

---

# Cornac Documentation

*Release 0.1.0.post5*

## Cornac Contributors

Mar 31, 2019



---

## Getting Started

---

<b>1 Installation</b>	<b>3</b>
<b>2 First example</b>	<b>5</b>
<b>3 Data</b>	<b>7</b>
<b>4 Models</b>	<b>21</b>
<b>5 Metrics</b>	<b>39</b>
<b>6 Evaluation methods</b>	<b>43</b>
<b>7 Experiment</b>	<b>47</b>
<b>8 Built-in datasets</b>	<b>49</b>
<b>Python Module Index</b>	<b>53</b>



**Cornac** is python recommender systems library for **easy**, **effective** and **efficient** experiments. Cornac is **simple** and **handy**. It is designed from the ground-up to faithfully reflect the standard steps taken by researchers to implement and evaluate personalized recommendation models. Moreover, contributing new recommender models, evaluation metrics, etc., to Cornac is very easy and smooth. For instance, if you already have a python implementation of your model, e.g., PMF, you will need to spend less than 5 minutes in average to integrate it to Cornac.



# CHAPTER 1

---

## Installation

---

Currently, we are supporting Python 3 (version 3.6 is recommended). There are several ways to install Cornac:

- **From PyPI (you may need a C++ compiler):**

```
pip3 install cornac
```

- **From Anaconda:**

```
conda install cornac -c qttruong -c pytorch
```

- **From the GitHub source (for latest updates):**

```
pip3 install Cython
git clone https://github.com/PreferredAI/cornac.git
cd cornac
python3 setup.py install
```

### Note:

Additional dependencies required by models are listed [here](#).

Some of the algorithms use *OpenMP* to support multi-threading. For OSX users, in order to run those algorithms efficiently, you might need to install *gcc* from Homebrew to have an OpenMP compiler:

```
brew install gcc | brew link gcc
```

If you want to utilize your GPUs, you might consider:

- TensorFlow installation instructions: <https://www.tensorflow.org/install/>
- PyTorch installation instructions: <https://pytorch.org/get-started/locally/>
- cuDNN: <https://docs.nvidia.com/deeplearning/sdk/cudnn-install/> (for Nvidia GPUs)



# CHAPTER 2

---

## First example

---

This example will show you how to run your very first experiment using Cornac.

```
import cornac as cn

# Load MovieLens 100K dataset
ml_100k = cn.datasets.movieLens.load_100k()

# Split data based on ratio
ratio_split = cn.eval_methods.RatioSplit(data=ml_100k,
                                         test_size=0.2,
                                         rating_threshold=4.0,
                                         seed=123)

# Here we are comparing: Biased MF, PMF, and BPR
mf = cn.models.MF(k=10, max_iter=25, learning_rate=0.01, lambda_reg=0.02, use_
    ↪bias=True)
pmf = cn.models.PMF(k=10, max_iter=100, learning_rate=0.001, lamda=0.001)
bpr = cn.models.BPR(k=10, max_iter=200, learning_rate=0.01, lambda_reg=0.01)

# Define metrics used to evaluate the models
mae = cn.metrics.MAE()
rmse = cn.metrics.RMSE()
rec_20 = cn.metrics.Recall(k=20)
ndcg_20 = cn.metrics.NDCG(k=20)
auc = cn.metrics.AUC()

# Put it together into an experiment and run
exp = cn.Experiment(eval_method=ratio_split,
                     models=[mf, pmf, bpr],
                     metrics=[mae, rmse, rec_20, ndcg_20, auc],
                     user_based=True)
exp.run()
```

Output:

	MAE	RMSE	Recall@20	NDCG@20	AUC	Train (s)	Test (s)
<b>MF</b>	0.7441	0.9007	0.0622	0.0534	0.2952	0.0736	8.1187
<b>PMF</b>	0.7493	0.9084	0.0835	0.0673	0.4749	358.7642	8.4184
<b>BPR</b>	1.5595	1.8864	0.0744	0.0657	0.5932	2.5395	8.5734

# CHAPTER 3

---

## Data

---

```
class cornac.data.FeatureModule(features=None, ids=None, normalized=False, **kwargs)
```

### Parameters

- **features** (*numpy.ndarray* or *scipy.sparse.csr\_matrix*, default = *None*) – Numpy 2d-array that the row indices are aligned with user/item in *ids*.
- **ids** (*List*, default = *None*) – List of user/item ids that the indices are aligned with *corpus*. If *None*, the indices of provided *features* will be used as *ids*.

**batch\_feature** (*batch\_ids*)

Return a matrix (batch of feature vectors) corresponding to provided *batch\_ids*

**build** (*id\_map*=*None*)

Build the feature matrix. Features will be swapped if the *id\_map* is provided

**feature\_dim**

Return the dimensionality of the feature vectors

**features**

Return the whole feature matrix

```
class cornac.data.TextModule(corpus: List[str] = None, ids: List = None, tokenizer: cornac.data.text.Tokenizer = None, vocab: cornac.data.text.Vocabulary = None, max_vocab: int = None, max_doc_freq: Union[float, int] = 1.0, min_freq: int = 1, stop_words: Union[List, str] = None, **kwargs)
```

Text module

### Parameters

- **corpus** (*List* [*str*]), default = *None*) – List of user/item texts that the indices are aligned with *ids*.
- **ids** (*List*, default = *None*) – List of user/item ids that the indices are aligned with *corpus*. If *None*, the indices of provided *corpus* will be used as *ids*.

- **tokenizer** (`Tokenizer`, *optional*, `default = None`) – Tokenizer for text splitting. If None, the `BaseTokenizer` will be used.
- **vocab** (`Vocabulary`, *optional*, `default = None`) – Vocabulary of tokens. It contains mapping between tokens to their integer ids and vice versa.
- **max\_vocab** (`int`, *optional*, `default = None`) – The maximum size of the vocabulary. If vocab is provided, this will be ignored.
- **max\_doc\_freq** (`Union[float, int]` = `1.0`) – The maximum frequency of tokens appearing in documents to be excluded from vocabulary. If float, the value represents a proportion of documents, int for absolute counts. If `vocab` is not None, this will be ignored.
- **min\_freq** (`int`, `default = 1`) – The minimum frequency of tokens to be included into vocabulary. If `vocab` is not None, this will be ignored.
- **stop\_words** (`Collection, str, default: None`) – Collection of stop words which will be ignored when building `Vocabulary`. If str, it indicates a built-in stop words list. Currently, only `english` is supported.

**batch\_seq**(`batch_ids, max_length=None`)

Return a numpy matrix of text sequences containing token ids with `size=(len(batch_ids), max_length)`. If `max_length=None`, it will be inferred based on retrieved sequences.

**batch\_tfidf**(`batch_ids`)

Return matrix of TF-IDF features corresponding to provided `batch_ids`

**build**(`id_map=None`)

Build the model based on provided list of ordered ids

**class** `cornac.data.ImageModule(**kwargs)`

Image module

**batch\_image**(`batch_ids, target_size=(256, 256), color_mode='rgb', interpolation='nearest'`)

Return batch of images corresponding to provided `batch_ids`

**build**(`id_map=None`)

Build the model based on provided list of ordered ids

**class** `cornac.data.GraphModule(**kwargs)`

Graph module

**batch**(`batch_ids`)

Return batch of vectors from the sparse adjacency matrix corresponding to provided `batch_ids`.

**Parameters** **batch\_ids** (`array, required`) – An array contains the ids of rows to be returned from the sparse adjacency matrix.

**build**(`id_map=None`)

Build the feature matrix. Features will be swapped if the `id_map` is provided

**get\_train\_triplet**(`train_row_ids, train_col_ids`)

Get the training tuples

**class** `cornac.data.TrainSet(uid_map, iid_map)`

Training Set

**Parameters**

- **uid\_map** (`defaultdict`) – The dictionary containing mapping from original ids to mapped ids of users.
- **iid\_map** (`defaultdict`) – The dictionary containing mapping from original ids to mapped ids of items.

```
get_iid(raw_iid)
    Return the mapped id of an item given a raw id

get_uid(raw_uid)
    Return the mapped id of a user given a raw id

static idx_iter(idx_range, batch_size=1, shuffle=False)
    Create an iterator over batch of indices

Parameters

- batch_size (int, optional, default = 1) –
- shuffle (bool, optional) – If True, orders of triplets will be randomized. If False, default orders kept

Returns iterator

Return type batch of indices (array of np.int)

iid_list
    Return the list of mapped item ids

is_unk_item(mapped_iid)
    Return whether or not an item is unknown given the mapped id

is_unk_user(mapped_uid)
    Return whether or not a user is unknown given the mapped id

num_items
    Return the number of items

num_users
    Return the number of users

raw_iid_list
    Return the list of raw item ids

raw_uid_list
    Return the list of raw user ids

uid_list
    Return the list of mapped user ids

class cornac.data.MatrixTrainSet(matrix, max_rating, min_rating, global_mean, uid_map, iid_map)
    Training set contains preference matrix

Parameters

- matrix (scipy.sparse.csr_matrix) – Preferences in the form of scipy sparse matrix.
- max_rating (float) – Maximum value of the preferences.
- min_rating (float) – Minimum value of the preferences.
- global_mean (float) – Average value of the preferences.
- uid_map (defaultdict) – The dictionary containing mapping from original ids to mapped ids of users.
- iid_map (defaultdict) – The dictionary containing mapping from original ids to mapped ids of items.

