dictionary** (*dictionary\_name=None, dictionary\_target\_name=None, class\_id=None, min\_df=None, max\_df=None, min\_df\_rate=None, max\_df\_rate=None, min\_tf=None, max\_tf=None, max\_dictionary\_size=None, args=None*)

**Parameters**

- **dictionary\_name** (*str*) – name of the dictionary in the core to filter
- **dictionary\_target\_name** (*str*) – name for the new filtered dictionary in the core
- **class\_id** (*str*) – class\_id to filter
- **min\_df** (*float*) – min df value to pass the filter
- **max\_df** (*float*) – max df value to pass the filter
- **min\_df\_rate** (*float*) – min df rate to pass the filter
- **max\_df\_rate** (*float*) – max df rate to pass the filter
- **min\_tf** (*float*) – min tf value to pass the filter
- **max\_tf** (*float*) – max tf value to pass the filter
- **max\_dictionary\_size** (*float*) – give an easy option to limit dictionary size; rare tokens will be excluded until dictionary reaches given size.
- **args** – an instance of `FilterDictionaryArgs`

**fit\_offline** (*batch\_filenames=None, batch\_weights=None, num\_collection\_passes=None, batches\_folder=None*)

**Parameters**

- **batch\_filenames** (*list of str*) – name of batches to process
- **batch\_weights** (*list of float*) – weights of batches to process
- **num\_collection\_passes** (*int*) – number of outer iterations
- **batches\_folder** (*str*) – folder containing batches to process

**fit\_online** (*batch\_filenames=None, batch\_weights=None, update\_after=None, apply\_weight=None, decay\_weight=None, async=None*)

**Parameters**

- **batch\_filenames** (*list of str*) – name of batches to process

- **batch\_weights** (*list of float*) – weights of batches to process
- **update\_after** (*list of int*) – number of batches to be passed for Phi synchronizations
- **apply\_weight** (*list of float*) – weight of applying new counters (len == len of update\_after)
- **decay\_weight** (*list of float*) – weight of applying old counters (len == len of update\_after)
- **async** (*bool*) – whether to use the async implementation of the EM-algorithm or not

**gather\_dictionary** (*dictionary\_target\_name=None, data\_path=None, cooc\_file\_path=None, vocab\_file\_path=None, symmetric\_cooc\_values=None, args=None*)

#### Parameters

- **dictionary\_target\_name** (*str*) – name of the dictionary in the core
- **data\_path** (*str*) – full path to batches folder
- **cooc\_file\_path** (*str*) – full path to the file with cooc info
- **vocab\_file\_path** (*str*) – full path to the file with vocabulary
- **symmetric\_cooc\_values** (*bool*) – whether the cooc matrix should be considered to be symmetric or not
- **args** – an instance of GatherDictionaryArgs

**get\_dictionary** (*dictionary\_name*)

Parameters **dictionary\_name** (*str*) – name of dictionary to get

**get\_info** ()

**get\_phi\_info** (*model*)

Parameters **model** (*str*) – name of matrix in BigARTM

Returns messages.TopicModel object

**get\_phi\_matrix** (*model, topic\_names=None, class\_ids=None, use\_sparse\_format=None*)

#### Parameters

- **model** (*str*) – name of matrix in BigARTM
- **topic\_names** (*list of str or None*) – list of topics to retrieve (None means all topics)
- **class\_ids** (*list of str or None*) – list of class ids to retrieve (None means all class ids)
- **use\_sparse\_format** (*bool*) – use sparsedense layout

Returns numpy.ndarray with Phi data (i.e., p(w|t) values)

**get\_score** (*score\_name*)

#### Parameters

- **score\_name** (*str*) – the user defined name of score to retrieve
- **score\_config** – reference to score data object

**get\_score\_array** (*score\_name*)

**Parameters**

- **score\_name** (*str*) – the user defined name of score to retrieve
- **score\_config** – reference to score data object

**get\_theta\_info**()

**Returns** messages.ThetaMatrix object

**get\_theta\_matrix** (*topic\_names=None*)

**Parameters** **topic\_names** (*list of str or None*) – list of topics to retrieve (None means all topics)

**Returns** numpy.ndarray with Theta data (i.e., p(t|d) values)

**import\_dictionary** (*filename, dictionary\_name*)

**Parameters**

- **filename** (*str*) – full name of dictionary file
- **dictionary\_name** (*str*) – name of imported dictionary

**import\_model** (*model, filename*)

**Parameters**

- **model** (*str*) – name of matrix in BigARTM
- **filename** (*str*) – the name of file to load model from binary format

**initialize\_model** (*model\_name=None, topic\_names=None, dictionary\_name=None, seed=None, args=None*)

**Parameters**

- **model\_name** (*str*) – name of pwt matrix in BigARTM
- **topic\_names** (*list of str*) – the list of names of topics to be used in model
- **dictionary\_name** (*str*) – name of imported dictionary
- **seed** (*unsigned int or -1, default None*) – seed for random initialization, None means no seed
- **args** – an instance of InitilaizeModelArgs

**merge\_model** (*models, nwt, topic\_names=None, dictionary\_name=None*)

Merge multiple nwt-increments together.

**Parameters**

- **models** (*dict*) – list of models with nwt-increments and their weights, key - nwt\_source\_name, value - source\_weight.
- **nwt** (*str*) – the name of target matrix to store combined nwt. The matrix will be created by this operation.
- **topic\_names** (*list of str*) – names of topics in the resulting model. By default model names are taken from the first model in the list.
- **dictionary\_name** – name of dictionary that defines which tokens to include in merged model

**normalize\_model** (*pwt, nwt, rwt=None*)

**Parameters**



- **pwt** (*str*) – name of pwt matrix in BigARTM
- **nwt** (*str*) – name of nwt matrix in BigARTM
- **rwt** (*str*) – name of rwt matrix in BigARTM

**process\_batches** (*pwt*, *nwt*=None, *num\_document\_passes*=None, *batches\_folder*=None, *batches*=None, *regularizer\_name*=None, *regularizer\_tau*=None, *class\_ids*=None, *class\_weights*=None, *find\_theta*=False, *reuse\_theta*=False, *find\_ptdw*=False, *predict\_class\_id*=None)

#### Parameters

- **pwt** (*str*) – name of pwt matrix in BigARTM
- **nwt** (*str*) – name of nwt matrix in BigARTM
- **num\_document\_passes** (*int*) – number of inner iterations during processing
- **batches\_folder** (*str*) – full path to data folder (alternative 1)
- **batches** (*list of str*) – full file names of batches to process (alternative 2)
- **regularizer\_name** (*list of str*) – list of names of Theta regularizers to use
- **regularizer\_tau** (*list of float*) – list of tau coefficients for Theta regularizers
- **class\_ids** (*list of str*) – list of class ids to use during processing
- **class\_weights** (*list of float*) – list of corresponding weights of class ids
- **find\_theta** (*bool*) – find theta matrix for ‘batches’ (if alternative 2)
- **reuse\_theta** (*bool*) – initialize by theta from previous collection pass
- **find\_ptdw** (*bool*) – calculate and return Ptdw matrix or not (works if find\_theta == False)
- **predict\_class\_id** (*str, default None*) – class\_id of a target modality to predict

#### Returns

- tuple (messages.ThetaMatrix, numpy.ndarray) — the info about Theta (if find\_theta == True)
- messages.ThetaMatrix — the info about Theta (if find\_theta == False)

**reconfigure** (*topic\_names*=None, *class\_ids*=None, *scores*=None, *regularizers*=None, *num\_processors*=None, *pwt\_name*=None, *nwt\_name*=None, *num\_document\_passes*=None, *reuse\_theta*=None, *cache\_theta*=None)

**reconfigure\_regularizer** (*name*, *config*=None, *tau*=None, *gamma*=None)

**reconfigure\_score** (*name*, *config*)

**reconfigure\_topic\_name** (*topic\_names*=None)

**regularize\_model** (*pwt*, *nwt*, *rwt*, *regularizer\_name*, *regularizer\_tau*, *regularizer\_gamma*=None)

#### Parameters

- **pwt** (*str*) – name of pwt matrix in BigARTM
- **nwt** (*str*) – name of nwt matrix in BigARTM
- **rwt** (*str*) – name of rwt matrix in BigARTM

- **regularizer\_name** (*list of str*) – list of names of Phi regularizers to use
- **regularizer\_tau** (*list of double*) – list of tau coefficients for Phi regularizers

**transform** (*batches=None, batch\_filenames=None, theta\_matrix\_type=None, predict\_class\_id=None*)

**Parameters**

- **batches** – list of Batch instances
- **batch\_weights** (*list of float*) – weights of batches to transform
- **theta\_matrix\_type** (*int*) – type of matrix to be returned
- **predict\_class\_id** (*int*) – type of matrix to be returned

**Returns** messages.ThetaMatrix object

---

## Release Notes

---

### Changes in Python API

This page describes recent changes in BigARTM's Python API. Note that the API might be affected by changes in the underlying protobuf messages. For this reason we recommend to review [Changes in Protobuf Messages](#).

For further reference about Python API refer to [ARTM model](#), [Q & A](#) or [tutorials](#).

#### v0.8.2

**Warning:** BigARTM 3rdparty dependency had been upgraded from protobuf 2.6.1 to protobuf 3.0.0. This may affect you upgrade from previous version of bigartm. Please report any issues at [bigartm-users@googlegroups.com](mailto:bigartm-users@googlegroups.com).

**Warning:** BigARTM now require you to install `tqdm` library to visualize progress bars. To install use `pip install tqdm` or `conda install -c conda-forge tqdm`.

- Add support for python 3.0
- Add hARTM class to support hierarchy model
- Add HierarchySparsingTheta for advanced inference of hierarchical models
- Enable replacing regularizers in ARTM-like models:

```
# using operator[]-like style
model.regularizers['somename'] = SomeRegularizer(...)
# using keyword argument overwrite in add function
model.regularizers.add(SomeRegularizer(name='somename', ...), overwrite=True)
```

- Better error reporting: raise exception in `fit_offline`, `fit_online` and `transform` if there is no data to process)
- Better support for changes in topic names, with `reconfigure()`, `initialize()` and `merge_model()`
- Show progress bars in `fit_offline`, `fit_online` and `transform`.
- Add `ARTM.reshape_topics` method to add/remove/reorder topics.
- Add `max_dictionary_size` parameter to `Dictionary.filter()`
- Add `class_ids` parameter to `BatchVectorizer.__init__()`
- Add `dictionary_name` parameter to `MasterComponent.merge_model()`

- Add `ARTM.transform_sparse()` and `ARTM.get_theta_sparse()` for sparse retrieval of theta matrix
- Add `ARTM.get_phi_sparse()` for sparse retrieval of phi matrix

### v0.8.1

- New source type 'bow\_n\_wd' was added into `BatchVectorizer` class. This type oriented on using the output of `CountVectorizer` and `TfidfVectorizer` classes from `sklearn`. New parameters of `BatchVectorizer` are: `n_wd` (`numpy.array`) and `vocabulary`(`dict`)
- LDA model was added as one of the public interfaces. It is a restricted ARTM model created to simplify BigARTM usage for new users with few experience in topic modeling.
- `BatchVectorizer` got a flag 'gather\_dictionary', which has default value 'True'. This means that BV would create dictionary and save it in the `BV.dictionary` field. For 'bow\_n\_wd' format the dictionary will be gathered whenever the flag was set to 'False' or to 'True'.
- Add relative regularization for Phi matrix

### v0.8.0

**Warning:** Note that your script can be affected by our changes in the default values for `num_document_passes` and `reuse_theta` parameters (see below). We recommend to use our new default settings, `num_document_passes = 10` and `reuse_theta = False`. However, if you choose to explicitly set `num_document_passes = 1` then make sure to also set `reuse_theta = True`, otherwise you will experience very slow convergence.

- all operations to work with dictionaries were moved into a separate class `artm.Dictionary`. (details in [the documentation](#)). The mapping between old and new methods is very straightforward: `ARTM.gather_dictionary` is replaced with `Dictionary.gather` method, which allows to gather a dictionary from a set of batches; `ARTM.filter_dictionary` is replaced with `Dictionary.filter` method, which allows to filter a dictionary based on term frequency and document frequency; `ARTM.load_dictionary` is replaced with `Dictionary.load` method, which allows to load a dictionary previously exported to disk in `Dictionary.save` method; `ARTM.create_dictionary` is replaced with `Dictionary.create` method, which allows to create a dictionary based on custom protobuf message `DictionaryData`, containing a set of dictionary entries; etc... The following code snippet gives a basic example:

```
my_dictionary = artm.Dictionary()
my_dictionary.gather(data_path='my_collection_batches', vocab_file_path='vocab.txt')
my_dictionary.save(dictionary_path='my_collection_batches/my_dictionary')
my_dictionary.load(dictionary_path='my_collection_batches/my_dictionary.dict')
model = artm.ARTM(num_topics=20, dictionary=my_dictionary)
model.scores.add(artm.PerplexityScore(name='my_fisrt_perplexity_score',
                                     use_unigram_document_model=False,
                                     dictionary=my_dictionary))
```

- added `library_version` property to `ARTM` class to query for the version of the underlying BigARTM library; returns a string in MAJOR.MINOR.PATCH format;
- `dictionary_name` argument had been renamed to `dictionary` in many places across python interface, including scores and regularizers. This is done because those arguments can now except not just a string, but also the `artm.Dictionary` class itself.

- with `Dictionary` class users no longer have to generate names for their dictionaries (e.g. the unique `dictionary_name` identifier that references the dictionary). You may use `Dictionary.name` field to access to the underlying name of the dictionary.
- added `dictionary` argument to `ARTM.__init__` constructor to let user initialize the model; note that we've change the behavior that model is automatically initialized whenever user calls `fit_offline` or `fit_online`. Now this is no longer the case, and we expect user to either pass a dictionary in `ARTM.__init__` constructor, or manually call `ARTM.initialize` method. If neither is performed then `ARTM.fit_offline` and `ARTM.fit_online` will throw an exception.
- added `seed` argument to `ARTM.__init__` constructor to let user randomly initialize the model;
- added new score and score tracker `BackgroundTokensRatio`
- remove the default value from `num_topics` argument in `ARTM.__init__` constructor, which previously was defaulting to `num_topics = 10`; now user must always specify the desired number of topics;
- moved argument `reuse_theta` from `fit_offline` method into `ARTM.__init__` constructor; the argument is still used to indicate that the previous theta matrix should be re-used on the next pass over the collection; setting `reuse_theta = True` in the constructor will now be applied to `fit_online`, which previously did not have this option.
- moved common argument `num_document_passes` from `ARTM.fit_offline`, `ARTM.fit_online`, `ARTM.transform` methods into `ARTM.__init__` constructor.
- changed the default value of `cache_theta` parameter from `True` to `False` (in `ARTM.__init__` constructor); this is done to avoid excessive memory usage due to caching of the entire Theta matrix; if caching is indeed required user has to manually turn it on by setting `cache_theta = True`.
- changed the default value of `reuse_theta` parameter from `True` to `False` (in `ARTM.__init__` constructor); the reason is the same as for changing the default for `cache_theta` parameter
- changed the default value of `num_document_passes` parameter from 1 to 10 (in `ARTM.__init__` constructor);
- added arguments `apply_weight`, `decay_weight` and `update_after` in `ARTM.fit_online` method; each argument accepts a list of floats; setting all three arguments will override the default behavior of the online algorithm that rely on a specific formula with  $\tau_0$ ,  $\kappa$  and `update_every`.
- added argument `async` (boolean flag) in `ARTM.fit_online` method for improved performance.
- added argument `theta_matrix_type` in `ARTM.transform` method; potential values are: `"dense_theta"`, `"dense_ptdw"`, `None`; default matrix type is `"dense_theta"`.
- introduced a separate method `ARTM.remove_theta` to clear cached theta matrix; remove corresponding boolean switch `remove_theta` from `ARTM.get_theta` method.
- removed `ARTM.fit_transform` method; note that the name was confusing because this method has never fitted the model; the purpose of `ARTM.fit_transform` was to retrieve Theta matrix after fitting the model (`ARTM.fit_offline` or `ARTM.fit_online`); same functionality is now available via `ARTM.get_theta` method.
- introduced `ARTM.get_score` method, which will exist in parallel to score tracking functionality; the goal for `ARTM.get_score(score_name)` is to always return the latest version of the score; for Phi scores this means to calculate them on fly; for Theta scores this means to return a score aggregated over last call to `ARTM.fit_offline`, `ARTM.fit_online` or `ARTM.transform` methods; opposite to `ARTM.get_score` the score tracking functionality returns the overall history of a score. For further details on score calculation refer to [Q&A section](#) in our wiki page.
- added `data_weight` in `BatchVectorizer.__init__` constructor to let user specify an individual weight for each batch

- score tracker classes had been rewritten, so you should make minor changes in the code that retrieves scores; for example:
- added an API to initialize logging with custom logging directory, log level, etc... Search out wiki page [Q&A](#) for more details.

```
# in v0.7.x
print model.score_tracker['Top100Tokens'].last_topic_info[topic_name].tokens

# in v0.8.0
last_tokens = model.score_tracker['Top100Tokens'].last_tokens
print last_tokens[topic_name]
```

### v0.7.x

See [BigARTM v0.7.X Release Notes](#).

## Changes in Protobuf Messages

### v0.8.2

- added `CollectionParserConfig.num_threads` to control the number of threads that perform parsing. At the moment the feature is only implemented for VW-format.
- added `CollectionParserConfig.class_id` (repeated string) to control which modalities should be parsed. If token's `class_id` is not from this list, it will be excluded from the resulting batches. When the list is empty, all modalities are included (this is the default behavior, as before).
- added `CollectionParserInfo` message to export diagnostics information from `ArtmParseCollection`
- added `FilterDictionaryArgs.max_dictionary_size` to give user an easy option to limit his dictionary size
- added `MergeModelArgs.dictionary_name` to define the set of tokens in the resulting matrix
- added `ThetaMatrix.num_values`, `TopicModel.num_values` to define number of non-zero elements in sparse format

### v0.8.0

**Warning:** New batches, created in *BigARTM v0.8*, **CAN NOT** be used in the previous versions of the library. Old batches, created prior to *BigARTM v0.8*, can still be used. See below for details.

- added `token_id` and `token_weight` field in `Item` message, and obsoleted `Item.field`. Internally the library will merge the content of `Field.token_id` and `Field.token_weight` across all fields, and store the result back into `Item.token_id`, `Item.token_weight`. New `Item` message is as follows:

```
message Item {
  optional int32 id = 1;
  repeated Field field = 2;  // obsolete in BigARTM v0.8.0
  optional string title = 3;
  repeated int32 token_id = 4;
```

```
repeated float token_weight = 5;
}
```

- renamed `topics_count` into `num_topics` across multiple messages (TopicModel, ThetaMatrix, etc)
- renamed `inner_iterations_count` into `num_document_passes` in `ProcessBatchesArgs`
- renamed `passes` into `num_collection_passes` in `FitOfflineMasterModelArgs`
- renamed `threads` into `num_processors` in `MasterModelConfig`
- renamed `topic_index` field into `topic_indices` in `TopicModel` and `ThetaMatrix` messages
- added messages `ScoreArray`, `GetScoreArrayArgs` and `ClearScoreArrayCacheArgs` to bring score tracking functionality down into BigARTM core
- added messages `BackgroundTokensRatioConfig` and `BackgroundTokensRatio` (new score)
- moved `model_name` from `GetScoreValueArgs` into `ScoreConfig`; this is done to support score tracking functionality in BigARTM core; each Phi score needs to know which model to use in calculation
- removed `topics_count` from `InitializeModelArgs`; users must specify topic names in `InitializeModelArgs.topic_name` field
- removed `topic_index` from `GetThetaMatrixArgs`; users must specify topic names to retrieve in `GetThetaMatrixArgs.topic_name`
- removed `batch` field in `GetThetaMatrixArgs` and `GetScoreValueArgs.batch` messages; users should use `ArtmRequestTransformMasterModel` or `ArtmRequestProcessBatches` to process new batches and calculate theta scores
- removed `reset_scores` flag in `ProcessBatchesArgs`; users should use new API `ArtmClearScoreCache`
- removed `clean_cache` flag in `GetThetaMatrixArgs`; users should use new API `ArtmClearThetaCache`
- removed `MasterComponentConfig`; users should use `ArtmCreateMasterModel` and pass `MasterModelConfig`
- removed obsolete fields in `CollectionParserConfig`; same arguments can be specified at `GatherDictionaryArgs` and passed to `ArtmGatherDictionary`
- removed `Filter` message in `InitializeModelArgs`; same arguments can be specified at `FilterDictionaryArgs` and passed to `ArtmFilterDictionary`
- removed `batch_name` from `ImportBatchesArgs`; the field is no longer needed; batches will be identified via their `Batch.id` identifier
- removed `use_v06_api` in `MasterModelConfig`
- removed `ModelConfig` message
- removed `SynchronizeModelArgs`, `AddBatchArgs`, `InvokeIterationArgs`, `WaitIdleArgs` messages; users should use new APIs based on `MasterModel`
- removed `GetRegularizerStateArgs`, `RegularizerInternalState`, `MultiLanguagePhiInternalState` messages
- removed `model_name` and `model_name_cache` in `ThetaMatrix`, `GetThetaMatrixArgs` and `ProcessBatchesArgs`; the code of master component is simplified to only handle one theta matrix, so there is no longer any reason to identify theta matrix with `model_name`

- removed `Stream` message, `MasterComponentConfig.stream` field, and all `stream_name` fields across several messages; train/test streaming functionality is fully removed; users are expected to manage their train and test collections (for example as separate folders with batches)
- removed `use_sparse_bow` field in several messages; the computation mode with dense matrices is no longer supported;
- renamed `item_count` into `num_items` in `ThetaSnippetScoreConfig`
- add global enum `ScoreType` as a replacement for enums `Type` from `ScoreConfig` and `ScoreData` messages
- add global enum `RegularizerType` as a replacement for enum `Type` from `RegularizerConfig` message
- add global enum `MatrixLayout` as a replacement for enum `MatrixLayout` from `GetThetaMatrixArgs` and `GetTopicModelArgs` messages
- add global enum `ThetaMatrixType` as a replacement for enum `ThetaMatrixType` from `ProcessBatchesArgs` and `TransformMasterModelArgs` messages
- renamed enum `Type` into `SmoothType` in `SmoothPtdwConfig` to avoid conflicts in C# messages
- renamed enum `Mode` into `SparseMode` in `SpecifiedSparsePhiConfig` to avoid conflicts in C# messages
- renamed enum `Format` into `CollectionFormat` in `CollectionParserConfig` to avoid conflicts in C# messages
- renamed enum `NameType` into `BatchNameType` in `CollectionParserConfig` to avoid conflicts in C# messages
- renamed field `transform_type` into `type` in `TransformConfig` to avoid conflicts in C# messages
- remove message `CopyRequestResultArgs`; this is a breaking change; please check that
  - all previous calls to `ArtmCopyRequestResult` are changed to `ArtmCopyRequestedMessage`
  - all previous calls to `ArtmCopyRequestResultEx` with request types `GetThetaSecondPass` and `GetModelSecondPass` are changed to `ArtmCopyRequestedObject`
  - all previous calls to `ArtmCopyRequestResultEx` with `DefaultRequestType` are changed to `ArtmCopyRequestedMessage`
- remove field `request_type` in `GetTopicModelArgs`; to request only topics and/or tokens users should set `GetTopicModelArgs.matrix_layout` to `MatrixLayout_Sparse`, and `GetTopicModelArgs.eps` = 1.001 (any number greather that 1.0).
- change optional `FloatArray` into repeated float in field `coherence` of `TopTokensScore`
- change optional `DoubleArray` into repeated double in fields `kernel_size`, `kernel_purity`, `kernel_contrast` and `coherence` of `TopicKernelScore`
- change optional `StringArray` into repeated string in field `topic_name` of `TopicKernelScore`

## v0.7.x

See [BigARTM v0.7.X Release Notes](#).