```

```
classmethod from_uir(data,           global_uid_map=None,           global_iid_map=None,
                      global_ui_set=None, verbose=False)
Constructing TrainSet from triplet data.
```

**Parameters**

- **data** (*array-like, shape: [n\_examples, 3]*) – Data in the form of triplets (user, item, rating)
- **global\_uid\_map** (`defaultdict`, optional, default: None) – The dictionary containing global mapping from original ids to mapped ids of users.
- **global\_iid\_map** (`defaultdict`, optional, default: None) – The dictionary containing global mapping from original ids to mapped ids of items.
- **global\_ui\_set** (`set`, optional, default: None) – The global set of tuples (user, item). This helps avoiding duplicate observations.
- **verbose** (`bool`, default: `False`) – The verbosity flag.

**Returns** `train_set` – MatrixTrainSet object.

**Return type** <`cornac.data.MatrixTrainSet`>

**item\_iter** (*batch\_size=1, shuffle=False*)

Create an iterator over item ids

**Parameters**

- **batch\_size** (`int`, optional, default = 1) –
- **shuffle** (`bool`, optional) – If True, orders of triplets will be randomized. If False, default orders kept

**Returns** `iterator`

**Return type** batch of item ids (array of `np.int`)

**uij\_iter** (*batch\_size=1, shuffle=False*)

Create an iterator over data yielding batch of users, positive items, and negative items

**Parameters**

- **batch\_size** (`int`, optional, default = 1) –
- **shuffle** (`bool`, optional) – If True, orders of triplets will be randomized. If False, default orders kept

**Returns** `iterator` – batch of negative items (array of `np.int`)

**Return type** batch of users (array of `np.int`), batch of positive items (array of `np.int`),

**uir\_iter** (*batch\_size=1, shuffle=False*)

Create an iterator over data yielding batch of users, items, and rating values

**Parameters**

- **batch\_size** (`int`, optional, default = 1) –
- **shuffle** (`bool`, optional) – If True, orders of triplets will be randomized. If False, default orders kept

**Returns** `iterator` – batch of ratings (array of `np.float`)

**Return type** batch of users (array of `np.int`), batch of items (array of `np.int`),

**user\_iter** (*batch\_size=1, shuffle=False*)

Create an iterator over user ids

**Parameters**

- **batch\_size** (`int`, optional, default = 1) –
- **shuffle** (`bool`, optional) – If True, orders of triplets will be randomized. If False, default orders kept

**Returns iterator****Return type** batch of user ids (array of np.int)

```
class cornac.data.MultimodalTrainSet(matrix, max_rating, min_rating, global_mean,
                                      uid_map, iid_map, **kwargs)
```

Multimodal training set

**Parameters**

- **matrix** (`scipy.sparse.csr_matrix`) – Preferences in the form of scipy sparse matrix.
- **max\_rating** (`float`) – Maximum value of the preferences.
- **min\_rating** (`float`) – Minimum value of the preferences.
- **global\_mean** (`float`) – Average value of the preferences.
- **uid\_map** (`defaultdict`) – The dictionary containing mapping from original ids to mapped ids of users.
- **iid\_map** (`defaultdict`) – The dictionary containing mapping from original ids to mapped ids of items.

```
class cornac.data.TestSet(user_ratings, uid_map, iid_map)
```

Test Set

**Parameters**

- **user\_ratings** (`defaultdict of list`) – The dictionary containing lists of tuples of the form (item, rating). The keys are user ids.
- **uid\_map** (`defaultdict`) – The dictionary containing mapping from original ids to mapped ids of users.
- **iid\_map** (`defaultdict`) – The dictionary containing mapping from original ids to mapped ids of items.

```
classmethod from_uir(data, global_uid_map, global_iid_map, global_ui_set, verbose=False)
```

Constructing TestSet from triplet data.

**Parameters**

- **data** (`array-like`, shape: [n\_examples, 3]) – Data in the form of triplets (user, item, rating)
- **global\_uid\_map** (`defaultdict`) – The dictionary containing global mapping from original ids to mapped ids of users.
- **global\_iid\_map** (`defaultdict`) – The dictionary containing global mapping from original ids to mapped ids of items.
- **global\_ui\_set** (`set`) – The global set of tuples (user, item). This helps avoiding duplicate observations.
- **verbose** (`bool`, default: `False`) – The verbosity flag.

**Returns** `test_set` – TestSet object.

**Return type** <cornac.data.TestSet>

**get\_iid**(raw\_iid)  
Return the mapped id of an item given a raw id

**get\_ratings**(mapped\_uid)  
Return a list of tuples of (item, rating) of given mapped user id

**get\_uid**(raw\_uid)  
Return the mapped id of a user given a raw id

**users**  
Return a list of users

**class** cornac.data.MultimodalTestSet(user\_ratings, uid\_map, iid\_map, \*\*kwargs)  
Test Set

#### Parameters

- **user\_ratings** (defaultdict of list) – The dictionary containing lists of tuples of the form (item, rating). The keys are user ids.
- **uid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of users.
- **iid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of items.

## 3.1 Train Set

@author: Quoc-Tuan Truong <tuantq.vnu@gmail.com>

**class** cornac.data.trainset.MatrixTrainSet(matrix, max\_rating, min\_rating, global\_mean, uid\_map, iid\_map)

Training set contains preference matrix

#### Parameters

- **matrix** (scipy.sparse.csr\_matrix) – Preferences in the form of scipy sparse matrix.
- **max\_rating** (*float*) – Maximum value of the preferences.
- **min\_rating** (*float*) – Minimum value of the preferences.
- **global\_mean** (*float*) – Average value of the preferences.
- **uid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of users.
- **iid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of items.

**classmethod** **from\_uir**(data, global\_uid\_map=None, global\_iid\_map=None, global\_ui\_set=None, verbose=False)

Constructing TrainSet from triplet data.

#### Parameters

- **data** (array-like, shape: [n\_examples, 3]) – Data in the form of triplets (user, item, rating)

- **global\_uid\_map** (defaultdict, optional, default: None) – The dictionary containing global mapping from original ids to mapped ids of users.
- **global\_iid\_map** (defaultdict, optional, default: None) – The dictionary containing global mapping from original ids to mapped ids of items.
- **global\_ui\_set** (set, optional, default: None) – The global set of tuples (user, item). This helps avoiding duplicate observations.
- **verbose** (bool, default: False) – The verbosity flag.

**Returns** `train_set` – MatrixTrainSet object.

**Return type** <cornac.data.MatrixTrainSet>

**item\_iter** (`batch_size=1, shuffle=False`)

Create an iterator over item ids

**Parameters**

- **batch\_size** (int, optional, default = 1) –
- **shuffle** (bool, optional) – If True, orders of triplets will be randomized. If False, default orders kept

**Returns** iterator

**Return type** batch of item ids (array of np.int)

**uij\_iter** (`batch_size=1, shuffle=False`)

Create an iterator over data yielding batch of users, positive items, and negative items

**Parameters**

- **batch\_size** (int, optional, default = 1) –
- **shuffle** (bool, optional) – If True, orders of triplets will be randomized. If False, default orders kept

**Returns** iterator – batch of negative items (array of np.int)

**Return type** batch of users (array of np.int), batch of positive items (array of np.int),

**uir\_iter** (`batch_size=1, shuffle=False`)

Create an iterator over data yielding batch of users, items, and rating values

**Parameters**

- **batch\_size** (int, optional, default = 1) –
- **shuffle** (bool, optional) – If True, orders of triplets will be randomized. If False, default orders kept

**Returns** iterator – batch of ratings (array of np.float)

**Return type** batch of users (array of np.int), batch of items (array of np.int),

**user\_iter** (`batch_size=1, shuffle=False`)

Create an iterator over user ids

**Parameters**

- **batch\_size** (int, optional, default = 1) –
- **shuffle** (bool, optional) – If True, orders of triplets will be randomized. If False, default orders kept

**Returns** iterator

**Return type** batch of user ids (array of np.int)

```
class cornac.data.trainset.MultimodalTrainSet(matrix, max_rating, min_rating,
                                              global_mean, uid_map, iid_map,
                                              **kwargs)
```

Multimodal training set

#### Parameters

- **matrix** (scipy.sparse.csr\_matrix) – Preferences in the form of scipy sparse matrix.
- **max\_rating** (*float*) – Maximum value of the preferences.
- **min\_rating** (*float*) – Minimum value of the preferences.
- **global\_mean** (*float*) – Average value of the preferences.
- **uid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of users.
- **iid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of items.

```
class cornac.data.trainset.TrainSet(uid_map, iid_map)
```

Training Set

#### Parameters

- **uid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of users.
- **iid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of items.

**get\_iid**(raw\_iid)

Return the mapped id of an item given a raw id

**get\_uid**(raw\_uid)

Return the mapped id of a user given a raw id

**static idx\_iter**(idx\_range, batch\_size=1, shuffle=False)

Create an iterator over batch of indices

#### Parameters

- **batch\_size** (*int*, optional, default = 1) –
- **shuffle** (*bool*, optional) – If True, orders of triplets will be randomized. If False, default orders kept

#### Returns iterator

**Return type** batch of indices (array of np.int)

**iid\_list**

Return the list of mapped item ids

**is\_unk\_item**(mapped\_iid)

Return whether or not an item is unknown given the mapped id

**is\_unk\_user**(mapped\_uid)

Return whether or not a user is unknown given the mapped id

**num\_items**

Return the number of items

---

**num\_users**  
Return the number of users

**raw\_iid\_list**  
Return the list of raw item ids

**raw\_uid\_list**  
Return the list of raw user ids

**uid\_list**  
Return the list of mapped user ids

## 3.2 Test Set

@author: Quoc-Tuan Truong <tuantq.vnu@gmail.com>

**class** cornac.data.testset.**MultimodalTestSet** (*user\_ratings*, *uid\_map*, *iid\_map*, \*\**kwargs*)  
Test Set

### Parameters

- **user\_ratings** (defaultdict of `list`) – The dictionary containing lists of tuples of the form (item, rating). The keys are user ids.
- **uid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of users.
- **iid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of items.

**class** cornac.data.testset.**TestSet** (*user\_ratings*, *uid\_map*, *iid\_map*)  
Test Set

### Parameters

- **user\_ratings** (defaultdict of `list`) – The dictionary containing lists of tuples of the form (item, rating). The keys are user ids.
- **uid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of users.
- **iid\_map** (defaultdict) – The dictionary containing mapping from original ids to mapped ids of items.

**classmethod** **from\_uir** (*data*, *global\_uid\_map*, *global\_iid\_map*, *global\_ui\_set*, *verbose=False*)  
Constructing TestSet from triplet data.

### Parameters

- **data** (*array-like*, *shape*: `[n_examples, 3]`) – Data in the form of triplets (user, item, rating)
- **global\_uid\_map** (defaultdict) – The dictionary containing global mapping from original ids to mapped ids of users.
- **global\_iid\_map** (defaultdict) – The dictionary containing global mapping from original ids to mapped ids of items.
- **global\_ui\_set** (`set`) – The global set of tuples (user, item). This helps avoiding duplicate observations.
- **verbose** (`bool`, *default*: `False`) – The verbosity flag.

**Returns** `test_set` – TestSet object.

**Return type** `<cornac.data.TestSet>`

**get\_iid** (`raw_iid`)  
Return the mapped id of an item given a raw id

**get\_ratings** (`mapped_uid`)  
Return a list of tuples of (item, rating) of given mapped user id

**get\_uid** (`raw_uid`)  
Return the mapped id of a user given a raw id

**users**  
Return a list of users

### 3.3 Graph Module

@author: Aghiles Salah <[asalah@smu.edu.sg](mailto:asalah@smu.edu.sg)>

**class** `cornac.data.graph.GraphModule(**kwargs)`  
Graph module

**batch** (`batch_ids`)  
Return batch of vectors from the sparse adjacency matrix corresponding to provided batch\_ids.

**Parameters** `batch_ids` (`array, required`) – An array contains the ids of rows to be returned from the sparse adjacency matrix.

**build** (`id_map=None`)  
Build the feature matrix. Features will be swapped if the id\_map is provided

**get\_train\_triplet** (`train_row_ids, train_col_ids`)  
Get the training tuples

### 3.4 Text Module

@author: Quoc-Tuan Truong <[tuantq.vnu@gmail.com](mailto:tuantq.vnu@gmail.com)>

**class** `cornac.data.text.Tokenizer`  
Generic class for other subclasses to extend from. This typically either splits text into word tokens or character tokens.

**batch\_tokenize** (`texts: List[str]`) → `List[List[str]]`  
Splitting a corpus with multiple text documents.

**Returns** `tokens`

**Return type** `List[List[str]]`

**tokenize** (`t: str`) → `List[str]`  
Splitting text into tokens.

**Returns** `tokens`

**Return type** `List[str]`

**class** `cornac.data.text.BaseTokenizer` (`sep: str = ''`, `pre_rules: List[Callable[str, str]] = None`, `stop_words: Union[List, str] = None`)  
A base tokenizer use a provided delimiter `sep` to split text.

---

**batch\_tokenize** (*texts*: *List[str]*) → *List[List[str]]*  
Splitting a corpus with multiple text documents.

**Returns tokens**

**Return type** *List[List[str]]*

**tokenize** (*t*: *str*) → *List[str]*  
Splitting text into tokens.

**Returns tokens**

**Return type** *List[str]*

**class** cornac.data.text.Vocabulary (*idx2tok*: *List[str]*, *use\_special\_tokens*: *bool* = *False*)

Vocabulary basically contains mapping between numbers and tokens and vice versa.

**classmethod from\_sequences** (*sequences*: *List[List[str]]*, *max\_vocab*: *int* = *None*,  
*min\_freq*: *int* = *1*, *use\_special\_tokens*: *bool* = *False*) →  
cornac.data.text.Vocabulary

Build a vocabulary from sequences (list of list of tokens).

**classmethod from\_tokens** (*tokens*: *List[str]*, *max\_vocab*: *int* = *None*, *min\_freq*: *int* = *1*,  
*use\_special\_tokens*: *bool* = *False*) → cornac.data.text.Vocabulary

Build a vocabulary from list of tokens.

**classmethod load** (*path*)

Load a vocabulary from *path* to a pickle file.

**save** (*path*)

Save *idx2tok* into a pickle file.

**to\_idx** (*tokens*: *List[str]*) → *List[int]*

Convert a list of *tokens* to their integer indices.

**to\_text** (*indices*: *List[int]*, *sep*=*' '*) → *List[str]*

Convert a list of integer *indices* to their tokens.

**class** cornac.data.text.CountVectorizer (*tokenizer*: cornac.data.text.Tokenizer = *None*,  
*vocab*: cornac.data.text.Vocabulary = *None*,  
*max\_doc\_freq*: Union[float, int] = *1.0*, *min\_freq*:  
*int* = *1*, *max\_features*: *int* = *None*, *stop\_words*:  
Union[List, str] = *None*, *binary*: *bool* = *False*)

Convert a collection of text documents to a matrix of token counts This implementation produces a sparse representation of the counts using `scipy.sparse.csr_matrix`.

### Parameters

- **tokenizer** (`Tokenizer`, optional, default = *None*) – Tokenizer for text splitting. If *None*, the `BaseTokenizer` will be used.
- **vocab** (`Vocabulary`, optional, default = *None*) – Vocabulary of tokens. It contains mapping between tokens to their integer ids and vice versa.
- **max\_doc\_freq** (`Union[float, int]` = *1.0*) – The maximum frequency of tokens appearing in documents to be excluded from vocabulary. If float, the value represents a proportion of documents, int for absolute counts. If *vocab* is not *None*, this will be ignored.
- **min\_freq** (`int`, default = *1*) – The minimum frequency of tokens to be included into vocabulary. If *vocab* is not *None*, this will be ignored.
- **max\_features** (`int`, default=*None*) – If not *None*, build a vocabulary that only consider the top *max\_features* ordered by term frequency across the corpus. If *vocab* is not *None*, this will be ignored.