## Changes in BigARTM CLI

### v0.8.2

- added option `--rand-seed` to initialize random number generator; without this options, RNG will be set using system time
- added option `--write-vw-corpus` to convert batches into plain text file in Vowpal Wabbit format
- change the naming scheme of the batches, saved with `--save-batches` option. Previously file names were guid-based, while new format will look like this: `abcde.batch`. New format ensures the ordering of the documents in the collection is be preserved, given that user scans batches alphabetically.
- added switch `--guid-batch-name` to enable old naming scheme of batches (guid-based names). This option is useful if you launch multiple instances of BigARTM CLI to concurrently generate batches.
- speedup parsing large files in VowpalWabbit format
- when `--use-modality` is specified, the batches saved with `--save-batches` will only include tokens from these modalities. Other tokens will be ignored during parsing. This option is implemented for both VW and UCI BOW formats.
- implement `TopicSelection`, `LabelRegularization`, `ImproveCoherence`, `Biterms` regularizer in BigARTM CLI
- added option `--dictionary-size` to give user an easy option to limit his dictionary size
- add more diagnostics information about dictionary size (before and after filtering)
- add strict verification of scores and regularizers; for example, BigARTM CLI will raise an exception for this input: `bigartm -t obj:10,back:5 --regularizer "0.5 SparsePhi #obj*"`. There shouldn't be star sign in `#obj*`.

### v0.8.0

- renamed `--passes` into `--num-collection-passes`
- renamed `--num-inner-iterations` into `--num-document-passes`
- removed `--model-v06` option
- removed `--use-dense-bow` option

### v0.7.x

See [BigARTM v0.7.X Release Notes](#).

## Changes in c\_interface

### v0.8.2

**Warning:** BigARTM 3rdparty dependency had been upgraded from `protobuf 2.6.1` to `protobuf 3.0.0`. This may affect you upgrade from previous version of bigartm. Please report any issues at [bigartm-users@googlegroups.com](mailto:bigartm-users@googlegroups.com).

- Change `ArtmParseCollection` to return `CollectionParserInfo` message
- Add APIs to enable JSON serialization for all input and output protobuf messages
  - `ArtmSetProtobufMessageFormatToJson()`
  - `ArtmSetProtobufMessageFormatToBinary()`
  - `ArtmProtobufMessageFormatIsJson()`

The default setting is, as before, to serialize all message into binary buffer. Note that for with json serialization one should use `RegularizerConfig.config_json`, `ScoreConfig.config_json` and `ScoreData.data_json` instead of `RegularizerConfig.config`, `ScoreConfig.config` and `ScoreData.data`.

- Revisit documentation for `c_interface`
- Change integer types in `c_interface` from `int` to `int64_t` (from `stdint.h`). This allows to validate 2 GB limit for protobuf messages, and also to passing larger objects in `ArtmCopyRequestedObject`.
- Add `ArtmReconfigureTopicName` method to add/remove/reorder topic names
- Support sparse format for external retrieval of theta and phi matrices

## v0.8.0

- Removed `ArtmCreateMasterComponent` and `ArtmReconfigureMasterComponent`
- Removed `ArtmCreateModel` and `ArtmReconfigureModel`
- Removed `ArtmAddBatch`, `ArtmInvokeIteration`, `ArtmWaitIdle`, `ArtmSynchronizeModel`
- Removed `ArtmRequestRegularizerState`
- Renamed `ArtmCopyRequestResult` into `ArtmCopyRequestedMessage`
- Renamed `ArtmCopyRequestResultEx` into `ArtmCopyRequestedObject`
- Added `ArtmClearThetaCache` and `ArtmClearScoreCache`
- Added `ArtmRequestScoreArray` and `ArtmClearScoreArrayCache`
- Added `GetArtmVersion` to query for the version; returns a string in “<MAJOR>.<MINOR>.<PATCH>” format

## v0.7.x

See [BigARTM v0.7.X Release Notes](#).

# BigARTM v0.7.X Release Notes

## BigARTM v0.7.0 Release notes

We are happy to introduce BigARTM v0.7.0, which brings you the following changes:

- New-style models
- Network modulus operandi is removed
- Coherence regularizer and scores (experimental)

## New-style models

BigARTM v0.7.0 exposes new APIs to give you additional control over topic model inference:

- `ProcessBatches`
- `MergeModel`
- `RegularizeModel`
- `NormalizeModel`

Besides being more flexible, new APIs bring many additional benefits:

- Fully deterministic inference, no dependency on threads scheduling or random numbers generation
- Less bottlenecks for performance (`DataLoader` and `Merger` threads are removed)
- Phi-matrix regularizers can be implemented externally
- Capability to output Phi matrices directly into your NumPy matrices (scheduled for BigARTM v0.7.2)
- Capability for store Phi matrices in sparse format (scheduled for BigARTM v0.7.3)
- Capability for async `ProcessBatches` and non-blocking online algorithm (BigARTM v0.7.4)
- Form solid foundation for high performance networking (BigARTM v0.8.X)

The picture below illustrates scalability of BigARTM v0.7.0 vs v0.6.4. Top chart (in green) corresponds to CPU usage at 28 cores on machine with 32 virtual cores (16 physical cores + hyper threading). As you see, new version is much more stable. In addition, new version consumes less memory.



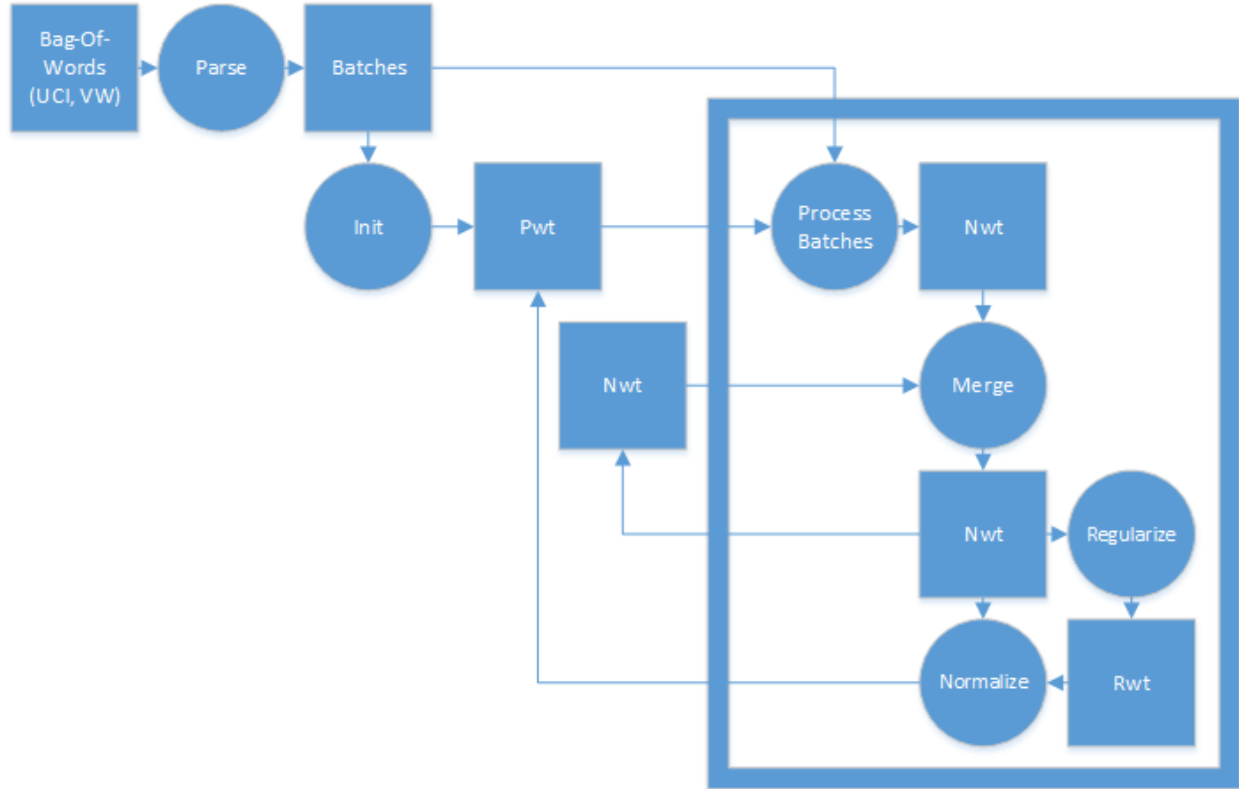
Refer to the following examples that demonstrate usage of new APIs for offline, online and regularized topic modelling:

- [example17\\_process\\_batches.py](#)
- [example18\\_merge\\_model.py](#)
- [example19\\_regularize\\_model.py](#)

Models, tuned with the new API are referred to as *new-style models*, as opposite to *old-style models* inferred with `AddBatch`, `InvokeIteration`, `WaitIdle` and `SynchronizeModel` APIs.

**Warning:** For BigARTM v0.7.X we will continue to support old-style models. However, you should consider upgrading to new-style models because old APIs (AddBatch, InvokeIteration, WaitIdle and SynchronizeModel) are likely to be removed in future releases.

The following flow chart gives a typical use-case on new APIs in online regularized algorithm:



### Notes on upgrading existing code to new-style models

1. New APIs can only read batches from disk. If your current script passes batches via memory (in `AddBatchArgs.batch` field) then you need to store batches on disk first, and then process them with `ProcessBatches` method.
2. Initialize your model as follows:
  - For `python_interface`: using `MasterComponent.InitializeModel` method
  - For `cpp_interface`: using `MasterComponent.InitializeModel` method
  - For `c_interface`: using `ArtnInitializeModel` method

Remember that you should not create `ModelConfig` in order to use this methods. Pass your `topics_count` (or `topic_name` list) as arguments to `InitializeModel` method.

3. Learn the difference between Phi and Theta scores, as well as between Phi and Theta regularizes. The following table gives an overview:

Object	Theta	Phi
Scores	<ul style="list-style-type: none"> <li>• Perplexity</li> <li>• SparsityTheta</li> <li>• ThetaSnippet</li> <li>• ItemsProcessed</li> </ul>	<ul style="list-style-type: none"> <li>• SparsityPhi</li> <li>• TopTokens</li> <li>• TopicKernel</li> </ul>
Regularizers	<ul style="list-style-type: none"> <li>• SmoothSparseTheta</li> </ul>	<ul style="list-style-type: none"> <li>• DecorrelatorPhi</li> <li>• ImproveCoherencePhi</li> <li>• LabelRegularizationPhi</li> <li>• SmoothSparsePhi</li> <li>• SpecifiedSparsePhi</li> </ul>

Phi regularizers needs to be calculated explicitly in `RegularizeModel`, and then applied in `NormalizeModel` (via optional `rwf` argument). Theta regularizers needs to be enabled in `ProcessBatchesArgs`. Then they will be automatically calculated and applied during `ProcessBatches`.

Phi scores can be calculated at any moment based on the new-style model (same as for old-style models). Theta scores can be retrieved in two equivalent ways:

```
pwt_model = "pwt"
master.ProcessBatches(pwt_model, batches, "nwt")
perplexity_score.GetValue(pwt_model).value
```

or

```
pwt_model = "pwt"
process_batches_result = master.ProcessBatches(pwt_model, batches, "nwt")
perplexity_score.GetValue(scores = process_batches_result).value
```

Second way is more explicit. However, the first way allows you to combine aggregate scores accross multiple `ProcessBatches` calls:

```
pwt_model = "pwt"
master.ProcessBatches(pwt_model, batches1, "nwt")
master.ProcessBatches(pwt_model, batches2, "nwt", reset_scores=False)
perplexity_score.GetValue(pwt_model).value
```

This works because BigARTM caches the result of `ProcessBatches` together (in association with `pwt_model`). The `reset_scores` switch disables the default behaviour, which is to reset the cache for `pwt_model` at the beginning of each `ProcessBatch` call.

4. Continue using `GetThetaMatrix` and `GetTopicModel` to retrieve results from the library. For `GetThetaMatrix` to work you still need to enable `cache_theta` in master component. Remember to use the same model in `GetThetaMatrix` as you used as the input to `ProcessBatches`. You may also omit “target\_nwt” argument in `ProcessBatches` if you are not interested in getting this output.

```
master.ProcessBatches("pwt", batches)
theta_matrix = master.GetThetaMatrix("pwt")
```

5. Stop using certain APIs:

- For `python_interface`: stop using class `Model` and `ModelConfig` message
- For `cpp_interface`: stop using class `Model` and `ModelConfig` message
- For `c_interface`: stop using methods `ArtmCreateModel`, `ArtmReconfigureModel`, `ArtmInvokeIteration`, `ArtmAddBatch`, `ArtmWaitIdle`, `ArtmSynchronizeModel`

**Notes on models handling (reusing, sharing input and output, etc)**

Is allowed to output the result of ProcessBatches, NormalizeModel, RegularizeModel and MergeModel into an existing model. In this case the existing model will be fully overwritten by the result of the operation. For all operations except ProcessBatches it is also allowed to use the same model in inputs and as an output. For example, typical usage of MergeModel involves combining “nwt” and “nwt\_hat” back into “nwt”. This scenario is fully supported. The output and input of ProcessBatches must refer to two different models. Finally, note that MergeModel will ignore all non-existing models in the input (and log a warning). However, if none of the input models exist then MergeModel will throw an error.

**Known differences**

1. Decorrelator regularizer will give slightly different result in the following scenario:

```
master.ProcessBatches("pwt", batches, "nwt")
master.RegularizeModel("pwt", "nwt", "rwt", phi_regularizers)
master.NormalizeModel("nwt", "pwt", "rwt")
```

To get the same result as from model.Synchronize() adjust your script as follows:

```
master.ProcessBatches("pwt", batches, "nwt")
master.NormalizeModel("nwt", "pwt_temp")
master.RegularizeModel("pwt_temp", "nwt", "rwt", phi_regularizers)
master.NormalizeModel("nwt", "pwt", "rwt")
```

2. You may use GetThetaMatrix(pwt) to retrieve Theta-matrix, previously calculated for new-style models inside ProcessBatches. However, you can not use GetThetaMatrix(pwt, batch) for new models. They do not have corresponding ModelConfig, and as a result you need to go through ProcessBatches to pass all parameters.

**Network modus operandi is removed**

Network modus operandi had been removed from BigARTM v0.7.0.

This decision had been taken because current implementation struggle from many issues, particularly from poor performance and stability. We expect to re-implement this functionality on top of new-style models.

Please, let us know if this caused issues for you, and we will consider to re-introduce networking in v0.8.0.

**Coherence regularizer and scores (experimental)**

Refer to example in [example16\\_coherence\\_score.py](#).

**BigARTM v0.7.1 Release notes**

We are happy to introduce BigARTM v0.7.1, which brings you the following changes:

- BigARTM notebooks — new source of information about BigARTM
- ArtmModel — a brand new Python API
- Much faster retrieval of Phi and Theta matrices from Python
- Much faster dictionary imports from Python
- Auto-detect and use all CPU cores by default
- Fixed Import/Export of topic models (was broken in v0.7.0)
- New capability to implement Phi-regularizers in Python code

- Improvements in Coherence score

Before you upgrade to BigARTM v0.7.1 please review the changes that *break backward compatibility*.

### BigARTM notebooks

BigARTM notebooks is your go-to links to read more ideas, examples and other information around BigARTM:

- [BigARTM notebooks in English](#)
- [BigARTM notebooks in Russian](#)

### ArtmModel

Best thing about ArtmModel is that this API had been designed by BigARTM users. Not by BigARTM programmers. This means that BigARTM finally has a nice, clean and easy-to-use programming interface for Python. Don't believe it? Just take a look and some examples:

- [ArtmModel examples in English](#)
- [ArtmModel examples in Russian](#)

That is cool, right? This new API allows you to load input data from several file formats, infer topic model, find topic distribution for new documents, visualize scores, apply regularizers, and perform many other actions. Each action typically takes one line to write, which allows you to work with BigARTM interactively from Python command line.

ArtmModel exposes most of BigARTM functionality, and it should be sufficiently powerful to cover 95% of all BigARTM use-cases. However, for the most advanced scenarios you might still need to go through the previous API ([artm.library](#)). When in doubt which API to use, ask [bigartm-users@googlegroups.com](mailto:bigartm-users@googlegroups.com) — we are there to help!

### Coding Phi-regularizers in Python code

This is of course one of those very advanced scenarios where you need to go down to the old API :) Take a look at this example:

- [example19\\_regularize\\_model](#)
- [example20\\_attach\\_model](#)

First one tells how to use Phi regularizers, built into BigARTM. Second one provides a new capability to manipulate Phi matrix from Python. We call this **Attach** numpy matrix to the model, because this is similar to attaching debugger (like gdb or Visual Studio) to a running application.

To implement your own Phi regularizer in Python you need to to **attach** to `rwt` model from the first example, and update its values.

### Other changes

**Fast retrieval of Phi and Theta matrices.** In BigARTM v0.7.1 dense Phi and Theta matrices will be retrieved to Python as numpy matrices. All copying work will be done in native C++ code. This is much faster comparing to current solution, where all data is transferred in a large Protobuf message which needs to be deserialized in Python. ArtmModel already takes advantage of this performance improvements.

**Fast dictionary import.** BigARTM core now supports importing dictionary files from disk, so you no longer have to load them to Python. ArtmModel already take advantage of this performance improvement.



**Auto-detect number of CPU cores.** You no longer need to specify `num_processors` parameter. By default BigARTM will detect the number of cores on your machine and load all of them. `num_processors` still can be used to limit CPU resources used by BigARTM.

**Fixed Import/Export of topic models.** Export and Import of topic models will now work. As simple as this:

```
master.ExportModel("pwt", "file_on_disk.model")
master.ImportModel("pwt", "file_on_disk.model")
```

This will also take care of very large models above 1 GB that does not fit into single protobuf message.

**Coherence scores.** Ask [bigartm-users@googlegroups.com](mailto:bigartm-users@googlegroups.com) if you are interested :)

## Breaking changes

- **Changes in Python methods** `MasterComponent.GetTopicModel` and `MasterComponent.GetThetaMatrix`

From BigARTM v0.7.1 and onwards method `MasterComponent.GetTopicModel` of the low-level Python API will return a tuple, where first argument is of type `TopicModel` (protobuf message), and second argument is a numpy matrix. `TopicModel` message will keep all fields as usual, except `token_weights` field which will become empty. Information from `token_weights` field had been moved to numpy matrix (rows = tokens, columns = topics).

Similarly, `MasterComponent.GetThetaMatrix` will also return a tuple, where first argument is of type `ThetaMatrix` (protobuf message), and second argument is a numpy matrix. `ThetaMatrix` message will keep all fields as usual, except `item_weights` field which will become empty. Information from `item_weights` field had been moved to numpy matrix (rows = items, columns = topics).

Updated examples:

- [example11\\_get\\_theta\\_matrix.py](#)
- [example12\\_get\\_topic\\_model](#)

**Warning:** Use the followign syntax to restore the old behaviour:

- `MasterComponent.GetTopicModel(use_matrix = False)`
- `MasterComponent.GetThetaMatrix(use_matrix = False)`

This will return a complete protobuf message, without numpy matrix.

- **Python method `ParseCollectionOrLoadDictionary` is now obsolete**
  - Use `ParseCollection` method to convert collection into a set of batches
  - Use `MasterComponent.ImportDictionary` to load dictionary into BigARTM
  - Updated example: [example06\\_use\\_dictionaries.py](#)

## BigARTM v0.7.2 Release notes

We are happy to introduce BigARTM v0.7.2, which brings you the following changes:

- Enhancements in high-level python API (`ArtemModel` -> `ARTM`)
- Enhancements in low-level python API (`library.py` -> `master_component.py`)
- Enhancements in CLI interface (`cpp_client`)
- Status and information retrievals from BigARTM

- Allow float token counts (`token_count` -> `token_weight`)
- Allow custom weights for each batch (`ProcessBatchesArgs.batch_weight`)
- Bug fixes and cleanup in the online documentation

### Enhancements in Python APIs

Note that `ArtmModel` had been renamed to `ARTM`. The naming conventions follow the same pattern as in [scikit learn](#) (e.g. `fit`, `transform` and `fit_transform` methods).

Also note that all input data is now handled by `BatchVectorizer` class.

Refer to notebooks in [English](#) and in [Russian](#) for further details about ARTM interface.

Also note that previous low-level python API `library.py` is superseded by a new API `master_component.py`. For now both APIs are available, but the old one will be removed in future releases. Refer to [this folder](#) for further examples of the new low-level python API.

Remember that any use of low-level APIs is discouraged. Our recommendation is to always use the high-level python API `ARTM`, and e-mail us know if some functionality is not exposed there.

### Enhancements in CLI interface

BigARTM command line interface `cpp_client` had been enhanced with the following options:

- `--load_model` - to load model from file before processing
- `--save_model` - to save the model to binary file after processing
- `--write_model_readable` - to output the model in a human-readable format (CSV)
- `--write_predictions` - to write prediction in a human-readable format (CSV)
- `--dictionary_min_df` - to filter out tokens present in less than N documents / less than P% of documents
- `--dictionary_max_df` - filter out tokens present in less than N documents / less than P% of documents
- `--tau0` - an option of the online algorithm, describing the weight parameter in the online update formula. Optional, defaults to 1024.
- `--kappa` - an option of the online algorithm, describing the exponent parameter in the online update formula. Optional, defaults to 0.7.

Note that for `--dictionary_min_df` and `--dictionary_max_df` can be treated as number, fraction, percent.