- **stop\_words** (*Collection, str, default: None*) – Collection of stop words which will be ignored when building *Vocabulary*. If str, it indicates a built-in stop words list. Currently, only *english* is supported.
- **binary** (*boolean, default=False*) – If True, all non zero counts are set to 1.

**fit** (*raw\_documents: List[str]*) → cornac.data.text.CountVectorizer

Build a vocabulary of all tokens in the raw documents.

**Parameters** **raw\_documents** (*iterable*) – An iterable which yields either str, unicode or file objects.

**Returns**

**Return type** self

**fit\_transform** (*raw\_documents: List[str]*) → (*typing.List[typing.List[str]]*, *<class 'scipy.sparse.csr.csr\_matrix'>*)

Build the vocabulary and return term-document matrix.

**Parameters** **raw\_documents** (*List[str]*) –

**Returns**

**sequences:** *List[List[str]]* Tokenized sequences of raw\_documents

**X:** array, [n\_samples, n\_features] Document-term matrix.

**Return type** (sequences, X)

**transform** (*raw\_documents: List[str]*) → (*typing.List[typing.List[str]]*, *<class 'scipy.sparse.csr.csr\_matrix'>*)

Transform documents to document-term matrix.

**Parameters** **raw\_documents** (*List[str]*) –

**Returns**

**sequences:** *List[List[str]]* Tokenized sequences of raw\_documents.

**X:** array, [n\_samples, n\_features] Document-term matrix.

**Return type** (sequences, X)

**class** cornac.data.text.TextModule (*corpus: List[str] = None, ids: List = None, tokenizer: cornac.data.text.Tokenizer = None, vocab: cornac.data.text.Vocabulary = None, max\_vocab: int = None, max\_doc\_freq: Union[float, int] = 1.0, min\_freq: int = 1, stop\_words: Union[List, str] = None, \*\*kwargs*)

Text module

**Parameters**

- **corpus** (*List[str], default = None*) – List of user/item texts that the indices are aligned with *ids*.
- **ids** (*List, default = None*) – List of user/item ids that the indices are aligned with *corpus*. If None, the indices of provided *corpus* will be used as *ids*.
- **tokenizer** (*Tokenizer, optional, default = None*) – Tokenizer for text splitting. If None, the BaseTokenizer will be used.
- **vocab** (*Vocabulary, optional, default = None*) – Vocabulary of tokens. It contains mapping between tokens to their integer ids and vice versa.

- **max\_vocab** (*int*, *optional*, *default* = *None*) – The maximum size of the vocabulary. If vocab is provided, this will be ignored.
- **max\_doc\_freq** (*Union[float, int]* = *1.0*) – The maximum frequency of tokens appearing in documents to be excluded from vocabulary. If float, the value represents a proportion of documents, int for absolute counts. If *vocab* is not None, this will be ignored.
- **min\_freq** (*int*, *default* = *1*) – The minimum frequency of tokens to be included into vocabulary. If *vocab* is not None, this will be ignored.
- **stop\_words** (*Collection, str, default: None*) – Collection of stop words which will be ignored when building *Vocabulary*. If str, it indicates a built-in stop words list. Currently, only *english* is supported.

**batch\_seq** (*batch\_ids, max\_length=None*)

Return a numpy matrix of text sequences containing token ids with size=(len(*batch\_ids*), *max\_length*). If *max\_length=None*, it will be inferred based on retrieved sequences.

**batch\_tfidf** (*batch\_ids*)

Return matrix of TF-IDF features corresponding to provided *batch\_ids*

**build** (*id\_map=None*)

Build the model based on provided list of ordered ids

## 3.5 Image Module

@author: Quoc-Tuan Truong <[tuantq.vnu@gmail.com](mailto:tuantq.vnu@gmail.com)>

**class** cornac.data.image.**ImageModule** (\*\**kwargs*)

Image module

**batch\_image** (*batch\_ids, target\_size=(256, 256), color\_mode='rgb', interpolation='nearest'*)

Return batch of images corresponding to provided *batch\_ids*

**build** (*id\_map=None*)

Build the model based on provided list of ordered ids

## 3.6 Reader

@author: Quoc-Tuan Truong <[tuantq.vnu@gmail.com](mailto:tuantq.vnu@gmail.com)>

**cornac.data.reader.read\_ui** (*fpather, value=1.0, sep='\t', skip\_lines=0*)

Read data in the form of implicit feedback user-items. Each line starts with user id followed by multiple of item ids.

### Parameters

- **fpather** (*str*) – Path to the data file
- **value** (*float*, *default: 1.0*) – Value for the feedback
- **sep** (*str, default: \t*) – The delimiter string.
- **skip\_lines** (*int, default: 0*) – Number of first lines to skip

**Returns** **triplets** – Data in the form of list of tuples of (user, item, 1).

**Return type** iterable

`cornac.data.reader.read_uir(fpath, u_col=0, i_col=1, r_col=2, sep='\t', skip_lines=0)`  
Read data in the form of triplets (user, item, rating).

**Parameters**

- **fpath** (`str`) – Path to the data file
- **u\_col** (`int`, `default`: 0) – Index of the user column
- **i\_col** (`int`, `default`: 1) – Index of the item column
- **r\_col** (`int`, `default`: 2) – Index of the rating column
- **sep** (`str`, `default`: :) – The delimiter string.
- **skip\_lines** (`int`, `default`: 0) – Number of first lines to skip

**Returns** `triplets` – Data in the form of list of tuples of (user, item, rating).

**Return type** `iterable`

# CHAPTER 4

---

## Models

---

### 4.1 Probabilistic Collaborative Representation Learning (PCRL)

@author: Aghiles Salah <asalah@smu.edu.sg>

```
class cornac.models.pcrl.recom_pcrl.PCRL(k=100,      z_dims=[300],      max_iter=300,
                                         batch_size=300,      learning_rate=0.001,
                                         name='pcrl',   trainable=True,   verbose=False,
                                         w_determinist=True, init_params={'G_r': None,
                                         'G_s': None, 'L_r': None, 'L_s': None})
```

Probabilistic Collaborative Representation Learning.

#### Parameters

- **k** (*int*, *optional*, *default*: 100) – The dimension of the latent factors.
- **z\_dims** (*Numpy 1d array*, *optional*, *default*: [300]) – The dimensions of the hidden intermediate layers ‘z’ in the order [dim(z\_L), …, dim(z\_1)], please refer to Figure 1 in the original paper for more details.
- **max\_iter** (*int*, *optional*, *default*: 300) – Maximum number of iterations (number of epochs) for variational PCRL.
- **batch\_size** (*int*, *optional*, *default*: 300) – The batch size for SGD.
- **learning\_rate** (*float*, *optional*, *default*: 0.001) – The learning rate for SGD.
- **aux\_info** (*see "cornac/examples/pcrl\_example.py" in the GitHub repo for an example of how to use cornac's graph module provide item auxiliary data (e.g., context, text, etc.) for PCRL.*) –
- **name** (*string*, *optional*, *default*: 'PCRL') – The name of the recommender model.

- **trainable** (*boolean, optional, default: True*) – When False, the model is not trained and Cornac assumes that the model already pre-trained (Theta, Beta and Xi are not None).
- **w\_determinist** (*boolean, optional, default: True*) – When True, determinist wheights “W” are used for the generator network, otherwise “W” is stochastic as in the original paper.
- **init\_params** (*dictionary, optional, default: { 'G\_s':None, 'G\_r':None, 'L\_s':None, 'L\_r':None }*) – List of initial parameters, e.g., init\_params = {'G\_s':G\_s, 'G\_r':G\_r, 'L\_s':L\_s, 'L\_r':L\_r}, where G\_s and G\_r are of type csc\_matrix or np.array with the same shape as Theta, see below). They represent respectively the “shape” and “rate” parameters of Gamma distribution over Theta. It is the same for L\_s, L\_r and Beta.
- **Theta** (*csc\_matrix, shape (n\_users, k)*) – The expected user latent factors.
- **Beta** (*csc\_matrix, shape (n\_items, k)*) – The expected item latent factors.

## References

- Salah, Aghiles, and Hady W. Lauw. Probabilistic Collaborative Representation Learning for Personalized Item Recommendation. In UAI 2018.

### **fit** (*train\_set*)

Fit the model to observations.

**Parameters** **train\_set** (*object of type TrainSet, required*) – An object containing the user-item preference in csr scipy sparse format, as well as some useful attributes such as mappings to the original user/item ids. Please refer to the class TrainSet in the “data” module for details.

### **score** (*user\_id, item\_id=None*)

Predict the scores/ratings of a user for a list of items.