- Use a percentage % sign to specify percentage value
- Use a floating value in `[0, 1)` range to specify a fraction
- Use an integer value (1 or greater) to indicate a number

### BigARTM v0.7.3 Release notes

BigARTM v0.7.3 releases the following changes:

- New command line tool for BigARTM
- Support for classification in `bigartm` CLI
- Support for asynchronous processing of batches
- Improvements in coherence regularizer and coherence score

- New *TopicMass* score for phi matrix
- Support for documents markup
- New API for importing batches through memory

## New command line tool for BigARTM

New CLI is named `bigartm` (or `bigrtm.exe` on Windows), and it supersedes previous CLI named `cpp_client`. New CLI has the following features:

- Parse collection in one of the [Formats](#)
- Load dictionary
- Initialize a new model, or import previously created model
- Perform EM-iterations to fit the model
- Export predicted probabilities for all documents into CSV file
- Export model into a file

All command-line options are listed [here](#), and you may see several examples on [BigARTM](#) page at github. At the moment full documentation is only available in [Russian](#).

## Support for classification in BigARTM CLI

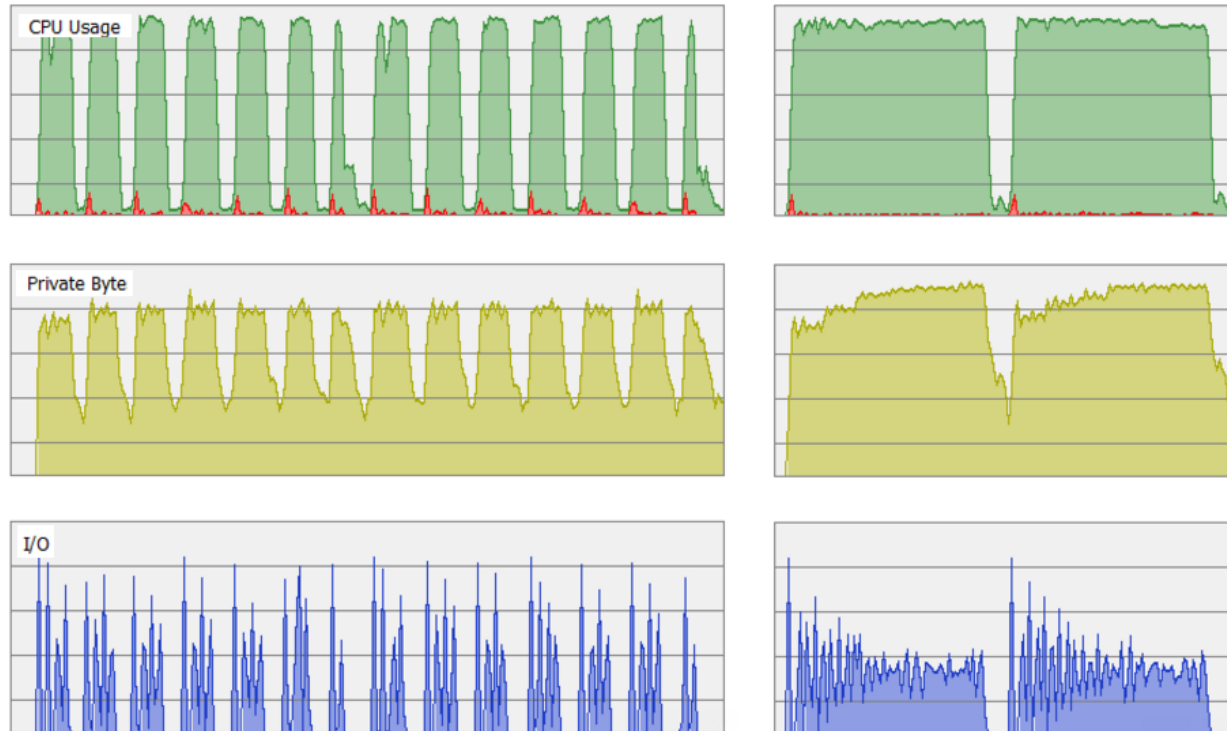
BigARTM CLI is now able to perform classification. The following example assumes that your batches have `target_class` modality in addition to the default modality (`@default_class`).

```
# Fit model
bigartm.exe --use-batches <your batches>
             --use-modality @default_class,target_class
             --topics 50
             --dictionary-min-df 10
             --dictionary-max-df 25%
             --save-model model.bin

# Apply model and output to text files
bigartm.exe --use-batches <your batches>
             --use-modality @default_class,target_class
             --topics 50
             --passes 0
             --load-model model.bin
             --predict-class target_class
             --write-predictions pred.txt
             --write-class-predictions pred_class.txt
             --csv-separator=tab
             --score ClassPrecision
```

## Support for asynchronous processing of batches

Asynchronous processing of batches enables applications to overlap EM-iterations better utilize CPU resources. The following chart shows CPU utilization of `bigartm.exe` with (left-hand side) and without async flag (right-hand side).



### TopicMass score for phi matrix

Topic mass score calculates cumulated topic mass for each topic. This is a useful metric to monitor balance between topics.

### Support for documents markup

Document markup provides topic distribution for each word in a document. Since BigARTM v0.7.3 it is possible to extract this information to use it. A potential application includes color-highlighted maps of the document, where every word is colored according to the most probable topic of the document.

In the code this feature is referred to as `ptdw` matrix. It is possible to extract and regularizer `ptdw` matrices. In future versions it will be also possible to calculate scores based on `ptdw` matrix.

### New API for importing batches through memory

New low-level APIs `ArtmImportBatches` and `ArtmDisposeBatches` allow to import batches from memory into BigARTM. Those batches are saved in BigARTM, and can be used for batches processing.

## BigARTM v0.7.4 Release notes

BigARTM v0.7.4 is a big release that includes major rework of dictionaries and [MasterModel](#).

### *bigartm/stable* branch

Up until now BigARTM has only one `master` branch, containing the latest code. This branch potentially includes untested code and unfinished features. We are now introducing `bigartm/stable` branch, and encourage all users

to stop using `master` and start fetching from `stable`. `stable` branch will be lagging behind `master`, and moved forward to `master` as soon as maintainers decide that it is ready. At the same point we will introduce a new tag (something like `v0.7.3`) and produce a new release for Windows. In addition, `stable` branch also might receive small urgent fixes in between releases, typically to address critical issues reported by our users. Such fixes will be also included in `master` branch.

## MasterModel

`MasterModel` is a new set of low-level APIs that allow users of C-interface to infer models and apply them to new data. The APIs are `ArtmCreateMasterModel`, `ArtmReconfigureMasterModel`, `ArtmFitOfflineMasterModel`, `ArtmFitOnlineMasterModel` and `ArtmRequestTransformMasterModel`, together with corresponding protobuf messages. For a usage example see `src/bigartm/srcmain.cc`.

This APIs should be easy to understand for the users who are familiar with Python interface. Basically, we take ARTM class in Python, and push it down to the core. Now users can create their model via `MasterModelConfig` (protobuf message), fit via `ArtmFitOfflineMasterModel` or `ArtmFitOnlineMasterModel`, and apply to the new data via `ArtmRequestTransformMasterModel`. This means that the user no longer has to orchestrate low-level building blocks such as `ArtmProcessBatches`, `ArtmMergeModel`, `ArtmRegularizeModel` and `ArtmNormalizeModel`.

`ArtmCreateMasterModel` is similar to `ArtmCreateMasterComponent` in a sence that it returns `master_id`, which can be later passed to all other APIs. This mean that most APIs will continue working as before. This applies to `ArtmRequestThetaMatrix`, `ArtmRequestTopicModel`, `ArtmRequestScore`, and many others.

## Rework of dictionaries

Previous implementation of the dictionaries was really messy, and we are trying to clean this up. This effort is not finished yet, however we decided to release current version because it is a major improvement comparing to the previous version. At the low-level (`c_interface`), we now have the following methods to work with dictionaries:

- `ArtmGatherDictionary` collects a dictionary based on a folder with batches,
- `ArtmFilterDictionary` filter tokens from the dictinoary based on their term frequency or document frequency,
- `ArtmCreateDictionary` creates a dictionary from a custom `DictionaryData` object (protobuf message),
- `ArtmRequestDictionary` retrieves a dictionary as `DictionaryData` object (protobuf message),
- `ArtmDisposeDictionary` deletes dictionary object from BigARTM,
- `ArtmImportDictionary` import dictionary from binary file,
- `ArtmExportDictionary` expor tdictionary into binary file.

All dictionaries are identified by a string ID (`dictionary_name`). Dictionaries can be used to initialize the model, in regularizers or in scores.

Note that `ArtmImportDictionary` and `ArtmExportDictionary` now uses a different format. For this reason we require that all imported or exported files end with `.dict` extension. This limitation is only introduced to make users aware of the change in binary format.

**Warning:** Please note that you have to re-generate all dictionaries, created in previous BigARTM versions. To force this limitation we decided that `ArtmImportDictionary` and `ArtmExportDictionary` will require all imported or exported files end with `.dict` extension. This limitation is only introduced to make users aware of the change in binary format.

Please note that in the next version (*BigARTM v0.8.0*) we are planing to break dictionary format once again. This is because we will introduce `boost.serialize` library for all import and export methods. From that point `boost.serialize` library will allow us to upgrade formats without breaking backwards compatibility.

The following example illustrate how to work with new dictionaries from Python.

```
# Parse collection in UCI format from D:\Datasets\docword.kos.txt and D:\Datasets\vocab.kos.txt
# and store the resulting batches into D:\Datasets\kos_batches
batch_vectorizer = artm.BatchVectorizer(data_format='bow_uci',
                                       data_path=r'D:\Datasets',
                                       collection_name='kos',
                                       target_folder=r'D:\Datasets\kos_batches')

# Initialize the model. For now dictionaries exist within the model,
# but we will address this in the future.
model = artm.ARTM(...)

# Gather dictionary named `dict` from batches.
# The resulting dictionary will contain all distinct tokens that occur
# in those batches, and their term frequencies
model.gather_dictionary("dict", "D:\Datasets\kos_batches")

# Filter dictionary by removing tokens with too high or too low term frequency
# Save the result as `filtered_dict`
model.filter_dictionary(dictionary_name='dict',
                       dictionary_target_name='filtered_dict',
                       min_df=10, max_df_rate=0.4)

# Initialize model from `diltered_dict`
model.initialize("filtered_dict")

# Import/export functionality
model.save_dictionary("filtered_dict", "D:\Datasets\kos.dict")
model.load_dictionary("filtered_dict2", "D:\Datasets\kos.dict")
```

### Changes in the infrastructure

- Static linkage for bigartm command-line executable on Linux. To disable static linkage use `cmake -DBUILD_STATIC_BIGARTM=OFF ..`
- Install BigARTM python API via `python setup.py install`

### Changes in core functionality

- Custom transform function for KL-div regularizers
- Ability to initialize the model with custom seed
- TopicSelection regularizers

- `PeakMemory` score (Windows only)
- Different options to name batches when parsing collection (`GUID` as today, and `CODE` for sequential numbering)

### Changes in Python API

- `ARTM.dispose()` method for managing native memory
- `ARTM.get_info()` method to retrieve internal state
- Performance fixes
- Expose class prediction functionality

### Changes in C++ interface

- Consume `MasterModel` APIs in C++ interface. Going forward this is the only C++ interface that we will support.

### Changes in console interface

- Better options to work with dictionaries
- `--write-dictionary-readable` to export dictionary
- `--force` switch to let user overwrite existing files
- `--help` generates much better examples
- `--model-v06` to experiment with old APIs (`ArtmInvokeIteration` / `ArtmWaitIdle` / `ArtmSynchronizeModel`)
- `--write-scores` switch to export scores into file
- `--time-limit` option to time-box model inference(as an alternative to `--passes` switch)





---

## BigARTM Developer's Guide

---

These pages describe the development process of BigARTM library. If your intent to use BigARTM as a typical user, please proceed to [Basic BigARTM tutorial for Windows users](#) or [Basic BigARTM tutorial for Linux and Mac OS-X users](#), depending on your operating system. If you intent is to contribute to the development BigARTM, please proceed to the links below.

### Downloads (Windows)

Download and install the following tools:

- **Git for Windows from <http://git-scm.com/download/win>**
  - <https://github.com/msysgit/msysgit/releases/download/Git-1.9.5-preview20141217/Git-1.9.5-preview20141217.exe>
- **Github for Windows from <https://windows.github.com/>**
  - <https://github-windows.s3.amazonaws.com/GitHubSetup.exe>
- Visual Studio 2013 Express for Windows Desktop from <https://www.visualstudio.com/en-us/products/visual-studio-express-vs.aspx>
- **CMake from <http://www.cmake.org/download/>**
  - <http://www.cmake.org/files/v3.0/cmake-3.0.2-win32-x86.exe>
- **Prebuilt Boost binaries from <http://sourceforge.net/projects/boost/files/boost-binaries/>, for example these two:**
  - [http://sourceforge.net/projects/boost/files/boost-binaries/1.57.0/boost\\_1\\_57\\_0-msvc-12.0-32.exe/download](http://sourceforge.net/projects/boost/files/boost-binaries/1.57.0/boost_1_57_0-msvc-12.0-32.exe/download)
  - [http://sourceforge.net/projects/boost/files/boost-binaries/1.57.0/boost\\_1\\_57\\_0-msvc-12.0-64.exe/download](http://sourceforge.net/projects/boost/files/boost-binaries/1.57.0/boost_1_57_0-msvc-12.0-64.exe/download)
- **Python from <https://www.python.org/downloads/>**
  - <https://www.python.org/ftp/python/2.7.9/python-2.7.9.amd64.msi>
  - <https://www.python.org/ftp/python/2.7.9/python-2.7.9.msi>
- (optional) If you plan to build documentation, download and install sphinx-doc as described here: <http://sphinx-doc.org/latest/index.html>
- (optional) 7-zip – <http://www.7-zip.org/a/7z920-x64.msi>
- (optional) Putty – <http://the.eearth.li/~sgtatham/putty/latest/x86/putty.exe>

All explicit links are given just for convenience if you are setting up new environment. You are free to choose other versions or tools, and most likely they will work just fine for BigARTM. Remember to match the following: \* Visual Studio version must match Boost binaries version, unless you build Boost yourself \* Use the same configuration (32 bit or 64 bit) for your Python and BigARTM binaries

## Source code

BigARTM is hosted in public GitHub repository:

<https://github.com/bigartm/bigartm>

We maintain two branches: `master` and `stable`. `master` branch is the latest source code, potentially including some unfinished features. `stable` branch will be lagging behind `master`, and moved forward to `master` as soon as maintainers decide that it is ready. Typically this should happen at the end of each month. At the same point we will introduce a new tag (something like `v0.7.3`) and produce a new release for Windows. In addition, `stable` branch also might receive small urgent fixes in between releases, typically to address critical issues reported by our users. Such fixes will be also included in `master` branch.

To contribute a fix you should `fork` the repository, code your fix and submit a `pull request` against `master` branch. All pull requests are regularly monitored by BigARTM maintainers and will be soon merged. Please, keep monitoring the status of your pull request on `travis`, which is a continuous integration system used by BigARTM project.

## Build C++ code on Windows

The following steps describe the procedure to build BigARTM's C++ code on Windows.

- Download and install [GitHub for Windows](#).
- Clone <https://github.com/bigartm/bigartm/> repository to any location on your computer. This location is further referred to as `$ (BIGARTM_ROOT)`.
- Download and install Visual Studio 2012 or any newer version. BigARTM will compile just fine with any edition, including any Visual Studio Express edition (available at [www.visualstudio.com](http://www.visualstudio.com)).
- Install `CMake` (tested with `cmake-3.0.1`, Win32 Installer).

Make sure that `CMake` executable is added to the `PATH` environmental variable. To achieve this either select the option “Add `CMake` to the system `PATH` for all users” during installation of `CMake`, or add it to the `PATH` manually.

- Download and install Boost 1.55 or any newer version.

We suggest to use the [Prebuilt Windows Binaries](#). Make sure to select version that match your version of Visual Studio. You may choose to work with either x64 or Win32 configuration, both of them are supported.

- Configure system variables `BOOST_ROOT` and `Boost_LIBRARY_DIR`.

If you have installed boost from the link above, and used the default location, then the setting should look similar to this:

```
setx BOOST_ROOT C:\local\boost_1_56_0
setx BOOST_LIBRARYDIR C:\local\boost_1_56_0\lib32-msvc-12.0
```

For all future details please refer to the documentation of [FindBoost module](#). We also encourage new `CMake` users to step through [CMake tutorial](#).

- Install Python 2.7 (tested with [Python 2.7.6](#)).

You may choose to work with either x64 or Win32 version of the Python, but make sure this matches the configuration of BigARTM you have choosed earlier. The x64 installation of python will be incompatible with 32 bit BigARTM, and virse versus.

- Use CMake to generate Visual Studio projects and solution files. To do so, open a command prompt, change working directory to `$(BIGARTM_ROOT)` and execute the following commands:

```
mkdir build
cd build
cmake ..
```

You might have to explicitly specify the `cmake` generator, especially if you are working with x64 configuration. To do so, use the following syntax:

```
cmake .. -G"Visual Studio 12 Win64"
```

CMake will generate Visual Studio under `$(BIGARTM_ROOT)/build/`.

- Open generated solution in Visual Studio and build it as you would usually build any other Visual Studio solution. You may also use MSBuild from Visual Studio command prompt.

The build will output result into the following folders:

- `$(BIGARTM_ROOT)/build/bin/[Debug|Release]` — binaries (.dll and .exe)
- `$(BIGARTM_ROOT)/build/lib/[Debug|Release]` — static libraries

At this point you should be able to run BigARTM tests, located here:  
`$(BIGARTM_ROOT)/build/bin/*/artm_tests.exe`.

## Python code on Windows

- Install Python 2.7 (this step is already done if you are following the instructions above),
- Add Python to the PATH environmental variable  
<http://stackoverflow.com/questions/6318156/adding-python-path-on-windows-7>
- Follow the instructions in README file in directory `$(BIGARTM_ROOT)/3rdparty/protobuf/python/`. In brief, this instructions ask you to run the following commands:

```
python setup.py build
python setup.py test
python setup.py install
```

On second step you fill see two failing tests:

```
Ran 216 tests in 1.252s
FAILED (failures=2)
```

This 2 failures are OK to ignore.

At this point you should be able to run BigARTM tests for Python, located under `$(BIGARTM_ROOT)/python/tests/`.

- [Optional] Download and add to MSVS Python Tools 2.0. All necessary instructions can be found at <https://pytools.codeplex.com/>. This will allow you debug you Python scripts using Visual Studio. You may start with the following solution: `$(BIGARTM_ROOT)/src/artm_vs2012.sln`.

## Compiling .proto files on Windows

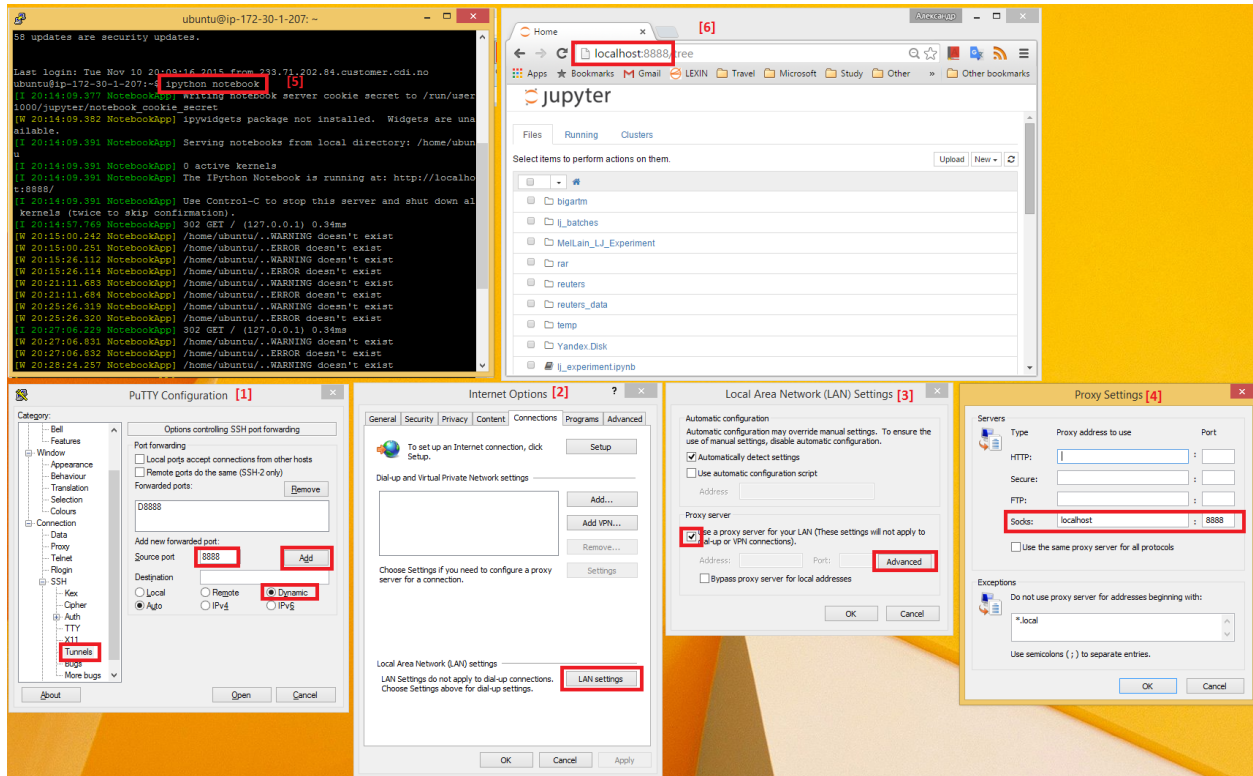
1. Open a new command prompt
2. Copy the following file into `$(BIGARTM_ROOT)/src/`
  - `$(BIGARTM_ROOT)/build/bin/CONFIG/protoc.exe`Here CONFIG can be either Debug or Release (both options will work equally well).
3. Change working directory to `$(BIGARTM_ROOT)/src/`
4. Run the following commands

```
.\protoc.exe --cpp_out=. --python_out=. .\artm\messages.proto  
.\protoc.exe --cpp_out=. .\artm\core\internals.proto
```

## Working with iPython notebooks remotely

It turned out to be common scenario to run BigARTM on a Linux server (for example on Amazon EC2), while connecting to it from Windows through `putty`. Here is a convenient way to use `ipython notebook` in this scenario:

1. Connect to the Linux machine via `putty`. Putty needs to be configured with dynamic tunnel for port 8888 as describe here on [this page](#) (port 8888 is a default port for `ipython notebook`). The same page describes how to configure internet properties:  
*Clicking on Settings in Internet Explorer, or Proxy Settings in Google Chrome, should open this dialogue. Navigate through to the Advanced Proxy section and add localhost:9090 as a SOCKS Proxy.*
2. Start `ipython notebook` in your `putty` terminal.
3. Open your favourite browser on Windows, and go to <http://localhost:8888>. Enjoy your notebook while the engine runs on remotely :)



## Build C++ code on Linux

Refer to [Basic BigARTM tutorial](#) for Linux and Mac OS-X users.

## Code style

### Configure Visual Studio

Open *Tools / Text Editor / All languages / Tabs* and configure as follows:

- Indenting - smart,
- Tab size - 2,
- Indent size - 2,
- Select “insert spaces”.

We also suggest to configure Visual Studio to [show space and tab crlf characters](#) (shortcut: Ctrl+R, Ctrl+W), and [enable vertical line at 120 characters](#).

In the code we follow [google code style](#) with the following changes:

- Exceptions are allowed
- Indentation must be 2 spaces. Tabs are not allowed.
- No lines should exceed 120 characters.

All .h and .cpp files under `$(BIGARTM_ROOT)/src/artm/` must be verified for code style with `cpplint.py` script. Files, generated by protobuf compiler, are the only exceptions from this rule.

To run the script you need some version of Python installed on your machine. Then execute the script like this:

```
python cpplint.py --linelength=120 <filename>
```

On Windows you may run this master-script to check all required files:

```
$ (BIGARTM_ROOT/utils/cpplint_all.bat.
```

Wiki pages:

- [Create New Regularizer](#)
- [Q & A](#)

---

## Legacy documentation pages

---

Legacy pages are kept to preserve existing user's links (favourites in browser, etc).

### Basic BigARTM tutorial for Linux and Mac OS-X users

Currently there is no distribution package of BigARTM for Linux. BigARTM had been tested on several Linux OS, and it is known to work well, but you have to get the source code and compile it locally on your machine.

#### Download sources and build

Clone the latest BigARTM code from our github repository, and build it via CMake as in the following script.

```
sudo apt-get install git make cmake build-essential libboost-all-dev
cd ~
git clone --branch=stable https://github.com/bigartm/bigartm.git
cd bigartm
mkdir build && cd build
cmake ..
make
```

#### Running BigARTM from command line

There is a simple utility `bigartm`, which allows you to run BigARTM from command line. To experiment with this tool you need a small dataset, which you can get via the following script. More datasets are available through [Downloads](#) page.

```
cd ~/bigartm
mkdir datasets && cd datasets
wget https://s3-eu-west-1.amazonaws.com/artm/docword.kos.txt.gz
wget https://s3-eu-west-1.amazonaws.com/artm/vocab.kos.txt
gunzip docword.kos.txt.gz
../build/src/bigartm/bigartm -d docword.kos.txt -v vocab.kos.txt
```

## Configure BigARTM Python API

For more advanced scenarios you need to configure Python interface for BigARTM. To use BigARTM from Python you need to use Google Protobuf. We recommend to use 'protobuf 2.5.1-pre', included in bigartm/3rdparty.

```
# Step 1 - add BigARTM python bindings to PYTHONPATH
export PYTHONPATH=~/bigartm/python:$PYTHONPATH

# Step 2 - install google protobuf
cd ~/bigartm
cp build/3rdparty/protobuf-cmake/protoc/protoc 3rdparty/protobuf/src/
cd 3rdparty/protobuf/python
python setup.py build
sudo python setup.py install

# Step 3 - point ARTM_SHARED_LIBRARY variable to libartm.so (libartm.dylib) location
export ARTM_SHARED_LIBRARY=~/bigartm/build/src/artm/libartm.so      # for linux
export ARTM_SHARED_LIBRARY=~/bigartm/build/src/artm/libartm.dylib  # for Mac OS X
```

At this point you may run examples under ~/bigartm/python/examples.

## Troubleshooting

```
>python setup.py build
File "setup.py", line 52
    print "Generating %s..." % output
SyntaxError: Missing parentheses in call to `print`
```

This error may happen during google protobuf installation. It indicates that you are using Python 3, which is not supported by BigARTM. (see [this question on StackOverflow](#) for more details on the error around *print*). Please use Python 2.7.9 to workaround this issue.

```
ubuntu@192.168.0.1:~/bigartm/python/examples$ python example01_synthetic_collection.py
Traceback (most recent call last):
  File "example01_synthetic_collection.py", line 6, in <module>
    import artm.messages_pb2, artm.library, random, uuid
ImportError: No module named artm.messages_pb2
```

This error indicate that python is unable to locate messages\_pb2.py and ``library.py files. Please verify if you executed Step #1 in the instructions above.

```
ubuntu@192.168.0.1:~/bigartm/python/examples$ python example01_synthetic_collection.py
Traceback (most recent call last):
  File "example01_synthetic_collection.py", line 6, in <module>
    import artm.messages_pb2, artm.library, random, uuid
  File "/home/ubuntu/bigartm/python/messages_pb2.py", line 4, in <module>
    from google.protobuf import descriptor as _descriptor
ImportError: No module named google.protobuf
```

This error indicated that python is unable to locate protobuf library. Please verify if you executed Step #2 in the instructions above. If you do not have permissions to execute `sudo python setup.py install` step, you may also try to update PYTHONPATH manually: `PYTHONPATH="/home/ubuntu/bigartm/3rdparty/protobuf/python:/home/ubuntu/bigartm/python:$PYTHONPATH"`



```
ubuntu@192.168.0.1:~/bigartm/python/examples$ python example01_synthetic_collection.py
libartm.so: cannot open shared object file: No such file or directory,
fall back to ARTM_SHARED_LIBRARY environment variable
Traceback (most recent call last):
  File "example01_synthetic_collection.py", line 27, in <module>
    with artm.library.MasterComponent() as master:
  File "/home/ubuntu/bigartm/python/artm/library.py", line 179, in __init__
    lib = Library().lib_
  File "/home/ubuntu/bigartm/python/artm/library.py", line 107, in __init__
    self.lib_ = ctypes.CDLL(os.environ['ARTM_SHARED_LIBRARY'])
  File "/usr/lib/python2.7/UserDict.py", line 23, in __getitem__
    raise KeyError(key)
KeyError: 'ARTM_SHARED_LIBRARY'
```

This error indicate that BigARTM's python interface can not locate libartm.so (libartm.dylib) files. Please verify if you executed Step #3 correctly.

## BigARTM on Travis-CI

To get a live usage example of BigARTM you may check BigARTM's [.travis.yml](#) script and the latest [continuous integration build](#).

## Basic BigARTM tutorial for Windows users

This tutorial gives guidelines for installing and running existing BigARTM examples via command-line interface and from Python environment.

### Download

Download latest binary distribution of BigARTM from <https://github.com/bigartm/bigartm/releases>. Explicit download links can be found at [Downloads](#) section (for 32 bit and 64 bit configurations).

The distribution will contain pre-build binaries, command-line interface and BigARTM API for Python. The distribution also contains a simple dataset and few python examples that we will be running in this tutorial. More datasets in BigARTM-compatible format are available in the [Downloads](#) section.

Refer to [Windows distribution](#) for details about other files, included in the binary distribution package.

### Running BigARTM from command line

No installation steps are required to run BigARTM from command line. After unpacking binary distribution simply open command prompt (cmd.exe), change current directory to bin folder inside BigARTM package, and run `cpp_client.exe` application as in the following example. As an optional step, we recommend to add bin folder of the BigARTM distribution to your PATH system variable.

```
>C:\BigARTM\bin>set PATH=%PATH%;C:\BigARTM\bin
>C:\BigARTM\bin>cpp_client.exe -v ../python/examples/vocab.kos.txt -d ../python/examples/docword.kos
Parsing text collection... OK.
Iteration 1 took 197 milliseconds.
  Test perplexity = 7108.35,
  Train perplexity = 7106.18,
```

```

Test sparsity theta = 0,
Train sparsity theta = 0,
Sparsity phi = 0.000144802,
Test items processed = 343,
Train items processed = 3087,
Kernel size = 5663,
Kernel purity = 0.958901,
Kernel contrast = 0.292389
Iteration 2 took 195 milliseconds.
Test perplexity = 2563.31,
Train perplexity = 2517.07,
Test sparsity theta = 0,
Train sparsity theta = 0,
Sparsity phi = 0.000144802,
Test items processed = 343,
Train items processed = 3087,
Kernel size = 5559.5,
Kernel purity = 0.956709,
Kernel contrast = 0.298198
...
#1: november(0.054) poll(0.015) bush(0.013) kerry(0.012) polls(0.012) governor(0.011)
#2: bush(0.0083) president(0.0059) republicans(0.0047) house(0.0042) people(0.0039) administration(0.0076)
#3: bush(0.031) iraq(0.018) war(0.012) kerry(0.0096) president(0.0078) administration(0.0076)
#4: kerry(0.018) democratic(0.013) dean(0.012) campaign(0.0097) poll(0.0095) race(0.0082)
ThetaMatrix (last 7 processed documents, ids = 1995,1996,1997,1998,1992,2000,1994):
Topic0: 0.02104 0.02155 0.00604 0.00835 0.00965 0.00006 0.91716
Topic1: 0.15441 0.76643 0.06484 0.11643 0.20409 0.00006 0.00957
Topic2: 0.00399 0.16135 0.00093 0.03890 0.10498 0.00001 0.00037
Topic3: 0.82055 0.05066 0.92819 0.83632 0.68128 0.99987 0.07289

```

We recommend to download larger datasets, available in [Downloads](#) section. All docword and vocab files can be consumed by BigARTM exactly as in the previous example.

Internally BigARTM always parses such files into batches format (for example, [enron\\_1k](#) (7.1 MB)). If you have downloaded such pre-parsed collection, you may feed it into BigARTM as follows:

```

>C:\BigARTM\bin>cpp_client.exe --batch_folder C:\BigARTM\enron
Reuse 40 batches in folder 'enron'
Loading dictionary file... OK.
Iteration 1 took 2502 milliseconds.

```

For more information about `cpp_client.exe` refer to [/ref/cpp\\_client](#) section.

## Configure BigARTM Python API

1. Install Python, for example from the following links:

- Python 2.7.9, 64 bit – <https://www.python.org/ftp/python/2.7.9/python-2.7.9.amd64.msi>, or
- Python 2.7.9, 32 bit – <https://www.python.org/ftp/python/2.7.9/python-2.7.9.msi>

Remember that the version of BigARTM package must match your version Python installed on your machine. If you have 32 bit operating system then you must select 32 bit for Python and BigARTM package. If you have 64 bit operating system then you are free to select either version. However, please note that memory usage of 32 bit processes is limited by 2 GB. For this reason we recommend to select 64 bit configurations.

Also you need to have several Python libraries to be installed on your machine:

- numpy >= 1.9.2
  - scipy >= 0.15.0
  - pandas >= 0.16.2
  - scikit-learn >= 0.16.1
2. Add C:\BigARTM\bin folder to your PATH system variable, and add C:\BigARTM\python to your PYTHONPATH system variable:

```
set PATH=%PATH%;C:\BigARTM\bin
set PATH=%PATH%;C:\Python27;C:\Python27\Scripts
set PYTHONPATH=%PYTHONPATH%;C:\BigARTM\Python
```

Remember to change C:\BigARTM and C:\Python27 with your local folders.

3. Setup *Google Protocol Buffers* library, included in the BigARTM release package.
  - Copy C:\BigARTM\bin\protoc.exe file into C:\BigARTM\protobuf\src folder
  - Run the following commands from command prompt

```
cd C:\BigARTM\protobuf\Python
python setup.py build
python setup.py install
```

Avoid python setup.py test step, as it produces several confusing errors. Those errors are harmless. For further details about protobuf installation refer to [protobuf/python/README](#).

If you are getting errors when configuring or using Python API, please refer to Troubleshooting chapter in [Basic BigARTM tutorial for Linux and Mac OS-X users](#). The list of issues is common between Windows and Linux.

## Running BigARTM from Python API

Refer to ARTM notebook ([in Russian](#) or [in English](#)), which describes high-level Python API of BigARTM.

## Enabling Basic BigARTM Regularizers

This paper describes the experiment with topic model regularization in BigARTM library using [experiment02\\_artm.py](#). The script provides the possibility to learn topic model with three regularizers (sparsing Phi, sparsing Theta and pairwise topic decorrelation in Phi). It also allows the monitoring of learning process by using quality measures as hold-out perplexity, Phi and Theta sparsity and average topic kernel characteristics.

**Warning:** Note that perplexity estimation can influence the learning process in the online algorithm, so we evaluate perplexity only once per 20 synchronizations to avoid this influence. You can change the frequency using `test_every` variable.

We suggest you to have BigARTM installed in \$YOUR\_HOME\_DIRECTORY. To proceed the experiment you need to execute the following steps:

1. Download the collection, represented as BigARTM batches:
  - [https://s3-eu-west-1.amazonaws.com/artm/enwiki-20141208\\_1k.7z](https://s3-eu-west-1.amazonaws.com/artm/enwiki-20141208_1k.7z)
  - [https://s3-eu-west-1.amazonaws.com/artm/enwiki-20141208\\_10k.7z](https://s3-eu-west-1.amazonaws.com/artm/enwiki-20141208_10k.7z)

This data represents a complete dump of the English Wikipedia (approximately 3.7 million documents). The size of one batch in first version is 1000 documents and 10000 in the second one. We used 10000. The decompressed folder with batches should be put into `$YOUR_HOME_DIRECTORY`. You also need to move there the dictionary file from the batches folder.

The batch, you'd like to use for hold-out perplexity estimation, also must be placed into `$YOUR_HOME_DIRECTORY`. In our experiment we used the batch named `243af5b8-beab-4332-bb42-61892df5b044.batch`.

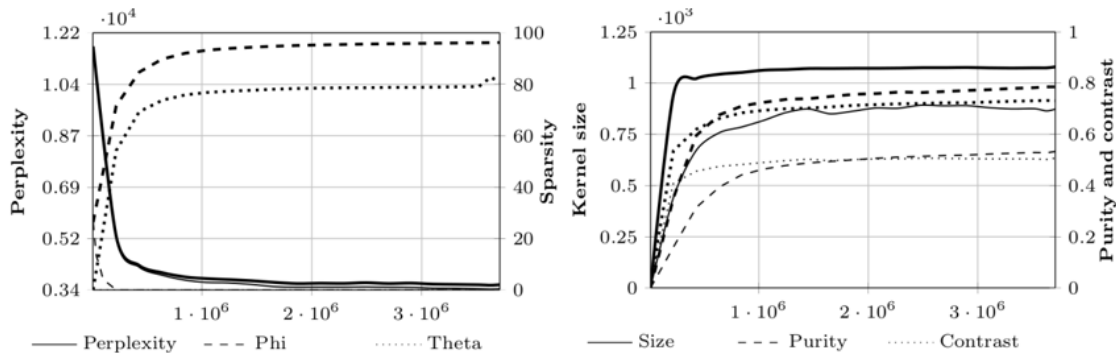
- The next step is the script preparation. Open it's code and find the declaration(-s) of variable(-s)
  - `home_folder` (line 8) and assign it the path `$YOUR_HOME_DIRECTORY`;
  - `batch_size` (line 28) and assign it the chosen size of batch;
  - `batches_disk_path` (line 36) and replace the string 'wiki\_10k' with the name of your directory with batches;
  - `test_batch_name` (line 43) and replace the string with direct batch's name with the name of your test batch;
  - `tau_decor`, `tau_phi` and `tau_theta` (lines 57-59) and substitute the values you'd like to use.
- If you want to estimate the final perplexity on another, larger test sample, put chosen batches into test folder (in `$YOUR_HOME_DIRECTORY` directory). Then find in the code of the script the declaration of variable `save_and_test_model` (line 30) and assign it `True`.
- After all launch the script. Current measures values will be printed into console. Note, that after synchronizations without perplexity estimation it's value will be replaced with string 'NO'. The results of synchronizations with perplexity estimation in addition will be put in corresponding files in results folder. The file format is general for all measures: the set of strings «(accumulated number of processed documents, measure value)»:

```
(10000, 0.018)
(220000, 0.41)
(430000, 0.456)
(640000, 0.475)
...
```

These files can be used for plot building.

If desired, you can easy change values of any variable in the code of script since it's sense is clearly commented. If you used all parameters and data identical our experiment you should get the results, close to these ones

Model/Functional	$\mathcal{P}_{10k}$	$\mathcal{P}_{100k}$	$S_{\Phi}$	$S_{\Theta}$	$\mathcal{K}_s$	$\mathcal{K}_p$	$\mathcal{K}_c$
LDA	3436	3801	0.0	0.0	873	0.533	0.507
ARTM	3577	3947	96.3	80.9	1079	0.785	0.731

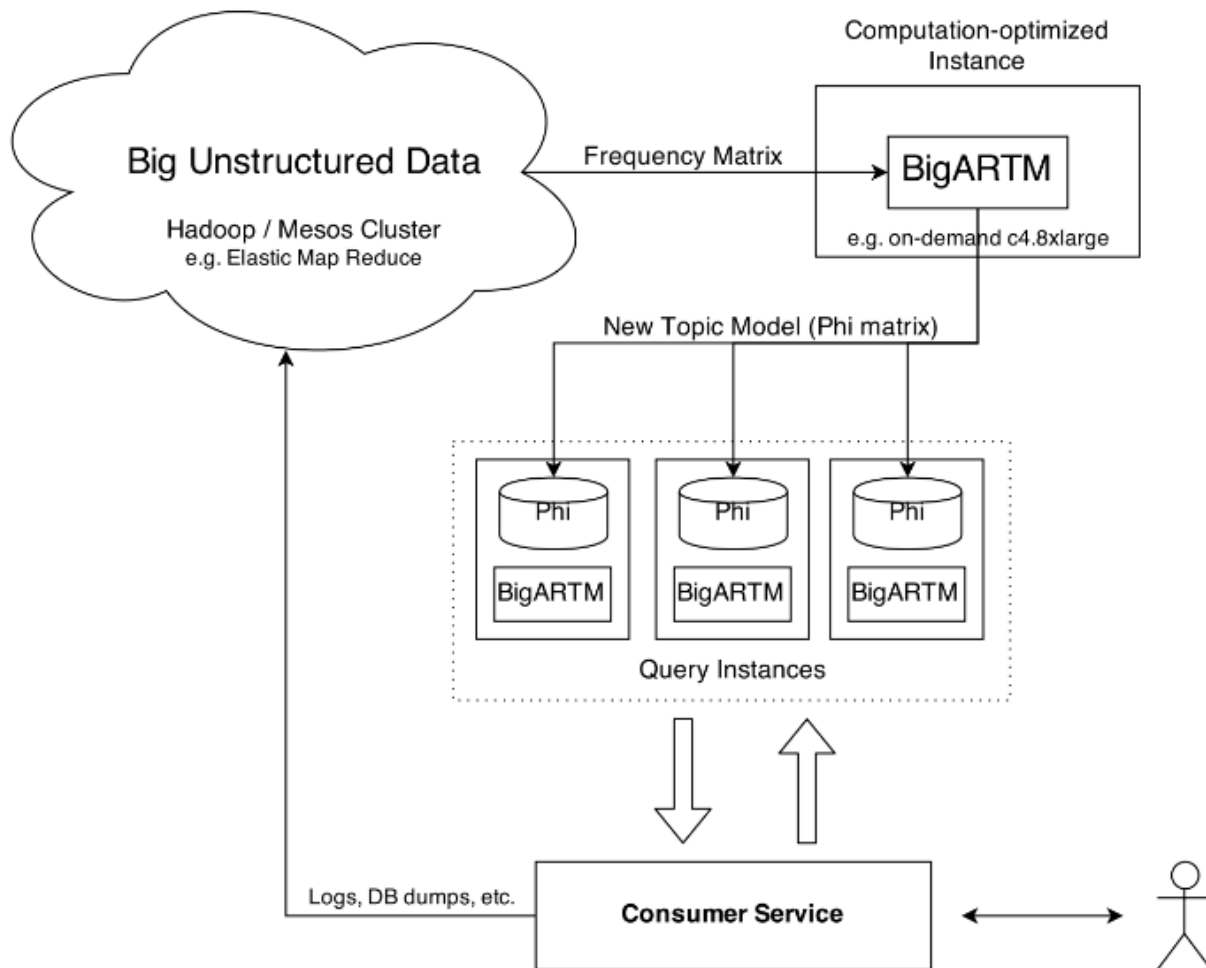


Here you can see the results of comparison between ARTM and LDA models. To make the experiment with LDA instead of ARTM you only need to change the values of variables `tau_decor`, `tau_phi` and `tau_theta` to 0,  $1 / \text{topics\_count}$  and  $1 / \text{topics\_count}$  respectively and run the script again.

**Warning:** Note, that we used machine with 8 cores and 15 Gb RAM for our experiment.

## BigARTM as a Service

The following diagram shows a suggested topology for a query service that involve topic modelling on Big Data.



Here the main use for Hadoop / MapReduce is to process your Big Unstructured Data into a compact bag-of-words representation. Due to out-of-core design and extreme performance BigARTM will be able to handle this data on a single compute-optimized node. The resulting topic model should be replicated on all query instances that serve user requests.

To avoid query-time dependency on BigARTM component you may want to infer topic distributions  $\theta_{\{td\}}$  for new documents in your code. This can be done as follows. Start from uniform topic assignment  $\theta_{\{td\}} = 1 / |\mathbf{T}|$  and update it in the following loop:

```

initialize  $\theta_{td}$  for all  $t \in T$ ;
repeat
   $Z_w := \sum_{t \in T} \phi_{wt}^{i-1} \theta_{td}$  for all  $w \in d$ ;
   $\theta_{td} := \frac{1}{n_d} \sum_{w \in d} n_{dw} \phi_{wt}^{i-1} \theta_{td} / Z_w$  for all  $t \in T$ ;
until  $\theta_d$  converges;

```

where  $n_{dw}$  is the number of word  $w$  occurrences in document  $d$ ,  $\phi_{wt}$  is an element of the Phi matrix. In BigARTM the loop is repeated `ModelConfig.inner_iterations_count` times (default to 10). To precisely replicate BigARTM behavior one needs to account for class weights and include regularizers. Please contact us if you need more details.

## BigARTM: The Algorithm Under The Hood

ToDo: link BigARTM to online batch PLSA algorithm.

ToDo: explain the notation in the algorithm.

ToDo: update the algorithm with regularization.