#### Parameters

- **user\_id** (*int, required*) – The index of the user for whom to perform score prediction.
- **item\_id** (*int, optional, default: None*) – The index of the item for that to perform score prediction. If None, scores for all known items will be returned.

**Returns** **res** – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.2 Collaborative Context Poisson Factorization (C2PF)

@author: Aghiles Salah <asalah@smu.edu.sg>

```
class cornac.models.c2pf.recom_c2pf.C2PF(k=100,      max_iter=100,      variant='c2pf',
                                           name=None,    trainable=True,   verbose=False,
                                           init_params={'G_r': None, 'G_s': None, 'L2_r':
                                           None, 'L2_s': None, 'L3_r': None, 'L3_s': None,
                                           'L_r': None, 'L_s': None})
```

Collaborative Context Poisson Factorization.

## Parameters

- **k** (*int*, optional, default: 100) – The dimension of the latent factors.
- **max\_iter** (*int*, optional, default: 100) – Maximum number of iterations for variational C2PF.
- **variant** (*string*, optional, default: 'c2pf') – C2pf's variant: c2pf: 'c2pf', 'tc2pf' (tied-c2pf) or 'rc2pf' (reduced-c2pf). Please refer to the original paper for details.
- **name** (*string*, optional, default: None) – The name of the recommender model. If None, then "variant" is used as the default name of the model.
- **trainable** (*boolean*, optional, default: True) – When False, the model is not trained and Cornac assumes that the model already pre-trained (Theta, Beta and Xi are not None).
- **Item\_context** (See "cornac/examples/c2pf\_example.py" in the GitHub repo for an example of how to use cornac's graph module to load and provide "item context" for C2PF.) –
- **init\_params** (dictionary, optional, default: {'G\_s':None, 'G\_r':None, 'L\_s':None, 'L\_r':None, 'L2\_s':None, 'L2\_r':None, 'L3\_s':None, 'L3\_r':None}) – List of initial parameters, e.g., init\_params = {'G\_s':G\_s, 'G\_r':G\_r, 'L\_s':L\_s, 'L\_r':L\_r, 'L2\_s':L2\_s, 'L2\_r':L2\_r, 'L3\_s':L3\_s, 'L3\_r':L3\_r}, where G\_s and G\_r are of type csc\_matrix or np.array with the same shape as Theta, see below). They represent respectively the "shape" and "rate" parameters of Gamma distribution over Theta. It is the same for L\_s, L\_r and Beta, L2\_s, L2\_r and Xi, L3\_s, L3\_r and Kappa.
- **Theta** (*csc\_matrix*, shape (n\_users, k)) – The expected user latent factors.
- **Beta** (*csc\_matrix*, shape (n\_items, k)) – The expected item latent factors.
- **Xi** (*csc\_matrix*, shape (n\_items, k)) – The expected context item latent factors multiplied by context effects Kappa, please refer to the paper below for details.

## References

- Salah, Aghiles, and Hady W. Lauw. A Bayesian Latent Variable Model of User Preferences with Item Context. In IJCAI, pp. 2667-2674. 2018.

**fit** (*train\_set*)

Fit the model to observations.

**Parameters** **train\_set** (*object of type TrainSet, required*) – An object containing the user-item preference in csr scipy sparse format, as well as some useful attributes such as mappings to the original user/item ids. Please refer to the class TrainSet in the "data" module for details.

**score** (*user\_id*, *item\_id=None*)

Predict the scores/ratings of a user for an item.

## Parameters

- **user\_id** (*int*, required) – The index of the user for whom to perform score prediction.
- **item\_id** (*int*, optional, default: None) – The index of the item for that to perform score prediction. If None, scores for all known items will be returned.

**Returns** `res` – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.3 Indexable Bayesian Personalized Ranking (IBPR)

@author: Dung D. Le (Andrew) <ddle.2015@smu.edu.sg>

```
class cornac.models.ibpr.recom_ibpr.IBPR(k=20, max_iter=100, learning_rate=0.05,
                                         lamda=0.001, batch_size=100, name='ibpr',
                                         trainable=True, verbose=False,
                                         init_params=None)
```

Indexable Bayesian Personalized Ranking.

### Parameters

- `k` (`int`, optional, default: 20) – The dimension of the latent factors.
- `max_iter` (`int`, optional, default: 100) – Maximum number of iterations or the number of epochs for SGD.
- `learning_rate` (`float`, optional, default: 0.05) – The learning rate for SGD.
- `lamda` (`float`, optional, default: 0.001) – The regularization parameter.
- `batch_size` (`int`, optional, default: 100) – The batch size for SGD.
- `name` (`string`, optional, default: 'IBPR') – The name of the recommender model.
- `trainable` (`boolean`, optional, default: True) – When False, the model is not trained and Cornac assumes that the model already pre-trained (U and V are not None).
- `verbose` (`boolean`, optional, default: False) – When True, some running logs are displayed.
- `init_params` (`dictionary`, optional, default: None) – List of initial parameters, e.g., `init_params = {'U':U, 'V':V}` please see below the definition of U and V.
- `U` (`csc_matrix`, shape (n\_users, k)) – The user latent factors, optional initialization via `init_params`.
- `V` (`csc_matrix`, shape (n\_items, k)) – The item latent factors, optional initialization via `init_params`.

### References

- Le, D. D., & Lauw, H. W. (2017, November). Indexable Bayesian personalized ranking for efficient top-k recommendation. In Proceedings of the 2017 ACM on Conference on Information and Knowledge Management (pp. 1389-1398). ACM.

`fit` (`train_set`)

Fit the model to observations.

**Parameters** `train_set` (*object of type TrainSet, required*) – An object containing the user-item preference in csr scipy sparse format, as well as some useful attributes such as mappings to the original user/item ids. Please refer to the class TrainSet in the “data” module for details.

**score** (*user\_id, item\_id=None*)

Predict the scores/ratings of a user for an item.

#### Parameters

- `user_id` (*int, required*) – The index of the user for whom to perform score prediction.
- `item_id` (*int, optional, default: None*) – The index of the item for that to perform score prediction. If None, scores for all known items will be returned.

**Returns** `res` – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.4 Online Indexable Bayesian Personalized Ranking (OIBPR)

@author: Dung D. Le (Andrew) <ddle.2015@smu.edu.sg>

```
class cornac.models.online_ibpr.recom_online_ibpr.OnlineIBPR(k=20,
                                                               max_iter=100,
                                                               learn-
                                                               ing_rate=0.05,
                                                               lamda=0.001,
                                                               batch_size=100,
                                                               name='online_ibpr',
                                                               trainable=True,
                                                               verbose=False,
                                                               init_params=None)
```

Online Indexable Bayesian Personalized Ranking.

#### Parameters

- `k` (*int, optional, default: 20*) – The dimension of the latent factors.
- `max_iter` (*int, optional, default: 100*) – Maximum number of iterations or the number of epochs for SGD.
- `learning_rate` (*float, optional, default: 0.05*) – The learning rate for SGD.
- `lamda` (*float, optional, default: 0.001*) – The regularization parameter.
- `batch_size` (*int, optional, default: 100*) – The batch size for SGD.
- `name` (*string, optional, default: 'IBPR'*) – The name of the recommender model.
- `trainable` (*boolean, optional, default: True*) – When False, the model is not trained and Cornac assumes that the model already pre-trained (U and V are not None).
- `verbose` (*boolean, optional, default: False*) – When True, some running logs are displayed.

- **init\_params** (*dictionary, optional, default: None*) – List of initial parameters, e.g., `init_params = {'U':U, 'V':V}` please see below the definition of U and V.
- **U** (*csc\_matrix, shape (n\_users, k)*) – The user latent factors, optional initialization via `init_params`.
- **V** (*csc\_matrix, shape (n\_items, k)*) – The item latent factors, optional initialization via `init_params`.

## References

- Le, D. D., & Lauw, H. W. (2017, November). Indexable Bayesian personalized ranking for efficient top-k recommendation. In Proceedings of the 2017 ACM on Conference on Information and Knowledge Management (pp. 1389-1398). ACM.

**fit** (*train\_set*)

Fit the model to observations.

**Parameters** **train\_set** (*object of type TrainSet, required*) – An object containing the user-item preference in csr scipy sparse format, as well as some useful attributes such as mappings to the original user/item ids. Please refer to the class TrainSet in the “data” module for details.

**score** (*user\_id, item\_id=None*)

Predict the scores/ratings of a user for an item.

**Parameters**

- **user\_id** (*int, required*) – The index of the user for whom to perform score prediction.
- **item\_id** (*int, optional, default: None*) – The index of the item for that to perform score prediction. If None, scores for all known items will be returned.

**Returns** **res** – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.5 Collaborative Ordinal Embedding (COE)

@author: Dung D. Le (Andrew) <ddle.2015@smu.edu.sg>

```
class cornac.models.coe.recom_coe.COE(k=20,      max_iter=100,      learning_rate=0.05,
                                         lamda=0.001,     batch_size=1000,     name='coe',
                                         trainable=True, verbose=False, init_params=None)
```

Collaborative Ordinal Embedding.

**Parameters**

- **k** (*int, optional, default: 20*) – The dimension of the latent factors.
- **max\_iter** (*int, optional, default: 100*) – Maximum number of iterations or the number of epochs for SGD.
- **learning\_rate** (*float, optional, default: 0.05*) – The learning rate for SGD.
- **lamda** (*float, optional, default: 0.001*) – The regularization parameter.

- **batch\_size** (*int*, *optional*, *default*: 100) – The batch size for SGD.
- **name** (*string*, *optional*, *default*: 'IBRP') – The name of the recommender model.
- **trainable** (*boolean*, *optional*, *default*: True) – When False, the model is not trained and Cornac assumes that the model already pre-trained (U and V are not None).
- **verbose** (*boolean*, *optional*, *default*: False) – When True, some running logs are displayed.
- **init\_params** (*dictionary*, *optional*, *default*: None) – List of initial parameters, e.g., init\_params = {'U':U, 'V':V} please see below the definition of U and V.
- **U** (*csc\_matrix*, *shape* (*n\_users*, *k*)) – The user latent factors, optional initialization via init\_params.
- **V** (*csc\_matrix*, *shape* (*n\_items*, *k*)) – The item latent factors, optional initialization via init\_params.

## References

- Le, D. D., & Lauw, H. W. (2016, June). Euclidean co-embedding of ordinal data for multi-type visualization. In Proceedings of the 2016 SIAM International Conference on Data Mining (pp. 396-404). Society for Industrial and Applied Mathematics.

**fit** (*train\_set*)

Fit the model to observations.

**Parameters** **train\_set** (*object of type TrainSet, required*) – An object containing the user-item preference in csr scipy sparse format, as well as some useful attributes such as mappings to the original user/item ids. Please refer to the class TrainSet in the “data” module for details.

**score** (*user\_id*, *item\_id=None*)

Predict the scores/ratings of a user for an item.

**Parameters**

- **user\_id** (*int*, *required*) – The index of the user for whom to perform score prediction.
- **item\_id** (*int*, *optional*, *default*: None) – The index of the item for that to perform score prediction. If None, scores for all known items will be returned.

**Returns** **res** – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.6 Visual Bayesian Personalized Ranking (VBPR)

@author: Guo Jingyao <jyguo@smu.edu.sg> Quoc-Tuan Truong <tuantq.vnu@gmail.com>

```
class cornac.models.vbpr.recom_vbpr.VBPR(k=10, k2=10, n_epochs=20, batch_size=100,
                                         learning_rate=0.001, lambda_w=0.01,
                                         lambda_b=0.01, lambda_e=0.0, use_gpu=False,
                                         trainable=True, init_params=None, **kwargs)
```

Visual Bayesian Personalized Ranking.

#### Parameters

- **k** (*int*, *optional*, *default*: 10) – The dimension of the gamma latent factors.
- **k2** (*int*, *optional*, *default*: 10) – The dimension of the theta latent factors.
- **n\_epochs** (*int*, *optional*, *default*: 20) – Maximum number of epochs for SGD.
- **batch\_size** (*int*, *optional*, *default*: 100) – The batch size for SGD.
- **learning\_rate** (*float*, *optional*, *default*: 0.001) – The learning rate for SGD.
- **lambda\_w** (*float*, *optional*, *default*: 0.01) – The regularization hyper-parameter for latent factor weights.
- **lambda\_b** (*float*, *optional*, *default*: 0.01) – The regularization hyper-parameter for biases.
- **lambda\_e** (*float*, *optional*, *default*: 0.0) – The regularization hyper-parameter for embedding matrix E and beta prime vector.
- **use\_gpu** (*boolean*, *optional*, *default*: True) – Whether or not to use GPU to speed up training.
- **trainable** (*boolean*, *optional*, *default*: True) – When False, the model is not trained and Cornac assumes that the model already pre-trained (U and V are not None).
- **init\_params** (*dictionary*, *optional*, *default*: None) –  
**Initial parameters, e.g.,** `init_params = {'Bi': beta_item, 'Gu': gamma_user, 'Gi': gamma_item, 'Tu': theta_user, 'E': emb_matrix, 'Bp': beta_prime}`

#### References

- HE, Ruining et MCAULEY, Julian. VBPR: Visual Bayesian Personalized Ranking from Implicit Feedback. In : AAAI. 2016. p. 144-150.

**fit** (*train\_set*)

Fit the model.

**Parameters** **train\_set** (*cornac.data.MultimodalTrainSet*) – Multimodal training set.

**score** (*user\_id*, *item\_id*=None)

Predict the scores/ratings of a user for an item.

#### Parameters

- **user\_id** (*int*, *required*) – The index of the user for whom to perform score prediction.
- **item\_id** (*int*, *optional*, *default*: None) – The index of the item for that to perform score prediction. If None, scores for all known items will be returned.

**Returns** `res` – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.7 Spherical k-means (Skmeans)

@author: Aghiles Salah <asalah@smu.edu.sg>

```
class cornac.models.skm.recom_skmeans.SKMeans (k=5, max_iter=100, name='Skmeans', trainable=True, tol=1e-06, verbose=True, init_par=None)
```

Spherical k-means based recommender.

### Parameters

- **k** (*int, optional, default: 5*) – The number of clusters.
- **max\_iter** (*int, optional, default: 100*) – Maximum number of iterations.
- **name** (*string, optional, default: 'Skmeans'*) – The name of the recommender model.
- **trainable** (*boolean, optional, default: True*) – When False, the model is not trained and Cornac assumes that the model is already trained.
- **tol** (*float, optional, default: 1e-6*) – Relative tolerance with regards to skmeans' criterion to declare convergence.
- **verbose** (*boolean, optional, default: False*) – When True, some running logs are displayed.
- **init\_par** (*numpy 1d array, optional, default: None*) – The initial object partition, 1d array containing the cluster label (int type starting from 0) of each object (user). If par = None, then skmeans is initialized randomly.
- **centroids** (*csc\_matrix, shape (k, n\_users)*) – The matrix of cluster centroids.

### References

- Salah, Aghiles, Nicoleta Rogovschi, and Mohamed Nadif. “A dynamic collaborative filtering system via a weighted clustering approach.” Neurocomputing 175 (2016): 206-215.

**fit** (*train\_set*)

Fit the model to observations.

**Parameters** `train_set` (*object of type TrainSet, required*) – An object containing the user-item preference in csr scipy sparse format, as well as some useful attributes such as mappings to the original user/item ids. Please refer to the class TrainSet in the “data” module for details.

**score** (*user\_id, item\_id=None*)

Predict the scores/ratings of a user for an item.

### Parameters

- **user\_id** (*int, required*) – The index of the user for whom to perform score prediction.

- **item\_id** (*int*, *optional*, *default*: *None*) – The index of the item for that to perform score prediction. If *None*, scores for all known items will be returned.

**Returns** **res** – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.8 Collaborative Deep Learning (CDL)

@author: Trieu Thi Ly

```
class cornac.models.cdl.recom_cdl.CDL(k=50,          text_information=None,           autoen-
                                         coder_structure=None,           lambda_u=0.1,
                                         lambda_v=0.01, lambda_w=0.01, lambda_n=0.01,
                                         a=1, b=0.01, autoencoder_corruption=0.3, learning_
                                         rate=0.001, keep_prob=1.0, batch_size=100,
                                         max_iter=100, name='CDL', trainable=True, verbose=False, init_params=None)
```

Collaborative Deep Learning.