---

### Algorithm 1 BigARTM's algorithm

---

```

1: Initialize  $\phi_{wt}^0$  for all  $w \in W$  and  $t \in T$ ;
2: for all  $i = 1, \dots, I$  do
3:    $n_{wt}^i := 0, \tilde{n}_t^i := 0$  for all  $w \in W$  and  $t \in T$ ;
4:   for all batches  $D_j, j = 1, \dots, J$  do
5:      $\tilde{n}_{wt} := 0, \tilde{n}_t := 0$  for all  $w \in W$  and  $t \in T$ ;
6:     for all  $d \in D_j$  do
7:       initialize  $\theta_{td}$  for all  $t \in T$ ;
8:       repeat
9:          $Z_w := \sum_{t \in T} \phi_{wt}^{i-1} \theta_{td}$  for all  $w \in d$ ;
10:         $\theta_{td} := \frac{1}{n_d} \sum_{w \in d} n_{dw} \phi_{wt}^{i-1} \theta_{td} / Z_w$  for all  $t \in T$ ;
11:       until  $\theta_d$  converges;
12:       increment  $\tilde{n}_{wt}, \tilde{n}_t$  by  $n_{dw} \phi_{wt}^{i-1} \theta_{td} / Z_w$  for all  $w \in W$  and  $t \in T$ ;
13:        $n_{wt}^i := n_{wt}^i + \tilde{n}_{wt}^i$  for all  $w \in W$  and  $t \in T$ ;
14:        $n_t^i := n_t^i + \tilde{n}_t^i$  for all  $t \in T$ ;
15:    $\phi_{wt}^i := \frac{n_{wt}^i}{n_t^i}$  for all  $w \in W$  and  $t \in T$ ;

```

---

In this algorithm most CPU resources are consumed on steps 8-11 to infer topic distribution for each document. This operation can be executed concurrently across documents or batches. In BigARTM this parallelization is done across batches to avoid splitting the work into too small junks.

Processing each batch produces counters  $\tilde{n}_{wt}$  and  $\tilde{n}_t$ , which should be then merged with the corresponding counters coming from other batches. Since this information is produced by multiple concurrent threads the merging process should be thread-safe and properly synchronised. Our solution is to store all counters  $\tilde{n}_{wt}$

$n_{wt}$  and  $n_t$  into a single queue, from where they can be picked up by a single *merger thread*. This thread will then accumulate the counters without any locking.

Further in this text the term *outer iteration loop* stands for the loop at the step 2, and the term *inner iteration loop* stands for the loop at step 8. Instead of “repeat until it converges” criteria current implementation uses a fixed number of iterations, which is configured manually by the user.

Step 15 is incorporated into all steps that require  $\phi_{wt}$  (e.g. into steps 9, 10 and 11). These steps utilize counters from the previous iteration ( $n^{i-1}_{wt}$  and  $n^{i-1}_t$ ), which are no longer updated by the merger thread, hence they represent read-only data and can be accessed from multiple threads without any synchronization. At the same time the merger thread will accumulate counters for  $n_{wt}$  and  $n_t$  for the current iteration, again in a lock-free manner.

## Messages

This document explains all protobuf messages that can be transferred between the user code and BigARTM library.

**Warning:** Remember that all fields is marked as *optional* to enhance backwards compatibility of the binary protobuf format. Some fields will result in run-time exception when not specified. Please refer to the documentation of each field for more details.

Note that we discourage any usage of fields marked as *obsolete*. Those fields will be removed in future releases.

### DoubleArray

`class messages_pb2.DoubleArray`

Represents an array of double-precision floating point values.

```
message DoubleArray {
  repeated double value = 1 [packed = true];
}
```

### FloatArray

`class messages_pb2.FloatArray`

Represents an array of single-precision floating point values.

```
message FloatArray {
  repeated float value = 1 [packed = true];
}
```

### BoolArray

`class messages_pb2.BoolArray`

Represents an array of boolean values.

```
message BoolArray {
  repeated bool value = 1 [packed = true];
}
```

## IntArray

**class** messages\_pb2.**IntArray**

Represents an array of integer values.

```
message IntArray {
  repeated int32 value = 1 [packed = true];
}
```

## Item

**class** messages\_pb2.**Item**

Represents a unit of textual information. A typical example of an item is a document that belongs to some text collection.

```
message Item {
  optional int32 id = 1;
  repeated Field field = 2;
  optional string title = 3;
}
```

**Item.id**

An integer identifier of the item.

**Item.field**

A set of all fields withing the item.

**Item.title**

An optional title of the item.

## Field

**class** messages\_pb2.**Field**

Represents a field withing an item. The idea behind fields is that each item might have its title, author, body, abstract, actual text, links, year of publication, etc. Each of this entities should be represented as a Field. The topic model defines how those fields should be taken into account when BigARTM infers a topic model. Currently each field is represented as “bag-of-words” — each token is listed together with the number of its occurrences. Note that each Field is always part of an Item, Item is part of a Batch, and a batch always contains a list of tokens. Therefore, each Field just lists the indexes of tokens in the Batch.

```
message Field {
  optional string name = 1 [default = "@body"];
  repeated int32 token_id = 2;
  repeated int32 token_count = 3;
  repeated int32 token_offset = 4;
```



```

optional string string_value = 5;
optional int64 int_value = 6;
optional double double_value = 7;
optional string date_value = 8;

repeated string string_array = 16;
repeated int64 int_array = 17;
repeated double double_array = 18;
repeated string date_array = 19;
}

```

## Batch

### class messages\_pb2.Batch

Represents a set of items. In BigARTM a batch is never split into smaller parts. When it comes to concurrency this means that each batch goes to a single processor. Two batches can be processed concurrently, but items in one batch are always processed sequentially.

```

message Batch {
  repeated string token = 1;
  repeated Item item = 2;
  repeated string class_id = 3;
  optional string description = 4;
  optional string id = 5;
}

```

#### Batch.token

A set value that defines all tokens than may appear in the batch.

#### Batch.item

A set of items of the batch.

#### Batch.class\_id

A set of values that define for classes (modalities) of tokens. This repeated field must have the same length as *token*. This value is optional, use an empty list indicate that all tokens belong to the default class.

#### Batch.description

An optional text description of the batch. You may describe for example the source of the batch, preprocessing technique and the structure of its fields.

#### Batch.id

Unique identifier of the batch in a form of a GUID (example: 4fb38197-3f09-4871-9710-392b14f00d2e). This field is required.

## Stream

### class messages\_pb2.Stream

Represents a configuration of a stream. Streams provide a mechanism to split the entire collection into virtual subsets (for example, the ‘train’ and ‘test’ streams).

```

message Stream {
  enum Type {
    Global = 0;

```

```
    ItemIdModulus = 1;
}

optional Type type = 1 [default = Global];
optional string name = 2 [default = "@global"];
optional int32 modulus = 3;
repeated int32 residuals = 4;
}
```

**Stream.type**

A value that defines the type of the stream.

Global	Defines a stream containing all items in the collection.
ItemIdModulus	Defines a stream containing all items with ID that matches modulus and residuals. An item belongs to the stream iff the modulo reminder of item ID is contained in the residuals field.

**Stream.name**

A value that defines the name of the stream. The name must be unique across all streams defined in the master component.

## MasterComponentConfig

**class** messages\_pb2.MasterComponentConfig

Represents a configuration of a master component.

```
message MasterComponentConfig {
  optional string disk_path = 2;
  repeated Stream stream = 3;
  optional bool compact_batches = 4 [default = true];
  optional bool cache_theta = 5 [default = false];
  optional int32 processors_count = 6 [default = 1];
  optional int32 processor_queue_max_size = 7 [default = 10];
  optional int32 merger_queue_max_size = 8 [default = 10];
  repeated ScoreConfig score_config = 9;
  optional bool online_batch_processing = 13 [default = false]; // obsolete in BigARTM v0.5.8
  optional string disk_cache_path = 15;
}
```

**MasterComponentConfig.disk\_path**

A value that defines the disk location to store or load the collection.

**MasterComponentConfig.stream**

A set of all data streams to configure in master component. Streams can overlap if needed.

**MasterComponentConfig.compact\_batches**

A flag indicating whether to compact batches in AddBatch() operation. Compaction is a process that shrinks the dictionary of each batch by removing all unused tokens.

**MasterComponentConfig.cache\_theta**

A flag indicating whether to cache theta matrix. Theta matrix defines the discrete probability distribution of each document across the topics in topic model. By default BigARTM infers this distribution every time it processes the document. Option 'cache\_theta' allows to cache this theta matrix and re-use theha values when the same document is processed on the next iteration. This option must be set to 'true' before calling method `ArtmRequestThetaMatrix()`.

**MasterComponentConfig.processors\_count**

A value that defines the number of concurrent processor components. The number of processors should normally not exceed the number of CPU cores.

**MasterComponentConfig.processor\_queue\_max\_size**

A value that defines the maximal size of the processor queue. Processor queue contains batches, prefetch from disk into memory. Recommendations regarding the maximal queue size are as follows:

- the queue size should be at least as large as the number of concurrent processors;

**MasterComponentConfig.merger\_queue\_max\_size**

A value that defines the maximal size of the merger queue. Merger queue size contains an incremental updates of topic model, produced by processor components. Try reducing this parameter if BigARTM consumes too much memory.

**MasterComponentConfig.score\_config**

A set of all scores, available for calculation.

**MasterComponentConfig.online\_batch\_processing**

Obsolete in BigARTM v0.5.8.

**MasterComponentConfig.disk\_cache\_path**

A value that defines a writable disk location where this master component can store some temporary files. This can reduce memory usage, particularly when `cache_theta` option is enabled. Note that on clean shutdown master component will be cleaned this folder automatically, but otherwise it is your responsibility to clean this folder to avoid running out of disk.

## ModelConfig

**class** `messages_pb2.ModelConfig`

Represents a configuration of a topic model.

```
message ModelConfig {
  optional string name = 1 [default = "@model"];
  optional int32 topics_count = 2 [default = 32];
  repeated string topic_name = 3;
  optional bool enabled = 4 [default = true];
  optional int32 inner_iterations_count = 5 [default = 10];
  optional string field_name = 6 [default = "@body"]; // obsolete in BigARTM v0.5.8
  optional string stream_name = 7 [default = "@global"];
  repeated string score_name = 8;
  optional bool reuse_theta = 9 [default = false];
  repeated string regularizer_name = 10;
  repeated double regularizer_tau = 11;
  repeated string class_id = 12;
  repeated float class_weight = 13;
```

```
optional bool use_sparse_bow = 14 [default = true];
optional bool use_random_theta = 15 [default = false];
optional bool use_new_tokens = 16 [default = true];
optional bool opt_for_avx = 17 [default = true];
}
```

**ModelConfig.name**

A value that defines the name of the topic model. The name must be unique across all models defined in the master component.

**ModelConfig.topics\_count**

A value that defines the number of topics in the topic model.

**ModelConfig.topic\_name**

A repeated field that defines the names of the topics. All topic names must be unique within each topic model. This field is optional, but either *topics\_count* or *topic\_name* must be specified. If both specified, then *topics\_count* will be ignored, and the number of topics in the model will be based on the length of *topic\_name* field. When *topic\_name* is not specified the names for all topics will be autogenerated.

**ModelConfig.enabled**

A flag indicating whether to update the model during iterations.

**ModelConfig.inner\_iterations\_count**

A value that defines the fixed number of iterations, performed to infer the theta distribution for each document.

**ModelConfig.field\_name**

Obsolete in BigARTM v0.5.8

**ModelConfig.stream\_name**

A value that defines which stream the model should use.

**ModelConfig.score\_name**

A set of names that defines which scores should be calculated for the model.

**ModelConfig.reuse\_theta**

A flag indicating whether the model should reuse theta values cached on the previous iterations. This option require *cache\_theta* flag to be set to 'true' in *MasterComponentConfig*.

**ModelConfig.regularizer\_name**

A set of names that define which regularizers should be enabled for the model. This repeated field must have the same length as *regularizer\_tau*.

**ModelConfig.regularizer\_tau**

A set of values that define the regularization coefficients of the corresponding regularizer. This repeated field must have the same length as *regularizer\_name*.

**ModelConfig.class\_id**

A set of values that define for which classes (modalities) to build topic model. This repeated field must have the same length as *class\_weight*.

**ModelConfig.class\_weight**

A set of values that define the weights of the corresponding classes (modalities). This repeated field must have the same length as *class\_id*. This value is optional, use an empty list to set equal weights for all classes.

**ModelConfig.use\_sparse\_bow**

A flag indicating whether to use sparse representation of the Bag-of-words data. The default setting (*use\_sparse\_bow = true*) is best suited for processing textual collections where every token is represented in a small fraction of all documents. Dense representation (*use\_sparse\_bow = false*) better fits for non-textual collections (for example for matrix factorization).

Note that *class\_weight* and *class\_id* must not be used together with *use\_sparse\_bow=false*.

**ModelConfig.use\_random\_theta**

A flag indicating whether to initialize  $p(t|d)$  distribution with random uniform distribution. The default setting (*use\_random\_theta* = *false*) sets  $p(t|d) = 1/T$ , where *T* stands for *topics\_count*. Note that *reuse\_theta* flag takes priority over *use\_random\_theta* flag, so that if *reuse\_theta* = *true* and there is a cache entry from previous iteration the cache entry will be used regardless of *use\_random\_theta* flag.

**ModelConfig.use\_new\_tokens**

A flag indicating whether to automatically include new tokens into the topic model. This setting is set to *True* by default. As a result, every new token observed in batches is automatically incorporated into topic model during the next model synchronization (*ArtemSynchronizeModel()*). The *n\_wt\_* weights for new tokens randomly generated from  $[0..1]$  range.

**ModelConfig.opt\_for\_avx**

An experimental flag that allows to disable AVX optimization in processor. By default this option is enabled as on average it adds ca. 40% speedup on physical hardware. You may want to disable this option if you are running on Windows inside virtual machine, or in situation when BigARTM performance degrades from iteration to iteration.

This option does not affect the results, and is only intended for advanced users experimenting with BigARTM performance.

## RegularizerConfig

**class** `messages_pb2.RegularizerConfig`

Represents a configuration of a general regularizer.

```
message RegularizerConfig {
  enum Type {
    SmoothSparseTheta = 0;
    SmoothSparsePhi = 1;
    DecorrelatorPhi = 2;
    LabelRegularizationPhi = 4;
  }

  optional string name = 1;
  optional Type type = 2;
  optional bytes config = 3;
}
```

**RegularizerConfig.name**

A value that defines the name of the regularizer. The name must be unique across all names defined in the master component.

**RegularizerConfig.type**

A value that defines the type of the regularizer.

SmoothSparseTheta	Smooth-sparse regularizer for theta matrix
SmoothSparsePhi	Smooth-sparse regularizer for phi matrix
DecorrelatorPhi	Decorrelator regularizer for phi matrix
LabelRegularizationPhi	Label regularizer for phi matrix

**RegularizerConfig.config**

A serialized protobuf message that describes regularizer config for the specific regularizer type.

## SmoothSparseThetaConfig

**class** messages\_pb2.SmoothSparseThetaConfig

Represents a configuration of a SmoothSparse Theta regularizer.

```
message SmoothSparseThetaConfig {  
  repeated string topic_name = 1;  
  repeated float alpha_iter = 2;  
}
```

SmoothSparseThetaConfig.**topic\_name**

A set of topic names that defines which topics in the model should be regularized. This value is optional, use an empty list to regularize all topics.

SmoothSparseThetaConfig.**alpha\_iter**

A field of the same length as *ModelConfig.inner\_iterations\_count* that defines relative regularization weight for every iteration inner iterations. The actual regularization value is calculated as product of *alpha\_iter[i]* and *ModelConfig.regularizer\_tau*.

To specify different regularization weight for different topics create multiple regularizers with different *topic\_name* set, and use different values of *ModelConfig.regularizer\_tau*.

## SmoothSparsePhiConfig

**class** messages\_pb2.SmoothSparsePhiConfig

Represents a configuration of a SmoothSparse Phi regularizer.

```
message SmoothSparsePhiConfig {  
  repeated string topic_name = 1;  
  repeated string class_id = 2;  
  optional string dictionary_name = 3;  
}
```

SmoothSparsePhiConfig.**topic\_name**

A set of topic names that defines which topics in the model should be regularized. This value is optional, use an empty list to regularize all topics.

SmoothSparsePhiConfig.**class\_id**

This set defines which classes in the model should be regularized. This value is optional, use an empty list to regularize all classes.

SmoothSparsePhiConfig.**dictionary\_name**

An optional value defining the name of the dictionary to use. The entries of the dictionary are expected to have *DictionaryEntry.key\_token*, *DictionaryEntry.class\_id* and *DictionaryEntry.value* fields. The actual regularization value will be calculated as a product of *DictionaryEntry.value* and *ModelConfig.regularizer\_tau*.

This value is optional, if no dictionary is specified than all tokens will be regularized with the same weight.

## DecorrelatorPhiConfig

**class** messages\_pb2.DecorrelatorPhiConfig

Represents a configuration of a Decorrelator Phi regularizer.

```
message DecorrelatorPhiConfig {
  repeated string topic_name = 1;
  repeated string class_id = 2;
}
```

#### DecorrelatorPhiConfig.topic\_name

A set of topic names that defines which topics in the model should be regularized. This value is optional, use an empty list to regularize all topics.

#### DecorrelatorPhiConfig.class\_id

This set defines which classes in the model should be regularized. This value is optional, use an empty list to regularize all classes.

## LabelRegularizationPhiConfig

### class messages\_pb2.LabelRegularizationPhiConfig

Represents a configuration of a Label Regularizer Phi regularizer.

```
message LabelRegularizationPhiConfig {
  repeated string topic_name = 1;
  repeated string class_id = 2;
  optional string dictionary_name = 3;
}
```

#### LabelRegularizationPhiConfig.topic\_name

A set of topic names that defines which topics in the model should be regularized.

#### LabelRegularizationPhiConfig.class\_id

This set defines which classes in the model should be regularized. This value is optional, use an empty list to regularize all classes.

#### LabelRegularizationPhiConfig.dictionary\_name

An optional value defining the name of the dictionary to use.

## RegularizerInternalState

### class messages\_pb2.RegularizerInternalState

Represents an internal state of a general regularizer.

```
message RegularizerInternalState {
  enum Type {
    MultiLanguagePhi = 5;
  }

  optional string name = 1;
  optional Type type = 2;
  optional bytes data = 3;
}
```

## DictionaryConfig

**class** messages\_pb2.DictionaryConfig

Represents a static dictionary.

```
message DictionaryConfig {
  optional string name = 1;
  repeated DictionaryEntry entry = 2;
  optional int32 total_token_count = 3;
  optional int32 total_items_count = 4;
}
```

DictionaryConfig.**name**

A value that defines the name of the dictionary. The name must be unique across all dictionaries defined in the master component.

DictionaryConfig.**entry**

A list of all entries of the dictionary.

DictionaryConfig.**total\_token\_count**

A sum of *DictionaryEntry.token\_count* across all entries in this dictionary. The value is optional and might be missing when all entries in the dictionary does not carry the *DictionaryEntry.token\_count* attribute.

DictionaryConfig.**total\_items\_count**

A sum of *DictionaryEntry.items\_count* across all entries in this dictionary. The value is optional and might be missing when all entries in the dictionary does not carry the *DictionaryEntry.items\_count* attribute.

## DictionaryEntry

**class** messages\_pb2.DictionaryEntry

Represents one entry in a static dictionary.

```
message DictionaryEntry {
  optional string key_token = 1;
  optional string class_id = 2;
  optional float value = 3;
  repeated string value_tokens = 4;
  optional FloatArray values = 5;
  optional int32 token_count = 6;
  optional int32 items_count = 7;
}
```

DictionaryEntry.**key\_token**

A token that defines the key of the entry.

DictionaryEntry.**class\_id**

The class of the *DictionaryEntry.key\_token*.

DictionaryEntry.**value**

An optional generic value, associated with the entry. The meaning of this value depends on the usage of the dictionary.

DictionaryEntry.**token\_count**

An optional value, indicating the overall number of token occurrences in some collection.



DictionaryEntry.**items\_count**

An optional value, indicating the overall number of documents containing the token.

## ScoreConfig

**class** messages\_pb2.**ScoreConfig**

Represents a configuration of a general score.

```
message ScoreConfig {
  enum Type {
    Perplexity = 0;
    SparsityTheta = 1;
    SparsityPhi = 2;
    ItemsProcessed = 3;
    TopTokens = 4;
    ThetaSnippet = 5;
    TopicKernel = 6;
  }

  optional string name = 1;
  optional Type type = 2;
  optional bytes config = 3;
}
```

**ScoreConfig.name**

A value that defines the name of the score. The name must be unique across all names defined in the master component.

**ScoreConfig.type**

A value that defines the type of the score.

Perplexity	Defines a config of the Perplexity score
SparsityTheta	Defines a config of the SparsityTheta score
SparsityPhi	Defines a config of the SparsityPhi score
ItemsProcessed	Defines a config of the ItemsProcessed score
TopTokens	Defines a config of the TopTokens score
ThetaSnippet	Defines a config of the ThetaSnippet score
TopicKernel	Defines a config of the TopicKernel score

**ScoreConfig.config**

A serialized protobuf message that describes score config for the specific score type.

## ScoreData

**class** messages\_pb2.**ScoreData**

Represents a general result of score calculation.

```
message ScoreData {
  enum Type {
    Perplexity = 0;
    SparsityTheta = 1;
    SparsityPhi = 2;
    ItemsProcessed = 3;
    TopTokens = 4;
  }
```

```
    ThetaSnippet = 5;
    TopicKernel = 6;
}

optional string name = 1;
optional Type type = 2;
optional bytes data = 3;
}
```

**ScoreData.name**

A value that describes the name of the score. This name will match the name of the corresponding score config.

**ScoreData.type**

A value that defines the type of the score.

Perplexity	Defines a Perplexity score data
SparsityTheta	Defines a SparsityTheta score data
SparsityPhi	Defines a SparsityPhi score data
ItemsProcessed	Defines a ItemsProcessed score data
TopTokens	Defines a TopTokens score data
ThetaSnippet	Defines a ThetaSnippet score data
TopicKernel	Defines a TopicKernel score data

**ScoreData.data**

A serialized protobuf message that provides the specific score result.

## PerplexityScoreConfig

**class** messages\_pb2.**PerplexityScoreConfig**

Represents a configuration of a perplexity score.

```
message PerplexityScoreConfig {
  enum Type {
    UnigramDocumentModel = 0;
    UnigramCollectionModel = 1;
  }

  optional string field_name = 1 [default = "@body"]; // obsolete in BigARTM v0.5.8
  optional string stream_name = 2 [default = "@global"];
  optional Type model_type = 3 [default = UnigramDocumentModel];
  optional string dictionary_name = 4;
  optional float theta_sparsity_eps = 5 [default = 1e-37];
  repeated string theta_sparsity_topic_name = 6;
}
```

**PerplexityScoreConfig.field\_name**

Obsolete in BigARTM v0.5.8

**PerplexityScoreConfig.stream\_name**

A value that defines which stream should be used in perplexity calculation.

## PerplexityScore

**class** messages\_pb2.**PerplexityScore**

Represents a result of calculation of a perplexity score.

```
message PerplexityScore {
  optional double value = 1;
  optional double raw = 2;
  optional double normalizer = 3;
  optional int32 zero_words = 4;
  optional double theta_sparsity_value = 5;
  optional int32 theta_sparsity_zero_topics = 6;
  optional int32 theta_sparsity_total_topics = 7;
}
```

**PerplexityScore.value**

A perplexity value which is calculated as  $\exp(-\text{raw}/\text{normalizer})$ .

**PerplexityScore.raw**

A numerator of perplexity calculation. This value is equal to the likelihood of the topic model.

**PerplexityScore.normalizer**

A denominator of perplexity calculation. This value is equal to the total number of tokens in all processed items.

**PerplexityScore.zero\_words**

A number of tokens that have zero probability  $p(w|t,d)$  in a document. Such tokens are evaluated based on to unigram document model or unigram collection model.

**PerplexityScore.theta\_sparsity\_value**

A fraction of zero entries in the theta matrix.

## SparsityThetaScoreConfig

**class messages\_pb2.SparsityThetaScoreConfig**

Represents a configuration of a theta sparsity score.

```
message SparsityThetaScoreConfig {
  optional string field_name = 1 [default = "@body"]; // obsolete in BigARTM v0.5.8
  optional string stream_name = 2 [default = "@global"];
  optional float eps = 3 [default = 1e-37];
  repeated string topic_name = 4;
}
```

**SparsityThetaScoreConfig.field\_name**

Obsolete in BigARTM v0.5.8

**SparsityThetaScoreConfig.stream\_name**

A value that defines which stream should be used in theta sparsity calculation.

**SparsityThetaScoreConfig.eps**

A small value that defines zero threshold for theta probabilities. Theta values below the threshold will be counted as zeros when calculating theta sparsity score.

**SparsityThetaScoreConfig.topic\_name**

A set of topic names that defines which topics should be used for score calculation. The names correspond to *ModelConfig.topic\_name*. This value is optional, use an empty list to calculate the score for all topics.

## SparsityThetaScore

**class** messages\_pb2.**SparsityThetaScoreConfig**

Represents a result of calculation of a theta sparsity score.

```
message SparsityThetaScore {
  optional double value = 1;
  optional int32 zero_topics = 2;
  optional int32 total_topics = 3;
}
```

**SparsityThetaScore.value**

A value of theta sparsity that is calculated as  $\text{zero\_topics} / \text{total\_topics}$ .

**SparsityThetaScore.zero\_topics**

A numerator of theta sparsity score. A number of topics that have zero probability in a topic-item distribution.

**SparsityThetaScore.total\_topics**

A denominator of theta sparsity score. A total number of topics in a topic-item distributions that are used in theta sparsity calculation.

## SparsityPhiScoreConfig

**class** messages\_pb2.**SparsityPhiScoreConfig**

Represents a configuration of a sparsity phi score.

```
message SparsityPhiScoreConfig {
  optional float eps = 1 [default = 1e-37];
  optional string class_id = 2;
  repeated string topic_name = 3;
}
```

**SparsityPhiScoreConfig.eps**

A small value that defines zero threshold for phi probabilities. Phi values below the threshold will be counted as zeros when calculating phi sparsity score.

**SparsityPhiScoreConfig.class\_id**

A value that defines the class of tokens to use for score calculation. This value corresponds to *ModelConfig.class\_id* field. This value is optional. By default the score will be calculated for the default class ('@default\_class').

**SparsityPhiScoreConfig.topic\_name**

A set of topic names that defines which topics should be used for score calculation. This value is optional, use an empty list to calculate the score for all topics.

## SparsityPhiScore

**class** messages\_pb2.**SparsityPhiScore**

Represents a result of calculation of a phi sparsity score.

```
message SparsityPhiScore {
  optional double value = 1;
  optional int32 zero_tokens = 2;
  optional int32 total_tokens = 3;
}
```

**SparsityPhiScore.value**

A value of phi sparsity that is calculated as  $\text{zero\_tokens} / \text{total\_tokens}$ .

**SparsityPhiScore.zero\_tokens**

A numerator of phi sparsity score. A number of tokens that have zero probability in a token-topic distribution.

**SparsityPhiScore.total\_tokens**

A denominator of phi sparsity score. A total number of tokens in a token-topic distributions that are used in phi sparsity calculation.

## ItemsProcessedScoreConfig

**class** messages\_pb2.ItemsProcessedScoreConfig

Represents a configuration of an items processed score.

```
message ItemsProcessedScoreConfig {
  optional string field_name = 1 [default = "@body"]; // obsolete in BigARTM v0.5.8
  optional string stream_name = 2 [default = "@global"];
}
```

**ItemsProcessedScoreConfig.field\_name**

Obsolete in BigARTM v0.5.8

**ItemsProcessedScoreConfig.stream\_name**

A value that defines which stream should be used in calculation of processed items.

## ItemsProcessedScore

**class** messages\_pb2.ItemsProcessedScore

Represents a result of calculation of an items processed score.

```
message ItemsProcessedScore {
  optional int32 value = 1;
}
```

**ItemsProcessedScore.value**

A number of items that belong to the stream *ItemsProcessedScoreConfig.stream\_name* and have been processed during iterations. Currently this number is aggregated throughout all iterations.

## TopTokensScoreConfig

**class** messages\_pb2.TopTokensScoreConfig

Represents a configuration of a top tokens score.

```
message TopTokensScoreConfig {
  optional int32 num_tokens = 1 [default = 10];
  optional string class_id = 2;
  repeated string topic_name = 3;
}
```

**TopTokensScoreConfig.num\_tokens**

A value that defines how many top tokens should be retrieved for each topic.

**TopTokensScoreConfig.class\_id**

A value that defines for which class of the model to collect top tokens. This value corresponds to *ModelConfig.class\_id* field.

This parameter is optional. By default tokens will be retrieved for the default class ('@default\_class').

**TopTokensScoreConfig.topic\_name**

A set of values that represent the names of the topics to include in the result. The names correspond to *ModelConfig.topic\_name*.

This parameter is optional. By default top tokens will be calculated for all topics in the model.

## TopTokensScore

**class messages\_pb2.TopTokensScore**

Represents a result of calculation of a top tokens score.

```
message TopTokensScore {
  optional int32 num_entries = 1;
  repeated string topic_name = 2;
  repeated int32 topic_index = 3;
  repeated string token = 4;
  repeated float weight = 5;
}
```

The data in this score is represented in a table-like format. sorted on *topic\_index*. The following code block gives a typical usage example. The loop below is guarantied to process all top-N tokens for the first topic, then for the second topic, etc.

```
for (int i = 0; i < top_tokens_score.num_entries(); i++) {
  // Gives a index from 0 to (model_config.topics_size() - 1)
  int topic_index = top_tokens_score.topic_index(i);

  // Gives one of the topN tokens for topic 'topic_index'
  std::string token = top_tokens_score.token(i);

  // Gives the weight of the token
  float weight = top_tokens_score.weight(i);
}
```

**TopTokensScore.num\_entries**

A value indicating the overall number of entries in the score. All the remaining repeated fiels in this score will have this length.

**TopTokensScore.token**

A repeated field of *num\_entries* elements, containing tokens with high probability.

`TopTokensScore.weight`

A repeated field of `num_entries` elements, containing the  $p(t|w)$  probabilities.

`TopTokensScore.topic_index`

A repeated field of `num_entries` elements, containing integers between 0 and `(ModelConfig.topics_count - 1)`.

`TopTokensScore.topic_name`

A repeated field of `num_entries` elements, corresponding to the values of `ModelConfig.topic_name` field.

## ThetaSnippetScoreConfig

`class messages_pb2.ThetaSnippetScoreConfig`

Represents a configuration of a theta snippet score.

```
message ThetaSnippetScoreConfig {
  optional string field_name = 1 [default = "@body"]; // obsolete in BigARTM v0.5.8
  optional string stream_name = 2 [default = "@global"];
  repeated int32 item_id = 3 [packed = true]; // obsolete in BigARTM v0.5.8
  optional int32 item_count = 4 [default = 10];
}
```

`ThetaSnippetScoreConfig.field_name`

Obsolete in BigARTM v0.5.8

`ThetaSnippetScoreConfig.stream_name`

A value that defines which stream should be used in calculation of a theta snippet.

`ThetaSnippetScoreConfig.item_id`

Obsolete in BigARTM v0.5.8.

`ThetaSnippetScoreConfig.item_count`

The number of items to retrieve. `ThetaSnippetScore` will select last `item_count` processed items and return their theta vectors.

## ThetaSnippetScore

`class messages_pb2.ThetaSnippetScore`

Represents a result of calculation of a theta snippet score.

```
message ThetaSnippetScore {
  repeated int32 item_id = 1;
  repeated FloatArray values = 2;
}
```

`ThetaSnippetScore.item_id`

A set of item ids for which theta snippet have been calculated. Items are identified by the item id.

`ThetaSnippetScore.values`

A set of values that define topic probabilities for each item. The length of these repeated values will match the number of item ids specified in `ThetaSnippetScore.item_id`. Each repeated field contains float array of topic probabilities in the natural order of topic ids.

## TopicKernelScoreConfig

**class** messages\_pb2.TopicKernelScoreConfig

Represents a configuration of a topic kernel score.

```
message TopicKernelScoreConfig {
  optional float eps = 1 [default = 1e-37];
  optional string class_id = 2;
  repeated string topic_name = 3;
  optional double probability_mass_threshold = 4 [default = 0.1];
}
```

- *Kernel* of a topic model is defined as the list of all tokens such that the probability  $p(t \mid w)$  exceeds probability mass threshold.
- *Kernel size* of a topic  $t$  is defined as the number of tokens in its kernel.
- *Topic purity* of a topic  $t$  is defined as the sum of  $p(w \mid t)$  across all tokens  $w$  in the kernel.
- *Topic contrast* of a topic  $t$  is defined as the sum of  $p(t \mid w)$  across all tokens  $w$  in the kernel defided by the size of the kernel.

TopicKernelScoreConfig.**eps**

Defines the minimum threshold on kernel size. In most cases this parameter should be kept at the default value.

TopicKernelScoreConfig.**class\_id**

A value that defines the class of tokens to use for score calculation. This value corresponds to `ModelConfig.class_id` field. This value is optional. By default the score will be calculated for the default class ('@default\_class').

TopicKernelScoreConfig.**topic\_name**

A set of topic names that defines which topics should be used for score calculation. This value is optional, use an empty list to calculate the score for all topics.

TopicKernelScoreConfig.**probability\_mass\_threshold**

Defines the probability mass threshold (see the definition of *kernel* above).

## TopicKernelScore

**class** messages\_pb2.TopicKernelScore

Represents a result of calculation of a topic kernel score.

```
message TopicKernelScore {
  optional DoubleArray kernel_size = 1;
  optional DoubleArray kernel_purity = 2;
  optional DoubleArray kernel_contrast = 3;
  optional double average_kernel_size = 4;
  optional double average_kernel_purity = 5;
  optional double average_kernel_contrast = 6;
}
```

TopicKernelScore.**kernel\_size**

Provides the kernel size for all requested topics. The length of this *DoubleArray* is always equal to the overall number of topics. The values of  $-1$  correspond to non-calculated topics. The remaining values carry the kernel size of the requested topics.



**TopicKernelScore.kernel\_purity**

Provides the kernel purity for all requested topics. The length of this *DoubleArray* is always equal to the overall number of topics. The values of -1 correspond to non-calculated topics. The remaining values carry the kernel size of the requested topics.

**TopicKernelScore.kernel\_contrast**

Provides the kernel contrast for all requested topics. The length of this *DoubleArray* is always equal to the overall number of topics. The values of -1 correspond to non-calculated topics. The remaining values carry the kernel contrast of the requested topics.

**TopicKernelScore.average\_kernel\_size**

Provides the average kernel size across all the requested topics.

**TopicKernelScore.average\_kernel\_purity**

Provides the average kernel purity across all the requested topics.

**TopicKernelScore.average\_kernel\_contrast**

Provides the average kernel contrast across all the requested topics.

## TopicModel

**class messages\_pb2.TopicModel**

Represents a topic model. This message can contain data in either dense or sparse format. The key idea behind sparse format is to avoid storing zero  $p(w|t)$  elements of the Phi matrix. Please refer to the description of *TopicModel.topic\_index* field for more details.

To distinguish between these two formats check whether repeated field *TopicModel.topic\_index* is empty. An empty field indicate a dense format, otherwise the message contains data in a sparse format. To request topic model in a sparse format set *GetTopicModelArgs.use\_sparse\_format* field to True when calling *ArtemRequestTopicModel()*.

```
message TopicModel {
  enum OperationType {
    Initialize = 0;
    Increment = 1;
    Overwrite = 2;
    Remove = 3;
    Ignore = 4;
  }

  optional string name = 1 [default = "@model"];
  optional int32 topics_count = 2;
  repeated string topic_name = 3;
  repeated string token = 4;
  repeated FloatArray token_weights = 5;
  repeated string class_id = 6;

  message TopicModelInternals {
    repeated FloatArray n_wt = 1;
    repeated FloatArray r_wt = 2;
  }

  optional bytes internals = 7; // obsolete in BigARTM v0.6.3
  repeated IntArray topic_index = 8;
  repeated OperationType operation_type = 9;
}
```

**TopicModel.name**

A value that describes the name of the topic model (*TopicModel.name*).

**TopicModel.topics\_count**

A value that describes the number of topics in this message.

**TopicModel.topic\_name**

A value that describes the names of the topics included in given *TopicModel* message. This values will represent a subset of topics, defined by *GetTopicModelArgs.topic\_name* message. In case of empty *GetTopicModelArgs.topic\_name* this values will correspond to the entire set of topics, defined in *ModelConfig.topic\_name* field.

**TopicModel.token**

The set of all tokens, included in the topic model.

**TopicModel.token\_weights**

A set of token weights. The length of this repeated field will match the length of the repeated field *TopicModel.token*. The length of each *FloatArray* will match the *TopicModel.topics\_count* field (in dense representation), or the length of the corresponding *IntArray* from *TopicModel.topic\_index* field (in sparse representation).

**TopicModel.class\_id**

A set values that specify the class (modality) of the tokens. The length of this repeated field will match the length of the repeated field *TopicModel.token*.

**TopicModel.internals**

Obsolete in BigARTM v0.6.3.

**TopicModel.topic\_index**

A repeated field used for sparse topic model representation. This field has the same length as *TopicModel.token*, *TopicModel.class\_id* and *TopicModel.token\_weights*. Each element in *topic\_index* is an instance of *IntArray* message, containing a list of values between 0 and the length of *TopicModel.topic\_name* field. This values correspond to the indices in *TopicModel.topic\_name* array, and tell which topics has non-zero  $p(w|t)$  probabilities for a given token. The actual  $p(w|t)$  values can be found in *TopicModel.token\_weights* field. The length of each *IntArray* message in *TopicModel.topic\_index* field equals to the length of the corresponding *FloatArray* message in *TopicModel.token\_weights* field.

**Warning:** Be careful with *TopicModel.topic\_index* when this message represents a subset of topics, defined by *GetTopicModelArgs.topic\_name*. In this case indices correspond to the selected subset of topics, which might not correspond to topic indices in the original *ModelConfig* message.

**TopicModel.operation\_type**

A set of values that define operation to perform on each token when topic model is used as an argument of *ArtemOverwriteTopicModel()*.

Initial	Indicates that a new token should be added to the topic model. Initial <i>n_wt</i> counter will be initialized with random value from [0, 1] range. <i>TopicModel.token_weights</i> is ignored. This operation is ignored if token already exists.
Increment	Indicates that <i>n_wt</i> counter of the token should be increased by values, specified in <i>TopicModel.token_weights</i> field. A new token will be created if it does not exist yet.
Overwrite	Indicates that <i>n_wt</i> counter of the token should be set to the value, specified in <i>TopicModel.token_weights</i> field. A new token will be created if it does not exist yet.
Remove	Indicates that the token should be removed from the topic model. <i>TopicModel.token_weights</i> is ignored.
Ignore	Indicates no operation for the token. The effect is the same as if the token is not present in this message.

## ThetaMatrix

### class messages\_pb2.ThetaMatrix

Represents a theta matrix. This message can contain data in either dense or sparse format. The key idea behind sparse format is to avoid storing zero  $p(t|d)$  elements of the Theta matrix. Sparse representation of Theta matrix is equivalent to sparse representation of Phi matrix. Please, refer to [TopicModel](#) for detailed description of the sparse format.

```
message ThetaMatrix {
  optional string model_name = 1 [default = "@model"];
  repeated int32 item_id = 2;
  repeated FloatArray item_weights = 3;
  repeated string topic_name = 4;
  optional int32 topics_count = 5;
  repeated string item_title = 6;
  repeated IntArray topic_index = 7;
}
```

#### ThetaMatrix.model\_name

A value that describes the name of the topic model. This name will match the name of the corresponding model config.

#### ThetaMatrix.item\_id

A set of item IDs corresponding to *Item.id* values.

#### ThetaMatrix.item\_weights

A set of item ID weights. The length of this repeated field will match the length of the repeated field *ThetaMatrix.item\_id*. The length of each *FloatArray* will match the *ThetaMatrix.topics\_count* field (in dense representation), or the length of the corresponding *IntArray* from *ThetaMatrix.topic\_index* field (in sparse representation).

#### ThetaMatrix.topic\_name

A value that describes the names of the topics included in given *ThetaMatrix* message. This values will represent a subset of topics, defined by *GetThetaMatrixArgs.topic\_name* message. In case of empty *GetTopicModelArgs.topic\_name* this values will correspond to the entire set of topics, defined in *ModelConfig.topic\_name* field.

#### ThetaMatrix.topics\_count

A value that describes the number of topics in this message.

#### ThetaMatrix.item\_title

A set of item titles, corresponding to *Item.title* values. Beware that this field might be empty (e.g. of zero length) if all items did not have title specified in *Item.title*.

#### ThetaMatrix.topic\_index

A repeated field used for sparse theta matrix representation. This field has the same length as *ThetaMatrix.item\_id*, *ThetaMatrix.item\_weights* and *ThetaMatrix.item\_title*. Each element in *topic\_index* is an instance of *IntArray* message, containing a list of values between 0 and the length of *TopicModel.topic\_name* field. This values correspond to the indices in *ThetaMatrix.topic\_name* array, and tell which topics has non-zero  $p(t|d)$  probabilities for a given item. The actual  $p(t|d)$  values can be found in *ThetaMatrix.item\_weights* field. The length of each *IntArray* message in *ThetaMatrix.topic\_index* field equals to the length of the corresponding *FloatArray* message in *ThetaMatrix.item\_weights* field.

**Warning:** Be careful with *ThetaMatrix.topic\_index* when this message represents a subset of topics, defined by *GetThetaMatrixArgs.topic\_name*. In this case indices correspond to the selected subset of topics, which might not correspond to topic indices in the original *ModelConfig* message.

## CollectionParserConfig

**class** messages\_pb2.**CollectionParserConfig**

Represents a configuration of a collection parser.

```
message CollectionParserConfig {
  enum Format {
    BagOfWordsUci = 0;
    MatrixMarket = 1;
  }

  optional Format format = 1 [default = BagOfWordsUci];
  optional string docword_file_path = 2;
  optional string vocab_file_path = 3;
  optional string target_folder = 4;
  optional string dictionary_file_name = 5;
  optional int32 num_items_per_batch = 6 [default = 1000];
  optional string cooccurrence_file_name = 7;
  repeated string cooccurrence_token = 8;
  optional bool use_unity_based_indices = 9 [default = true];
}
```

CollectionParserConfig.**format**

A value that defines the format of a collection to be parsed.

BagOfWordsUci	<p>A bag-of-words collection, stored in UCI format. UCI format must have two files - <i>vocab.*.txt</i> and <i>docword.*.txt</i>, defined by <i>docword_file_path</i> and <i>vocab_file_path</i>. The format of the <i>docword.*.txt</i> file is 3 header lines, followed by NNZ triples:</p> <pre>D W NNZ docID wordID count docID wordID count ... docID wordID count</pre> <p>The file must be sorted on docID. Values of wordID must be unity-based (not zero-based). The format of the <i>vocab.*.txt</i> file is line containing wordID=n. Note that words must not have spaces or tabs. In <i>vocab.*.txt</i> file it is also possible to specify <i>Batch.class_id</i> for tokens, as it is shown in this example:</p> <pre>token1 @default_class token2 custom_class token3 @default_class token4</pre> <p>Use space or tab to separate token from its class. Token that are not followed by class label automatically get “@default_class” as a label (see “token4” in the example).</p>
MatrixMarket	<p>See the description at <a href="http://math.nist.gov/MatrixMarket/formats.html">http://math.nist.gov/MatrixMarket/formats.html</a> In this mode parameter <i>docword_file_path</i> must refer to a file in Matrix Market format. Parameter <i>vocab_file_path</i> is also required and must refer to a dictionary file exported in <i>gensim</i> format (<code>dictionary.save_as_text()</code>).</p>

CollectionParserConfig.docword\_file\_path

A value that defines the disk location of a `docword.*.txt` file (the bag of words file in sparse format).

`CollectionParserConfig.vocab_file_path`

A value that defines the disk location of a `vocab.*.txt` file (the file with the vocabulary of the collection).

`CollectionParserConfig.target_folder`

A value that defines the disk location where to stores all the results after parsing the colleciton. Usually the resulting location will contain a set of *batches*, and a *DictionaryConfig* that contains all unique tokens occured in the collection. Such location can be further passed MasterComponent via *MasterComponentConfig.disk\_path*.

`CollectionParserConfig.dictionary_file_name`

A file name where to save the *DictionaryConfig* message that contains all unique tokens occured in the collection. The file will be created in *target\_folder*.

This parameter is optional. The dictionary will be still collected even when this parameter is not provided, but the resulting dictionary will be only returned as the result of `ArtmRequestParseCollection`, but it will not be stored to disk.

In the resulting dictionary each entry will have the following fields:

- *DictionaryEntry.key\_token* - the textual representation of the token,
- *DictionaryEntry.class\_id* - the label of the default class (“@DefaultClass”),
- *DictionaryEntry.token\_count* - the overall number of occurrences of the token in the collection,
- *DictionaryEntry.items\_count* - the number of documents in the collection, containing the token.
- *DictionaryEntry.value* - the ratio between *token\_count* and *total\_token\_count*.

Use `ArtmRequestLoadDictionary` method to load the resulting dictionary.

`CollectionParserConfig.num_items_per_batch`

A value indicating the desired number of items per batch.

`CollectionParserConfig.cooccurrence_file_name`

A file name where to save the *DictionaryConfig* message that contains information about co-occurrence of all pairs of tokens in the collection. The file will be created in *target\_folder*.

This parameter is optional. No cooccurrence information will be collected if the filename is not provided.

In the resulting dictionary each entry will correspond to two tokens (‘<first>’ and ‘<second>’), and carry the information about co-occurrence of this tokens in the collection.

- *DictionaryEntry.key\_token* - a string of the form ‘<first>~<second>’, produced by concatenation of two tokens together via the tilde symbol (‘~’). <first> tokens is guarantied lexicographic less than the <second> token.
- *DictionaryEntry.class\_id* - the label of the default class (“@DefaultClass”).
- *DictionaryEntry.items\_count* - the number of documents in the collection, containing both tokens (‘<first>’ and ‘<second>’)

Use `ArtmRequestLoadDictionary` method to load the resulting dictionary.