### Parameters

- **k** (*int*, *optional*, *default*: *50*) – The dimension of the latent factors.
- **max\_iter** (*int*, *optional*, *default*: *100*) – Maximum number of iterations or the number of epochs for SGD.
- **shape** (*n\_items*, *n\_vocabularies*), *optional*, *default* (*text\_informationndarray*,) – Bag-of-words features of items
- **optional**, **default** (*autoencoder\_structurearray*,) – The number of neurons of encoder/ decoder layer for SDAE
- **learning\_rate** (*float*, *optional*, *default*: *0.001*) – The learning rate for AdamOptimizer.
- **lambda\_u** (*float*, *optional*, *default*: *0.1*) – The regularization parameter for users.
- **lambda\_v** (*float*, *optional*, *default*: *10*) – The regularization parameter for items.
- **lambda\_w** (*float*, *optional*, *default*: *0.1*) – The regularization parameter for SDAE weights.
- **lambda\_n** (*float*, *optional*, *default*: *1000*) – The regularization parameter for SDAE output.
- **a** (*float*, *optional*, *default*: *1*) – The confidence of observed ratings.
- **b** (*float*, *optional*, *default*: *0.01*) – The confidence of unseen ratings.
- **autoencoder\_corruption** (*float*, *optional*, *default*: *0.3*) – The corruption ratio for SDAE.
- **keep\_prob** (*float*, *optional*, *default*: *1.0*) – The probability that each element is kept in dropout of SDAE.
- **batch\_size** (*int*, *optional*, *default*: *100*) – The batch size for SGD.
- **name** (*string*, *optional*, *default*: *'CDL'*) – The name of the recommender model.

- **trainable** (*boolean, optional, default: True*) – When False, the model is not trained and Cornac assumes that the model already pre-trained (U and V are not None).
- **init\_params** (*dictionary, optional, default: None*) – List of initial parameters, e.g., init\_params = {'U':U, 'V':V} please see below the definition of U and V.
- **U** (*ndarray, shape (n\_users, k)*) – The user latent factors, optional initialization via init\_params.
- **V** (*ndarray, shape (n\_items, k)*) – The item latent factors, optional initialization via init\_params.

## References

- Hao Wang, Naiyan Wang, Dit-Yan Yeung. CDL: Collaborative Deep Learning for Recommender Systems. In : SIGKDD. 2015. p. 1235-1244.

### `fit(train_set)`

Fit the model to observations.

**Parameters** `train_set` (*object of type TrainSet, required*) – An object containing the user-item preference in csr scipy sparse format, as well as some useful attributes such as mappings to the original user/item ids. Please refer to the class TrainSet in the “data” module for details.

### `score(user_id, item_id=None)`

Predict the scores/ratings of a user for an item.

#### Parameters

- **user\_id** (*int, required*) – The index of the user for whom to perform score prediction.
- **item\_id** (*int, optional, default: None*) – The index of the item for that to perform score prediction. If None, scores for all known items will be returned.

**Returns** `res` – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.9 Hierarchical Poisson Factorization (HPF)

@author: Aghiles Salah <asalah@smu.edu.sg>

```
class cornac.models.hpf.recom_hpf.HPF(k=5,      max_iter=100,      name='HPF',      trainable=True,      verbose=False,      hierarchical=True,      init_params={'G_r': None, 'G_s': None, 'L_r': None, 'L_s': None})
```

Hierarchical Poisson Factorization.

#### Parameters

- **k** (*int, optional, default: 5*) – The dimension of the latent factors.
- **max\_iter** (*int, optional, default: 100*) – Maximum number of iterations.

- **name** (*string, optional, default: 'HPF'*) – The name of the recommender model.
- **trainable** (*boolean, optional, default: True*) – When False, the model is not trained and Cornac assumes that the model is already pre-trained (Theta and Beta are not None).
- **verbose** (*boolean, optional, default: False*) – When True, some running logs are displayed.
- **hierarchical** (*boolean, optional, default: True*) – When False, PF is used instead of HPF.
- **init\_params** (*dictionary, optional, default: {'G\_s':None, 'G\_r':None, 'L\_s':None, 'L\_r':None}*) – List of initial parameters, e.g., init\_params = {'G\_s':G\_s, 'G\_r':G\_r, 'L\_s':L\_s, 'L\_r':L\_r}, where G\_s and G\_r are of type csc\_matrix or np.array with the same shape as Theta, see below). They represent respectively the “shape” and “rate” parameters of Gamma distribution over Theta. Similarly, L\_s, L\_r are the shape and rate parameters of the Gamma over Beta.
- **Theta** (*csc\_matrix, shape (n\_users, k)*) – The expected user latent factors.
- **Beta** (*csc\_matrix, shape (n\_items, k)*) – The expected item latent factors.

## References

- Gopalan, Prem, Jake M. Hofman, and David M. Blei. Scalable Recommendation with Hierarchical Poisson Factorization. In UAI, pp. 326-335. 2015.

**fit** (*train\_set*)

Fit the model to observations.

**Parameters** **train\_set** (*object of type TrainSet, required*) – An object containing the user-item preference in csr scipy sparse format, as well as some useful attributes such as mappings to the original user/item ids. Please refer to the class TrainSet in the “data” module for details.

**score** (*user\_id, item\_id=None*)

Predict the scores/ratings of a user for an item.

**Parameters**

- **user\_id** (*int, required*) – The index of the user for whom to perform score prediction.
- **item\_id** (*int, optional, default: None*) – The index of the item for that to perform score prediction. If None, scores for all known items will be returned.

**Returns** **res** – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.10 Bayesian Personalized Ranking (BPR)

```
class cornac.models.bpr.recom_bpr.BPR
    Bayesian Personalized Ranking.
```

**Parameters**

- **k** (*int*, optional, default: 10) – The dimension of the latent factors.
- **max\_iter** (*int*, optional, default: 100) – Maximum number of iterations or the number of epochs for SGD.
- **learning\_rate** (*float*, optional, default: 0.001) – The learning rate for SGD.
- **lambda\_reg** (*float*, optional, default: 0.001) – The regularization hyper-parameter.
- **num\_threads** (*int*, optional, default: 0) – Number of parallel threads for training. If 0, all CPU cores will be utilized.
- **trainable** (*boolean*, optional, default: True) – When False, the model will not be re-trained, and input of pre-trained parameters are required.
- **verbose** (*boolean*, optional, default: True) – When True, some running logs are displayed.
- **init\_params** (*dictionary*, optional, default: None) – Initial parameters, e.g., init\_params = {'U': user\_factors, 'V': item\_factors, 'Bi': item\_biases}

## References

- Rendle, Steffen, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. BPR: Bayesian personalized ranking from implicit feedback. In UAI, pp. 452-461. 2009.

### fit

Fit the model to observations.

**Parameters** **train\_set** (*object of type TrainSet, required*) – An object contains the user-item preference in csr scipy sparse format, as well as some useful attributes such as mappings to the original user/item ids. Please refer to the class TrainSet in the “data” module for details.

### score

Predict the scores/ratings of a user for an item.

#### Parameters

- **user\_id** (*int*, required) – The index of the user for whom to perform score prediction.
- **item\_id** (*int*, optional, default: None) – The index of the item for that to perform score prediction. If None, scores for all known items will be returned.

**Returns** **res** – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.11 Probabilistic Matrix Factorization (PMF)

@author: Aghiles Salah

```
class cornac.models.pmf.recom_pmf.PMF(k=5,      max_iter=100,      learning_rate=0.001,
                                         gamma=0.9,    lamda=0.001,    name='PMF',   vari-
                                         ant='non_linear', trainable=True, verbose=False,
                                         init_params={'U': None, 'V': None})
```

Probabilistic Matrix Factorization.

#### Parameters

- **k** (*int*, *optional*, *default*: 5) – The dimension of the latent factors.
- **max\_iter** (*int*, *optional*, *default*: 100) – Maximum number of iterations or the number of epochs for SGD.
- **learning\_rate** (*float*, *optional*, *default*: 0.001) – The learning rate for SGD\_RMSProp.
- **gamma** (*float*, *optional*, *default*: 0.9) – The weight for previous/current gradient in RMSProp.
- **lamda** (*float*, *optional*, *default*: 0.001) – The regularization parameter.
- **name** (*string*, *optional*, *default*: 'PMF') – The name of the recommender model.
- **variant** ({ "linear", "non\_linear"}, *optional*, *default*: 'non\_linear') – Pmf variant. If 'non\_linear', the Gaussian mean is the output of a Sigmoid function. If 'linear' the Gaussian mean is the output of the identity function.
- **trainable** (*boolean*, *optional*, *default*: True) – When False, the model is not trained and Cornac assumes that the model already pre-trained (U and V are not None).
- **verbose** (*boolean*, *optional*, *default*: False) – When True, some running logs are displayed.
- **init\_params** (*dictionary*, *optional*, *default*: {'U':None, 'V':None}) – List of initial parameters, e.g., init\_params = {'U':U, 'V':V}. U: a csc\_matrix of shape (n\_users,k), containing the user latent factors. V: a csc\_matrix of shape (n\_items,k), containing the item latent factors.