`CollectionParserConfig.cooccurrence_token`

A list of tokens to collect cooccurrence information. A cooccurrence of the pair <first>~<second> will be collected only when both tokens are present in *CollectionParserConfig.cooccurrence\_token*.

`CollectionParserConfig.use_unity_based_indices`

A flag indicating whether to interpret indices in docword file as unity-based or as zero-based. By default ‘*use\_unity\_based\_indices* = *True*’, as required by UCI bag-of-words format.

## SynchronizeModelArgs

**class** messages\_pb2.**SynchronizeModelArgs**

Represents an argument of synchronize model operation.

```
message SynchronizeModelArgs {
  optional string model_name = 1;
  optional float decay_weight = 2 [default = 0.0];
  optional bool invoke_regularizers = 3 [default = true];
  optional float apply_weight = 4 [default = 1.0];
}
```

**SynchronizeModelArgs.model\_name**

The name of the model to be synchronized. This value is optional. When not set, all models will be synchronized with the same decay weight.

**SynchronizeModelArgs.decay\_weight**

The decay weight and *apply\_weight* define how to combine existing topic model with all increments, calculated since the last `ArtmSynchronizeModel()`. This is best described by the following formula:

$$n\_wt\_new = n\_wt\_old * decay\_weight + n\_wt\_inc * apply\_weight,$$

where *n\_wt\_old* describe current topic model, *n\_wt\_inc* describe increment calculated since last `ArtmSynchronizeModel()`, *n\_wt\_new* define the resulting topic model.

Expected values of both parameters are between 0.0 and 1.0. Here are some examples:

- Combination of *decay\_weight=0.0* and *apply\_weight=1.0* states that the previous Phi matrix of the topic model will be disregarded completely, and the new Phi matrix will be formed based on new increments gathered since last model synchronize.
- Combination of *decay\_weight=1.0* and *apply\_weight=1.0* states that new increments will be appended to the current Phi matrix without any decay.
- Combination of *decay\_weight=1.0* and *apply\_weight=0.0* states that new increments will be disregarded, and current Phi matrix will stay unchanged.
- To reproduce Online variational Bayes for LDA algorithm by Matthew D. Hoffman set *decay\_weight = 1 - rho* and *apply\_weight = rho*, where parameter *rho* is defined as  $\rho = \exp(\tau + t, -\kappa)$ . See [Online Learning for Latent Dirichlet Allocation](#) for further details.

**SynchronizeModelArgs.apply\_weight**

See *decay\_weight* for the description.

**SynchronizeModelArgs.invoke\_regularizers**

A flag indicating whether to invoke all phi-regularizers.

## InitializeModelArgs

**class** messages\_pb2.**InitializeModelArgs**

Represents an argument of `ArtmInitializeModel()` operation. Please refer to [example14\\_initialize\\_topic\\_model.py](#) for further information.

```
message InitializeModelArgs {
  enum SourceType {
    Dictionary = 0;
    Batches = 1;
  }
```

```
}

message Filter {
  optional string class_id = 1;
  optional float min_percentage = 2;
  optional float max_percentage = 3;
  optional int32 min_items = 4;
  optional int32 max_items = 5;
  optional int32 min_total_count = 6;
  optional int32 min_one_item_count = 7;
}

optional string model_name = 1;
optional string dictionary_name = 2;
optional SourceType source_type = 3 [default = Dictionary];

optional string disk_path = 4;
repeated Filter filter = 5;
}
```

`InitializeModelArgs.model_name`

The name of the model to be initialized.

`InitializeModelArgs.dictionary_name`

The name of the dictionary containing all tokens that should be initialized.

## GetTopicModelArgs

Represents an argument of `ArtmRequestTopicModel()` operation.

```
message GetTopicModelArgs {
  enum RequestType {
    Pwt = 0;
    Nwt = 1;
  }

  optional string model_name = 1;
  repeated string topic_name = 2;
  repeated string token = 3;
  repeated string class_id = 4;
  optional bool use_sparse_format = 5;
  optional float eps = 6 [default = 1e-37];
  optional RequestType request_type = 7 [default = Pwt];
}
```

`GetTopicModelArgs.model_name`

The name of the model to be retrieved.

`GetTopicModelArgs.topic_name`

The list of topic names to be retrieved. This value is optional. When not provided, all topics will be retrieved.

`GetTopicModelArgs.token`

The list of tokens to be retrieved. The length of this field must match the length of `class_id` field. This field is optional. When not provided, all tokens will be retrieved.



**GetTopicModelArgs.class\_id**

The list of classes corresponding to all tokens. The length of this field must match the length of *token* field. This field is only required together with *token*, otherwise it is ignored.

**GetTopicModelArgs.use\_sparse\_format**

An optional flag that defines whether to use sparse format for the resulting *TopicModel* message. See *TopicModel* message for additional information about the sparse format. Note that setting *use\_sparse\_format* = *true* results in empty *TopicModel.internals* field.

**GetTopicModelArgs.eps**

A small value that defines zero threshold for  $p(w|t)$  probabilities. This field is only used in sparse format.  $p(w|t)$  below the threshold will be excluded from the resulting Phi matrix.

**GetTopicModelArgs.request\_type**

An optional value that defines what kind of data to retrieve in this operation.

Pwt	Indicates that the resulting <i>TopicModel</i> message should contain $p(w t)$ probabilities. This values are normalized to form a probability distribution ( $\sum_w p(w t) = 1$ for all topics $t$ ).
Nwt	Indicates that the resulting <i>TopicModel</i> message should contain internal $n_{wt}$ counters of the topic model. This values represent an internal state of the topic model.

Default setting is to retrieve  $p(w|t)$  probabilities. This probabilities are sufficient to infer  $p(t|d)$  distributions using this topic model.

$n_{wt}$  counters allow you to restore the precise state of the topic model. By passing this values in *ArtmOverwriteTopicModel()* operation you are guarantied to get the model in the same state as you retrieved it. As the result you may continue topic model inference from the point you have stopped it last time.

$p(w|t)$  values can be also restored via *c::func:ArtmOverwriteTopicModel* operation. The resulting model will give the same  $p(t|d)$  distributions, however you should consider this model as *read-only*, and do not call *ArtmSynchronizeModel()* on it.

## GetThetaMatrixArgs

Represents an argument of *ArtmRequestThetaMatrix()* operation.

```
message GetThetaMatrixArgs {
  optional string model_name = 1;
  optional Batch batch = 2;
  repeated string topic_name = 3;
  repeated int32 topic_index = 4;
  optional bool clean_cache = 5 [default = false];
  optional bool use_sparse_format = 6 [default = false];
  optional float eps = 7 [default = 1e-37];
}
```

**GetThetaMatrixArgs.model\_name**

The name of the model to retrieved theta matrix for.

**GetThetaMatrixArgs.batch**

The *Batch* to classify with the model.

**GetThetaMatrixArgs.topic\_name**

The list of topic names, describing which topics to include in the Theta matrix. The values of this field should correspond to values in *ModelConfig.topic\_name*. This field is optional, by default all topics will be included.

**GetThetaMatrixArgs.topic\_index**

The list of topic indices, describing which topics to include in the Theta matrix. The values of this field should be an integers between 0 and (*ModelConfig.topics\_count* - 1). This field is optional, by default all topics will be included.

Note that this field acts similar to *GetThetaMatrixArgs.topic\_name*. It is not allowed to specify both *topic\_index* and *topic\_name* at the same time. The recommendation is to use *topic\_name*.

**GetThetaMatrixArgs.clean\_cache**

An optional flag that defines whether to clear the theta matrix cache after this operation. Setting this value to *True* will clear the cache for a topic model, defined by *GetThetaMatrixArgs.model\_name*. This value is only applicable when *MasterComponentConfig.cache\_theta* is set to *True*.

**GetThetaMatrixArgs.use\_sparse\_format**

An optional flag that defines whether to use sparse format for the resulting ThetaMatrix message. See ThetaMatrix message for additional information about the sparse format.

**GetThetaMatrixArgs.eps**

A small value that defines zero threshold for  $p(t|d)$  probabilities. This field is only used in sparse format.  $p(t|d)$  below the threshold will be excluded from the resulting Theta matrix.

## GetScoreValueArgs

Represents an argument of get score operation.

```
message GetScoreValueArgs {
  optional string model_name = 1;
  optional string score_name = 2;
  optional Batch batch = 3;
}
```

**GetScoreValueArgs.model\_name**

The name of the model to retrieved score for.

**GetScoreValueArgs.score\_name**

The name of the score to retrieved.

**GetScoreValueArgs.batch**

The *Batch* to calculate the score. This option is only applicable to cumulative scores. When not provided the score will be reported for all batches processed since last *ArtmInvokeIteration()*.

## AddBatchArgs

Represents an argument of *ArtmAddBatch()* operation.

```
message AddBatchArgs {
  optional Batch batch = 1;
  optional int32 timeout_milliseconds = 2 [default = -1];
  optional bool reset_scores = 3 [default = false];
  optional string batch_file_name = 4;
}
```

**AddBatchArgs.batch**

The *Batch* to add.

AddBatchArgs.**timeout\_milliseconds**

Timeout in milliseconds for this operation.

AddBatchArgs.**reset\_scores**

An optional flag that defines whether to reset all scores before this operation.

AddBatchArgs.**batch\_file\_name**

An optional value that defines disk location of the batch to add. You must choose between parameters *batch\_file\_name* or *batch* (either of them has to be specified, but not both at the same time).

## InvokeIterationArgs

Represents an argument of `ArtmInvokeIteration()` operation.

```
message InvokeIterationArgs {
  optional int32 iterations_count = 1 [default = 1];
  optional bool reset_scores = 2 [default = true];
  optional string disk_path = 3;
}
```

InvokeIterationArgs.**iterations\_count**

An integer value describing how many iterations to invoke.

InvokeIterationArgs.**reset\_scores**

An optional flag that defines whether to reset all scores before this operation.

InvokeIterationArgs.**disk\_path**

A value that defines the disk location with batches to process on this iteration.

## WaitIdleArgs

Represents an argument of `ArtmWaitIdle()` operation.

```
message WaitIdleArgs {
  optional int32 timeout_milliseconds = 1 [default = -1];
}
```

WaitIdleArgs.**timeout\_milliseconds**

Timeout in milliseconds for this operation.

## ExportModelArgs

Represents an argument of `ArtmExportModel()` operation.

```
message ExportModelArgs {
  optional string file_name = 1;
  optional string model_name = 2;
}
```

ExportModelArgs.**file\_name**

A target file name where to store topic model.

**ExportModelArgs.model\_name**

A value that describes the name of the topic model. This name will match the name of the corresponding model config.

## ImportModelArgs

Represents an argument of `ArtmImportModel()` operation.

```
message ImportModelArgs {
  optional string file_name = 1;
  optional string model_name = 2;
}
```

**ImportModelArgs.file\_name**

A target file name from where to load topic model.

**ImportModelArgs.model\_name**

A value that describes the name of the topic model. This name will match the name of the corresponding model config.

## C++ interface

BigARTM C++ interface is currently not documented. The main entry point is `MasterModel` class from `src/artm/cpp_interface.cc`. Please refer to `src/bigartm/srcmain.cc` for usage examples, and ask questions at [bigartm-users](#) or open a new [issue](#).

```
class MasterModel {
public:
  explicit MasterModel(const MasterModelConfig& config);
  ~MasterModel();

  int id() const { return id_; }
  MasterComponentInfo info() const; // misc. diagnostics information

  const MasterModelConfig& config() const { return config_; }
  MasterModelConfig* mutable_config() { return &config_; }
  void Reconfigure(); // apply MasterModel::config()

  // Operations to work with dictionary through disk
  void GatherDictionary(const GatherDictionaryArgs& args);
  void FilterDictionary(const FilterDictionaryArgs& args);
  void ImportDictionary(const ImportDictionaryArgs& args);
  void ExportDictionary(const ExportDictionaryArgs& args);
  void DisposeDictionary(const std::string& dictionary_name);

  // Operations to work with dictionary through memory
  void CreateDictionary(const DictionaryData& args);
  DictionaryData GetDictionary(const GetDictionaryArgs& args);

  // Operations to work with batches through memory
  void ImportBatches(const ImportBatchesArgs& args);
  void DisposeBatch(const std::string& batch_name);

  // Operations to work with model
```

```

void InitializeModel(const InitializeModelArgs& args);
void ImportModel(const ImportModelArgs& args);
void ExportModel(const ExportModelArgs& args);
void FitOnlineModel(const FitOnlineMasterModelArgs& args);
void FitOfflineModel(const FitOfflineMasterModelArgs& args);

// Apply model to batches
ThetaMatrix Transform(const TransformMasterModelArgs& args);
ThetaMatrix Transform(const TransformMasterModelArgs& args, Matrix* matrix);

// Retrieve operations
TopicModel GetTopicModel(const GetTopicModelArgs& args);
TopicModel GetTopicModel(const GetTopicModelArgs& args, Matrix* matrix);
ThetaMatrix GetThetaMatrix(const GetThetaMatrixArgs& args);
ThetaMatrix GetThetaMatrix(const GetThetaMatrixArgs& args, Matrix* matrix);

// Retrieve scores
ScoreData GetScore(const GetScoreValueArgs& args);
template <typename T>
T GetScoreAs(const GetScoreValueArgs& args);

```

**Warning:** What follows below in this page is really outdated.

In addition to this page consider to look at [Low-level API in C](#), [python\\_interface](#) or [Messages](#). These documentation files are also to certain degree relevant for C++ interface, because C++ interface is quite similar to Python interface and share the same Protobuf messages.

## MasterComponent

class **MasterComponent**

**MasterComponent** (const MasterComponentConfig &config)

Creates a master component with configuration defined by *MasterComponentConfig* message.

void **Reconfigure** (const MasterComponentConfig &config)

Updates the configuration of the master component.

const MasterComponentConfig &**config** () const

Returns current configuration of the master component.

MasterComponentConfig \***mutable\_config** ()

Returns mutable configuration of the master component. Remember to call *Reconfigure()* to propagate your changes to master component.

void **InvokeIteration** (int iterations\_count = 1)

Invokes certain number of iterations.

bool **AddBatch** (const Batch &batch, bool reset\_scores)

Adds batch to the processing queue.

bool **WaitIdle** (int timeout = -1)

Waits for iterations to be completed. Returns true if BigARTM completed before the specific timeout, otherwise false.

std::shared\_ptr<TopicModel> **GetTopicModel** (const std::string &model\_name)

Retrieves Phi matrix of a specific topic model. The resulting message *TopicModel* will contain information

about token weights distribution across topics.

std::shared\_ptr<TopicModel> **GetTopicModel** (const GetTopicModelArgs &args)

Retrieves Phi matrix based on extended parameters, specified in *GetTopicModelArgs* message. The resulting message *TopicModel* will contain information about token weights distribution across topics.

std::shared\_ptr<ThetaMatrix> **GetThetaMatrix** (const std::string &model\_name)

Retrieves Theta matrix of a specific topic model. The resulting message *ThetaMatrix* will contain information about items distribution across topics. Remember to set *MasterComponentConfig.cache\_theta* prior to the last iteration in order to gather Theta matrix.

std::shared\_ptr<ThetaMatrix> **GetThetaMatrix** (const GetThetaMatrixArgs &args)

Retrieves Theta matrix based on extended parameters, specified in *GetThetaMatrixArgs* message. The resulting message *ThetaMatrix* will contain information about items distribution across topics.

std::shared\_ptr<T> **GetScoreAs**<T> (const *Model* &model, const std::string &score\_name)

Retrieves given score for a specific model. Template argument must match the specific *ScoreData* type of the score (for example, *PerplexityScore*).

## Model

class **Model**

**Model** (const *MasterComponent* &master\_component, const ModelConfig &config)

Creates a topic model defined by *ModelConfig* inside given *MasterComponent*.

void **Reconfigure** (const ModelConfig &config)

Updates the configuration of the model.

const std::string &**name** () const

Returns the name of the model.

const ModelConfig &**config** () const

Returns current configuration of the model.

ModelConfig \***mutable\_config** ()

Returns mutable configuration of the model. Remember to call *Reconfigure()* to propagate your changes to the model.

void **Overwrite** (const TopicModel &topic\_model, bool commit = true)

Updates the model with new Phi matrix, defined by *topic\_model*. This operation can be used to provide an explicit initial approximation of the topic model, or to adjust the model in between iterations.

Depending on the *commit* flag the change can be applied immediately (*commit = true*) or queued (*commit = false*). The default setting is to use *commit = true*. You may want to use *commit = false* if your model is too big to be updated in a single protobuf message. In this case you should split your model into parts, each part containing subset of all tokens, and then submit each part in separate Overwrite operation with *commit = false*. After that remember to call *MasterComponent::WaitIdle()* and *Synchronize()* to propagate your change.

void **Initialize** (const *Dictionary* &dictionary)

Initialize topic model based on the *Dictionary*. Each token from the dictionary will be included in the model with randomly generated weight.

void **Export** (const string &file\_name)

Exports topic model into a file.

void **Import** (const string &file\_name)  
Imports topic model from a file.

void **Synchronize** (double decay\_weight, double apply\_weight, bool invoke\_regularizers)  
Synchronize the model.

This operation updates the Phi matrix of the topic model with all model increments, collected since the last call to `Synchronize()` method. The weights in the Phi matrix are set according to `decay_weight` and `apply_weight` values (refer to `SynchronizeModelArgs.decay_weight` for more details). Depending on `invoke_regularizers` parameter this operation may also invoke all regularizers.

Remember to call `Model::Synchronize()` operation every time after calling `MasterComponent::WaitIdle()`.

void **Synchronize** (const SynchronizeModelArgs &args)  
Synchronize the model based on extended arguments `SynchronizeModelArgs`.

## Regularizer

class **Regularizer**

**Regularizer** (const `MasterComponent` &master\_component, const RegularizerConfig &config)  
Creates a regularizer defined by `RegularizerConfig` inside given `MasterComponent`.

void **Reconfigure** (const RegularizerConfig &config)  
Updates the configuration of the regularizer.

const RegularizerConfig &**config** () const  
Returns current configuration of the regularizer.

RegularizerConfig \***mutable\_config** ()  
Returns mutable configuration of the regularizer. Remember to call `Reconfigure()` to propagate your changes to the regularizer.

## Dictionary

class **Dictionary**

**Dictionary** (const `MasterComponent` &master\_component, const DictionaryConfig &config)  
Creates a dictionary defined by `DictionaryConfig` inside given `MasterComponent`.

void **Reconfigure** (const DictionaryConfig &config)  
Updates the configuration of the dictionary.

const std::string **name** () const  
Returns the name of the dictionary.

const DictionaryConfig &**config** () const  
Returns current configuration of the dictionary.

## Utility methods

void **SaveBatch** (const Batch &batch, const std::string &disk\_path)  
Saves `Batch` into a specific folder. The name of the resulting file will be autogenerated, and the extension set to `.batch`

`std::shared_ptr<DictionaryConfig> LoadDictionary (const std::string &filename)`

Loads the *DictionaryConfig* message from a specific file on disk. *filename* must represent full disk path to the dictionary file.

`std::shared_ptr<Batch> LoadBatch (const std::string &filename)`

Loads the *Batch* message from a specific file on disk. *filename* must represent full disk path to the batch file, including *.batch* extension.

`std::shared_ptr<DictionaryConfig> ParseCollection (const CollectionParserConfig &config)`

Parses a text collection as defined by *CollectionParserConfig* message. Returns an instance of *DictionaryConfig* which carry all unique words in the collection and their frequencies.

## Windows distribution

This chapter describes content of BigARTM distribution package for Windows, available at <https://github.com/bigartm/bigartm/releases>.



bin/	Precompiled binaries of BigARTM for Windows. This folder must be added to PATH system variable.
bin/artm.dll	Core functionality of the BigARTM library.
bin/cpp_client.exe	Command line utility allows to perform simple experiments with BigARTM. Remember that not all BigARTM features are available through cpp_client, but it can serve as a good starting point to learn basic functionality. For further details refer to /ref/cpp_client.
protobuf/	A minimalistic version of Google Protocol Buffers ( <a href="https://code.google.com/p/protobuf/">https://code.google.com/p/protobuf/</a> ) library, required to run BigARTM from Python. To setup this package follow the instructions in protobuf/python/README file.
python/artm/	Python programming interface to BigARTM library. This folder must be added to PYTHONPATH system variable.
library.py	Implements all classes of BigARTM python interface.
messages_pb2.py	Contains all protobuf messages that can be transferred in and out BigARTM core library. Most common features are exposed with their own API methods, so normally you do not use python protobuf messages to operate BigARTM.
python/examples/	Python examples of how to use BigARTM:  Files docword.kos.txt and vocab.kos.txt represent a simple collection of text files in Bag-Of-Words format.  The files are taken from UCI Machine Learning Repository
<b>9.8. Windows distribution</b>	( <a href="https://archive.ics.uci.edu/ml/datasets/Bag+of+Words">https://archive.ics.uci.edu/ml/datasets/Bag+of+Words</a> ). <b>125</b>
src/	Source code and scripts of BigARTM library.



## a

`artm`, [43](#)

`artm.score_tracker`, [45](#)



## Symbols

- `__init__()` (artm.ARTM method), 25
  - `__init__()` (artm.BackgroundTokensRatioScore method), 45
  - `__init__()` (artm.BatchVectorizer method), 36
  - `__init__()` (artm.ClassPrecisionScore method), 45
  - `__init__()` (artm.DecorrelatorPhiRegularizer method), 40
  - `__init__()` (artm.Dictionary method), 37
  - `__init__()` (artm.ImproveCoherencePhiRegularizer method), 42
  - `__init__()` (artm.ItemsProcessedScore method), 43
  - `__init__()` (artm.KIFunctionInfo method), 39
  - `__init__()` (artm.LDA method), 30
  - `__init__()` (artm.LabelRegularizationPhiRegularizer method), 41
  - `__init__()` (artm.MasterComponent method), 49
  - `__init__()` (artm.PerplexityScore method), 43
  - `__init__()` (artm.SmoothPtdwRegularizer method), 42
  - `__init__()` (artm.SmoothSparsePhiRegularizer method), 39
  - `__init__()` (artm.SmoothSparseThetaRegularizer method), 40
  - `__init__()` (artm.SparsityPhiScore method), 43
  - `__init__()` (artm.SparsityThetaScore method), 44
  - `__init__()` (artm.SpecifiedSparsePhiRegularizer method), 41
  - `__init__()` (artm.ThetaSnippetScore method), 44
  - `__init__()` (artm.TopTokensScore method), 44
  - `__init__()` (artm.TopicKernelScore method), 44
  - `__init__()` (artm.TopicMassPhiScore method), 45
  - `__init__()` (artm.TopicSegmentationPtdwRegularizer method), 43
  - `__init__()` (artm.TopicSelectionThetaRegularizer method), 42
  - `__init__()` (artm.hARTM method), 33
  - `__init__()` (artm.score\_tracker.BackgroundTokensRatioScoreTracker method), 48
  - `__init__()` (artm.score\_tracker.ClassPrecisionScoreTracker method), 48
  - `__init__()` (artm.score\_tracker.ItemsProcessedScoreTracker method), 47
  - `__init__()` (artm.score\_tracker.PerplexityScoreTracker method), 46
  - `__init__()` (artm.score\_tracker.SparsityPhiScoreTracker method), 45
  - `__init__()` (artm.score\_tracker.SparsityThetaScoreTracker method), 46
  - `__init__()` (artm.score\_tracker.ThetaSnippetScoreTracker method), 47
  - `__init__()` (artm.score\_tracker.TopTokensScoreTracker method), 46
  - `__init__()` (artm.score\_tracker.TopicKernelScoreTracker method), 47
  - `__init__()` (artm.score\_tracker.TopicMassPhiScoreTracker method), 48
- ## A
- `add_level()` (artm.hARTM method), 34
  - `alpha_iter` (SmoothSparseThetaConfig attribute), 98
  - `apply_weight` (SynchronizeModelArgs attribute), 115
  - ARTM (class in artm), 25
  - artm (module), 25, 30, 33, 36, 37, 39, 43, 48
  - artm.score\_tracker (module), 45
  - artm::Dictionary (C++ class), 123
  - artm::Dictionary::config (C++ function), 123
  - artm::Dictionary::Dictionary (C++ function), 123
  - artm::Dictionary::name (C++ function), 123
  - artm::Dictionary::Reconfigure (C++ function), 123
  - artm::LoadBatch (C++ function), 124
  - artm::LoadDictionary (C++ function), 123
  - artm::MasterComponent (C++ class), 121
  - artm::MasterComponent::AddBatch (C++ function), 121
  - artm::MasterComponent::config (C++ function), 121
  - artm::MasterComponent::GetScoreAs<T> (C++ function), 122
  - artm::MasterComponent::GetThetaMatrix (C++ function), 122
  - artm::MasterComponent::GetTopicModel (C++ function), 121, 122
  - artm::MasterComponent::InvokeIteration (C++ function), 121

artm::MasterComponent::MasterComponent (C++ function), 121

artm::MasterComponent::mutable\_config (C++ function), 121

artm::MasterComponent::Reconfigure (C++ function), 121

artm::MasterComponent::WaitIdle (C++ function), 121

artm::Model (C++ class), 122

artm::Model::config (C++ function), 122

artm::Model::Export (C++ function), 122

artm::Model::Import (C++ function), 122

artm::Model::Initialize (C++ function), 122

artm::Model::Model (C++ function), 122

artm::Model::mutable\_config (C++ function), 122

artm::Model::name (C++ function), 122

artm::Model::Overwrite (C++ function), 122

artm::Model::Reconfigure (C++ function), 122

artm::Model::Synchronize (C++ function), 123

artm::ParseCollection (C++ function), 124

artm::Regularizer (C++ class), 123

artm::Regularizer::config (C++ function), 123

artm::Regularizer::mutable\_config (C++ function), 123

artm::Regularizer::Reconfigure (C++ function), 123

artm::Regularizer::Regularizer (C++ function), 123

artm::SaveBatch (C++ function), 123

ARTM\_ARGUMENT\_OUT\_OF\_RANGE (C macro), 23

ARTM\_CORRUPTED\_MESSAGE (C macro), 23

ARTM\_DISK\_READ\_ERROR (C macro), 23

ARTM\_DISK\_WRITE\_ERROR (C macro), 23

ARTM\_INTERNAL\_ERROR (C macro), 23

ARTM\_INVALID\_MASTER\_ID (C macro), 23

ARTM\_INVALID\_OPERATION (C macro), 23

ARTM\_STILL\_WORKING (C macro), 23

ARTM\_SUCCESS (C macro), 23

ArtemGetLastError (C function), 22

attach\_model() (artm.MasterComponent method), 49

average\_kernel\_contrast (TopicKernelScore attribute), 109

average\_kernel\_purity (TopicKernelScore attribute), 109

average\_kernel\_size (TopicKernelScore attribute), 109

## B

BackgroundTokensRatioScore (class in artm), 45

BackgroundTokensRatioScoreTracker (class in artm.score\_tracker), 48

batch (AddBatchArgs attribute), 118

batch (GetScoreValueArgs attribute), 118

batch (GetThetaMatrixArgs attribute), 117

batch\_file\_name (AddBatchArgs attribute), 119

batch\_size (artm.BatchVectorizer attribute), 37

batches\_list (artm.BatchVectorizer attribute), 37

BatchVectorizer (class in artm), 36

## C

cache\_theta (MasterComponentConfig attribute), 95

class\_id (Batch attribute), 93

class\_id (DecorrelatorPhiConfig attribute), 99

class\_id (DictionaryEntry attribute), 100

class\_id (GetTopicModelArgs attribute), 116

class\_id (LabelRegularizationPhiConfig attribute), 99

class\_id (ModelConfig attribute), 96

class\_id (SmoothSparsePhiConfig attribute), 98

class\_id (SparsityPhiScoreConfig attribute), 104

class\_id (TopicKernelScoreConfig attribute), 108

class\_id (TopicModel attribute), 110

class\_id (TopTokensScoreConfig attribute), 106

class\_weight (ModelConfig attribute), 96

ClassPrecisionScore (class in artm), 45

ClassPrecisionScoreTracker (class in artm.score\_tracker), 48

clean\_cache (GetThetaMatrixArgs attribute), 118

clear\_score\_array\_cache() (artm.MasterComponent method), 49

clear\_score\_cache() (artm.MasterComponent method), 49

clear\_theta\_cache() (artm.MasterComponent method), 49

compact\_batches (MasterComponentConfig attribute), 94

config (RegularizerConfig attribute), 97

config (ScoreConfig attribute), 101

cooccurrence\_file\_name (CollectionParserConfig attribute), 114

cooccurrence\_token (CollectionParserConfig attribute), 114

copy() (artm.Dictionary method), 38

create() (artm.Dictionary method), 38

create\_dictionary() (artm.MasterComponent method), 49

create\_regularizer() (artm.MasterComponent method), 49

create\_score() (artm.MasterComponent method), 50

## D

data (ScoreData attribute), 102

data\_path (artm.BatchVectorizer attribute), 37

decay\_weight (SynchronizerModelArgs attribute), 115

DecorrelatorPhiRegularizer (class in artm), 40

del\_level() (artm.hARTM method), 34

description (Batch attribute), 93

dictionary (artm.BatchVectorizer attribute), 37

Dictionary (class in artm), 37

dictionary\_file\_name (CollectionParserConfig attribute), 114

dictionary\_name (InitializeModelArgs attribute), 116

dictionary\_name (LabelRegularizationPhiConfig attribute), 99

dictionary\_name (SmoothSparsePhiConfig attribute), 98

disk\_cache\_path (MasterComponentConfig attribute), 95

disk\_path (InvokeIterationArgs attribute), 119

disk\_path (MasterComponentConfig attribute), 94

dispose() (artm.ARTM method), 26  
 dispose() (artm.hARTM method), 34  
 docword\_file\_path (CollectionParserConfig attribute), 113

## E

enabled (ModelConfig attribute), 96  
 entry (DictionaryConfig attribute), 100  
 eps (GetThetaMatrixArgs attribute), 118  
 eps (GetTopicModelArgs attribute), 117  
 eps (SparsityPhiScoreConfig attribute), 104  
 eps (SparsityThetaScoreConfig attribute), 103  
 eps (TopicKernelScoreConfig attribute), 108  
 export\_dictionary() (artm.MasterComponent method), 50  
 export\_model() (artm.MasterComponent method), 50

## F

field (Item attribute), 92  
 field\_name (ItemsProcessedScoreConfig attribute), 105  
 field\_name (ModelConfig attribute), 96  
 field\_name (PerplexityScoreConfig attribute), 102  
 field\_name (SparsityThetaScoreConfig attribute), 103  
 field\_name (ThetaSnippetScoreConfig attribute), 107  
 file\_name (ExportModelArgs attribute), 119  
 file\_name (ImportModelArgs attribute), 120  
 filter() (artm.Dictionary method), 38  
 filter\_dictionary() (artm.MasterComponent method), 50  
 fit\_offline() (artm.ARTM method), 26  
 fit\_offline() (artm.hARTM method), 34  
 fit\_offline() (artm.LDA method), 31  
 fit\_offline() (artm.MasterComponent method), 50  
 fit\_online() (artm.ARTM method), 26  
 fit\_online() (artm.LDA method), 31  
 fit\_online() (artm.MasterComponent method), 50  
 format (CollectionParserConfig attribute), 112

## G

gather() (artm.Dictionary method), 38  
 gather\_dictionary() (artm.MasterComponent method), 51  
 get\_dictionary() (artm.MasterComponent method), 51  
 get\_info() (artm.MasterComponent method), 51  
 get\_level() (artm.hARTM method), 34  
 get\_phi() (artm.ARTM method), 27  
 get\_phi() (artm.hARTM method), 35  
 get\_phi\_info() (artm.MasterComponent method), 51  
 get\_phi\_matrix() (artm.MasterComponent method), 51  
 get\_phi\_sparse() (artm.ARTM method), 27  
 get\_score() (artm.ARTM method), 27  
 get\_score() (artm.MasterComponent method), 51  
 get\_score\_array() (artm.MasterComponent method), 51  
 get\_theta() (artm.ARTM method), 28  
 get\_theta() (artm.hARTM method), 35  
 get\_theta() (artm.LDA method), 31  
 get\_theta\_info() (artm.MasterComponent method), 52

get\_theta\_matrix() (artm.MasterComponent method), 52  
 get\_theta\_sparse() (artm.ARTM method), 28  
 get\_top\_tokens() (artm.LDA method), 31

## H

hARTM (class in artm), 33

## I

id (Batch attribute), 93  
 id (Item attribute), 92  
 import\_dictionary() (artm.MasterComponent method), 52  
 import\_model() (artm.MasterComponent method), 52  
 ImproveCoherencePhiRegularizer (class in artm), 42  
 info (artm.ARTM attribute), 28  
 initialize() (artm.ARTM method), 28  
 initialize() (artm.LDA method), 32  
 initialize\_model() (artm.MasterComponent method), 52  
 inner\_iterations\_count (ModelConfig attribute), 96  
 internals (TopicModel attribute), 110  
 invoke\_regularizers (SynchronizeModelArgs attribute), 115  
 item (Batch attribute), 93  
 item\_count (ThetaSnippetScoreConfig attribute), 107  
 item\_id (ThetaMatrix attribute), 111  
 item\_id (ThetaSnippetScore attribute), 107  
 item\_id (ThetaSnippetScoreConfig attribute), 107  
 item\_title (ThetaMatrix attribute), 111  
 item\_weights (ThetaMatrix attribute), 111  
 items\_count (DictionaryEntry attribute), 101  
 ItemsProcessedScore (class in artm), 43  
 ItemsProcessedScoreTracker (class in artm.score\_tracker), 47  
 iterations\_count (InvokeIterationArgs attribute), 119

## K

kernel\_contrast (TopicKernelScore attribute), 109  
 kernel\_purity (TopicKernelScore attribute), 108  
 kernel\_size (TopicKernelScore attribute), 108  
 key\_token (DictionaryEntry attribute), 100  
 KIFunctionInfo (class in artm), 39

## L

LabelRegularizationPhiRegularizer (class in artm), 41  
 LDA (class in artm), 30  
 library\_version (artm.ARTM attribute), 28  
 library\_version (artm.hARTM attribute), 35  
 load() (artm.ARTM method), 28  
 load() (artm.Dictionary method), 39  
 load() (artm.hARTM method), 35  
 load() (artm.LDA method), 32  
 load\_text() (artm.Dictionary method), 39

## M

MasterComponent (class in artm), 48

`merge_model()` (`artm.MasterComponent` method), 52  
`merger_queue_max_size` (`MasterComponentConfig` attribute), 95  
`messages_pb2.Batch` (built-in class), 93  
`messages_pb2.BoolArray` (built-in class), 91  
`messages_pb2.CollectionParserConfig` (built-in class), 112  
`messages_pb2.DecorrelatorPhiConfig` (built-in class), 98  
`messages_pb2.DictionaryConfig` (built-in class), 100  
`messages_pb2.DictionaryEntry` (built-in class), 100  
`messages_pb2.DoubleArray` (built-in class), 91  
`messages_pb2.Field` (built-in class), 92  
`messages_pb2.FloatArray` (built-in class), 91  
`messages_pb2.InitializeModelArgs` (built-in class), 115  
`messages_pb2.IntArray` (built-in class), 92  
`messages_pb2.Item` (built-in class), 92  
`messages_pb2.ItemsProcessedScore` (built-in class), 105  
`messages_pb2.ItemsProcessedScoreConfig` (built-in class), 105  
`messages_pb2.LabelRegularizationPhiConfig` (built-in class), 99  
`messages_pb2.MasterComponentConfig` (built-in class), 94  
`messages_pb2.ModelConfig` (built-in class), 95  
`messages_pb2.PerplexityScore` (built-in class), 102  
`messages_pb2.PerplexityScoreConfig` (built-in class), 102  
`messages_pb2.RegularizerConfig` (built-in class), 97  
`messages_pb2.RegularizerInternalState` (built-in class), 99  
`messages_pb2.ScoreConfig` (built-in class), 101  
`messages_pb2.ScoreData` (built-in class), 101  
`messages_pb2.SmoothSparsePhiConfig` (built-in class), 98  
`messages_pb2.SmoothSparseThetaConfig` (built-in class), 98  
`messages_pb2.SparsityPhiScore` (built-in class), 104  
`messages_pb2.SparsityPhiScoreConfig` (built-in class), 104  
`messages_pb2.SparsityThetaScoreConfig` (built-in class), 103, 104  
`messages_pb2.Stream` (built-in class), 93  
`messages_pb2.SynchronizeModelArgs` (built-in class), 115  
`messages_pb2.ThetaMatrix` (built-in class), 111  
`messages_pb2.ThetaSnippetScore` (built-in class), 107  
`messages_pb2.ThetaSnippetScoreConfig` (built-in class), 107  
`messages_pb2.TopicKernelScore` (built-in class), 108  
`messages_pb2.TopicKernelScoreConfig` (built-in class), 108  
`messages_pb2.TopicModel` (built-in class), 109  
`messages_pb2.TopTokensScore` (built-in class), 106

`messages_pb2.TopTokensScoreConfig` (built-in class), 105  
`model_name` (`ExportModelArgs` attribute), 119  
`model_name` (`GetScoreValueArgs` attribute), 118  
`model_name` (`GetThetaMatrixArgs` attribute), 117  
`model_name` (`GetTopicModelArgs` attribute), 116  
`model_name` (`ImportModelArgs` attribute), 120  
`model_name` (`InitializeModelArgs` attribute), 116  
`model_name` (`SynchronizeModelArgs` attribute), 115  
`model_name` (`ThetaMatrix` attribute), 111

## N

`name` (`DictionaryConfig` attribute), 100  
`name` (`ModelConfig` attribute), 96  
`name` (`RegularizerConfig` attribute), 97  
`name` (`ScoreConfig` attribute), 101  
`name` (`ScoreData` attribute), 102  
`name` (`Stream` attribute), 94  
`name` (`TopicModel` attribute), 109  
`normalize_model()` (`artm.MasterComponent` method), 52  
`normalizer` (`PerplexityScore` attribute), 103  
`num_batches` (`artm.BatchVectorizer` attribute), 37  
`num_entries` (`TopTokensScore` attribute), 106  
`num_items_per_batch` (`CollectionParserConfig` attribute), 114  
`num_tokens` (`TopTokensScoreConfig` attribute), 106

## O

`online_batch_processing` (`MasterComponentConfig` attribute), 95  
`operation_type` (`TopicModel` attribute), 110  
`opt_for_avx` (`ModelConfig` attribute), 97

## P

`PerplexityScore` (class in `artm`), 43  
`PerplexityScoreTracker` (class in `artm.score_tracker`), 46  
`probability_mass_threshold` (`TopicKernelScoreConfig` attribute), 108  
`process_batches()` (`artm.MasterComponent` method), 53  
`processor_queue_max_size` (`MasterComponentConfig` attribute), 95  
`processors_count` (`MasterComponentConfig` attribute), 95

## R

`raw` (`PerplexityScore` attribute), 103  
`reconfigure()` (`artm.MasterComponent` method), 53  
`reconfigure_regularizer()` (`artm.MasterComponent` method), 53  
`reconfigure_score()` (`artm.MasterComponent` method), 53  
`reconfigure_topic_name()` (`artm.MasterComponent` method), 53  
`regularize_model()` (`artm.MasterComponent` method), 53  
`regularizer_name` (`ModelConfig` attribute), 96



regularizer\_tau (ModelConfig attribute), 96  
 remove\_theta() (artm.ARTM method), 29  
 remove\_theta() (artm.LDA method), 32  
 request\_type (GetTopicModelArgs attribute), 117  
 reset\_scores (AddBatchArgs attribute), 119  
 reset\_scores (InvokeIterationArgs attribute), 119  
 reshape\_topics() (artm.ARTM method), 29  
 reuse\_theta (ModelConfig attribute), 96

## S

save() (artm.ARTM method), 29  
 save() (artm.Dictionary method), 39  
 save() (artm.hARTM method), 36  
 save() (artm.LDA method), 32  
 save\_text() (artm.Dictionary method), 39  
 score\_config (MasterComponentConfig attribute), 95  
 score\_name (GetScoreValueArgs attribute), 118  
 score\_name (ModelConfig attribute), 96  
 SmoothPtdwRegularizer (class in artm), 42  
 SmoothSparsePhiRegularizer (class in artm), 39  
 SmoothSparseThetaRegularizer (class in artm), 40  
 SparsityPhiScore (class in artm), 43  
 SparsityPhiScoreTracker (class in artm.score\_tracker), 45  
 SparsityThetaScore (class in artm), 44  
 SparsityThetaScoreTracker (class in artm.score\_tracker), 46  
 SpecifiedSparsePhiRegularizer (class in artm), 41  
 stream (MasterComponentConfig attribute), 94  
 stream\_name (ItemsProcessedScoreConfig attribute), 105  
 stream\_name (ModelConfig attribute), 96  
 stream\_name (PerplexityScoreConfig attribute), 102  
 stream\_name (SparsityThetaScoreConfig attribute), 103  
 stream\_name (ThetaSnippetScoreConfig attribute), 107

## T

target\_folder (CollectionParserConfig attribute), 114  
 theta\_sparsity\_value (PerplexityScore attribute), 103  
 ThetaSnippetScore (class in artm), 44  
 ThetaSnippetScoreTracker (class in artm.score\_tracker), 47  
 timeout\_milliseconds (AddBatchArgs attribute), 118  
 timeout\_milliseconds (WaitIdleArgs attribute), 119  
 title (Item attribute), 92  
 token (Batch attribute), 93  
 token (GetTopicModelArgs attribute), 116  
 token (TopicModel attribute), 110  
 token (TopTokensScore attribute), 106  
 token\_count (DictionaryEntry attribute), 100  
 token\_weights (TopicModel attribute), 110  
 topic\_index (GetThetaMatrixArgs attribute), 117  
 topic\_index (ThetaMatrix attribute), 111  
 topic\_index (TopicModel attribute), 110  
 topic\_index (TopTokensScore attribute), 107  
 topic\_name (DecorrelatorPhiConfig attribute), 99

topic\_name (GetThetaMatrixArgs attribute), 117  
 topic\_name (GetTopicModelArgs attribute), 116  
 topic\_name (LabelRegularizationPhiConfig attribute), 99  
 topic\_name (ModelConfig attribute), 96  
 topic\_name (SmoothSparsePhiConfig attribute), 98  
 topic\_name (SmoothSparseThetaConfig attribute), 98  
 topic\_name (SparsityPhiScoreConfig attribute), 104  
 topic\_name (SparsityThetaScoreConfig attribute), 103  
 topic\_name (ThetaMatrix attribute), 111  
 topic\_name (TopicKernelScoreConfig attribute), 108  
 topic\_name (TopicModel attribute), 110  
 topic\_name (TopTokensScore attribute), 107  
 topic\_name (TopTokensScoreConfig attribute), 106  
 topic\_names (artm.ARTM attribute), 29  
 TopicKernelScore (class in artm), 44  
 TopicKernelScoreTracker (class in artm.score\_tracker), 47  
 TopicMassPhiScore (class in artm), 45  
 TopicMassPhiScoreTracker (class in artm.score\_tracker), 48  
 topics\_count (ModelConfig attribute), 96  
 topics\_count (ThetaMatrix attribute), 111  
 topics\_count (TopicModel attribute), 110  
 TopicSegmentationPtdwRegularizer (class in artm), 43  
 TopicSelectionThetaRegularizer (class in artm), 42  
 TopTokensScore (class in artm), 44  
 TopTokensScoreTracker (class in artm.score\_tracker), 46  
 total\_items\_count (DictionaryConfig attribute), 100  
 total\_token\_count (DictionaryConfig attribute), 100  
 total\_tokens (SparsityPhiScore attribute), 105  
 total\_topics (SparsityThetaScore attribute), 104  
 transform() (artm.ARTM method), 29  
 transform() (artm.hARTM method), 36  
 transform() (artm.LDA method), 32  
 transform() (artm.MasterComponent method), 54  
 transform\_sparse() (artm.ARTM method), 30  
 type (RegularizerConfig attribute), 97  
 type (ScoreConfig attribute), 101  
 type (ScoreData attribute), 102  
 type (Stream attribute), 94

## U

use\_new\_tokens (ModelConfig attribute), 97  
 use\_random\_theta (ModelConfig attribute), 96  
 use\_sparse\_bow (ModelConfig attribute), 96  
 use\_sparse\_format (GetThetaMatrixArgs attribute), 118  
 use\_sparse\_format (GetTopicModelArgs attribute), 117  
 use\_unity\_based\_indices (CollectionParserConfig attribute), 114

## V

value (DictionaryEntry attribute), 100  
 value (ItemsProcessedScore attribute), 105  
 value (PerplexityScore attribute), 103

[value \(SparsityPhiScore attribute\), 105](#)  
[value \(SparsityThetaScore attribute\), 104](#)  
[values \(ThetaSnippetScore attribute\), 107](#)  
[vocab\\_file\\_path \(CollectionParserConfig attribute\), 114](#)

## W

[weight \(TopTokensScore attribute\), 106](#)  
[weights \(artm.BatchVectorizer attribute\), 37](#)

## Z

[zero\\_tokens \(SparsityPhiScore attribute\), 105](#)  
[zero\\_topics \(SparsityThetaScore attribute\), 104](#)  
[zero\\_words \(PerplexityScore attribute\), 103](#)