#### References

- Mnih, Andriy, and Ruslan R. Salakhutdinov. Probabilistic matrix factorization. In NIPS, pp. 1257-1264. 2008.

#### **fit** (*train\_set*)

Fit the model to observations.

**Parameters** **train\_set** (*object of type TrainSet, required*) – An object containing the user-item preference in csr scipy sparse format, as well as some useful attributes such as mappings to the original user/item ids. Please refer to the class TrainSet in the “data” module for details.

#### **score** (*user\_id*, *item\_id=None*)

Predict the scores/ratings of a user for an item.

#### Parameters

- **user\_id** (*int*, *required*) – The index of the user for whom to perform score prediction.

- **item\_id** (*int, optional, default: None*) – The index of the item for that to perform score prediction. If None, scores for all known items will be returned.

**Returns** `res` – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.12 Matrix Factorization (MF)

@author: Quoc-Tuan Truong <tuantq.vnu@gmail.com>

**class** `cornac.models.mf.recom_mf.MF`

Matrix Factorization.

### Parameters

- **k** (*int, optional, default: 10*) – The dimension of the latent factors.
- **max\_iter** (*int, optional, default: 100*) – Maximum number of iterations or the number of epochs for SGD.
- **learning\_rate** (*float, optional, default: 0.01*) – The learning rate.
- **lambda\_reg** (*float, optional, default: 0.001*) – The lambda value used for regularization.
- **use\_bias** (*boolean, optional, default: True*) – When True, user, item, and global biases are used.
- **early\_stop** (*boolean, optional, default: False*) – When True, delta loss will be checked after each iteration to stop learning earlier.
- **trainable** (*boolean, optional, default: True*) – When False, the model will not be re-trained, and input of pre-trained parameters are required.
- **verbose** (*boolean, optional, default: True*) – When True, running logs are displayed.
- **init\_params** (*dictionary, optional, default: None*) – Initial parameters, e.g., `init_params = {'U': user_factors, 'V': item_factors, 'Bu': user_biases, 'Bi': item_biases}`

### References

- Koren, Y., Bell, R., & Volinsky, C. Matrix factorization techniques for recommender systems. In Computer, (8), 30-37. 2009.

### fit

Fit the model to observations.

**Parameters** `train_set` (*object of type TrainSet, required*) – An object contains the user-item preference in csr scipy sparse format, as well as some useful attributes such as mappings to the original user/item ids. Please refer to the class TrainSet in the “data” module for details.

### score

Predict the scores/ratings of a user for an item.

### Parameters

- **user\_id** (*int*, *required*) – The index of the user for whom to perform score prediction.
- **item\_id** (*int*, *optional*, *default*: *None*) – The index of the item for that to perform score prediction. If *None*, scores for all known items will be returned.

**Returns** **res** – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array

## 4.13 Convolutional Matrix Factorization (ConvMF)

@author: Tran Thanh Binh

```
class cornac.models.conv_mf.recom_convmf.ConvMF(give_item_weight=True,
                                                n_epochs=50,           lambda_u=1,
                                                lambda_v=100,          k=50,
                                                name='convmf',         trainable=True,
                                                verbose=False,         dropout_rate=0.2,
                                                emb_dim=200,           max_len=300,
                                                num_kernel_per_ws=100,
                                                init_params=None)
```

### Parameters

- **k** (*int*, *optional*, *default*: *50*) – The dimension of the user and item latent factors.
- **n\_epochs** (*int*, *optional*, *default*: *50*) – Maximum number of epochs for training.
- **lambda\_u** (*float*, *optional*, *default*: *1.0*) – The regularization hyper-parameter for user latent factor.
- **lambda\_v** (*float*, *optional*, *default*: *100.0*) – The regularization hyper-parameter for item latent factor.
- **emb\_dim** (*int*, *optional*, *default*: *200*) – The embedding size of each word. One word corresponds with [1 x emb\_dim] vector in the embedding space
- **max\_len** (*int*, *optional*, *default* *300*) – The maximum length of item's document
- **num\_kernel\_per\_ws** (*int*, *optional*, *default*: *100*) – The number of kernel filter in convolutional layer
- **dropout\_rate** (*float*, *optional*, *default*: *0.2*) – Dropout rate while training CNN
- **give\_item\_weight** (*boolean*, *optional*, *default*: *True*) – When True, each item will be weighted base on the number of user who have rated this item
- **init\_params** (*dict*, *optional*, *default*: *{'U':None, 'V':None, 'W': None}*) – Initial U and V matrix and initial weight for embedding layer W
- **trainable** (*boolean*, *optional*, *default*: *True*) – When False, the model is not trained and Cornac assumes that the model already pre-trained (U and V are not *None*).

## References

- Donghyun Kim1, Chanyoung Park1. ConvMF: Convolutional Matrix Factorization for Document Context-Aware Recommendation. In :10th ACM Conference on Recommender Systems Pages 233-240

**fit** (*train\_set*)

Fit the model.

**Parameters** **train\_set** (*cornac.data.MultimodalTrainSet*) – Multimodal training set.

**score** (*user\_id*, *item\_id=None*)

Predict the scores/ratings of a user for an item.

**Parameters**

- **user\_id** (*int*, *required*) – The index of the user for whom to perform score prediction.
- **item\_id** (*int*, *optional, default: None*) – The index of the item for that to perform score prediction. If None, scores for all known items will be returned.

**Returns** **res** – Relative scores that the user gives to the item or to all known items

**Return type** A scalar or a Numpy array



# CHAPTER 5

---

## Metrics

---

### 5.1 Area Under the Curve (AUC)

```
class cornac.metrics.AUC  
Area Under the ROC Curve (AUC).
```

#### References

<https://arxiv.org/ftp/arxiv/papers/1205/1205.2618.pdf>

### 5.2 Normalized Discount Cumulative Gain (NDCG)

```
class cornac.metrics.NDCG(k=-1)  
Normalized Discount Cumulative Gain.  
  
Parameters k (int, optional, default: -1 (all)) – The number of items in the  
top@k list. If None, all items will be considered.
```

#### References

[https://en.wikipedia.org/wiki/Discounted\\_cumulative\\_gain](https://en.wikipedia.org/wiki/Discounted_cumulative_gain)

### 5.3 Normalized Cumulative Reciprocal Rank (NCRR)

```
class cornac.metrics.NCRR(k=-1)  
Normalized Cumulative Reciprocal Rank.  
  
Parameters k (int, optional, default: -1 (all)) – The number of items in the  
top@k list. If None, all items will be considered.
```

## 5.4 Mean Reciprocal Rank (MRR)

```
class cornac.metrics.MRR
```

Mean Reciprocal Rank.

**Parameters** `k` (*int, optional, default: -1 (all)*) – The number of items in the top@k list. If None, all items will be considered.

### References

[https://en.wikipedia.org/wiki/Mean\\_reciprocal\\_rank](https://en.wikipedia.org/wiki/Mean_reciprocal_rank)

## 5.5 Precision

```
class cornac.metrics.Precision(k=-1)
```

Precision@K.

**Parameters** `k` (*int, optional, default: -1 (all)*) – The number of items in the top@k list. If None, all items will be considered.

## 5.6 Recall

```
class cornac.metrics.Recall(k=-1)
```

Recall@K.

**Parameters** `k` (*int, optional, default: -1 (all)*) – The number of items in the top@k list. If None, all items will be considered.

## 5.7 Fmeasure (F1)

```
class cornac.metrics.FMeasure(k=-1)
```

F-measure@K@.

**Parameters** `k` (*int, optional, default: -1 (all)*) – The number of items in the top@k list. If None, all items will be considered.

## 5.8 Mean Absolute Error (MAE)

```
class cornac.metrics.MAE
```

Mean Absolute Error.

**name**

Name of the measure.

**Type** string, value: ‘MAE’

## 5.9 Root Mean Squared Error (RMSE)

```
class cornac.metrics.RMSE
```

Root Mean Squared Error.

**name**

Name of the measure.

**Type** string, value: ‘RMSE’



# CHAPTER 6

---

## Evaluation methods

---

### 6.1 Base Method

@author: Quoc-Tuan Truong <tuantq.vnu@gmail.com>

```
class cornac.eval_methods.base_method.BaseMethod(data=None, fmt='UIR',
                                                 rating_threshold=1.0, exclude_unknowns=False,
                                                 verbose=False, **kwargs)
```

Base Evaluation Method

#### Parameters

- **data** (*array-like*) – The original data.
- **data\_format** (*str*, *default*: ‘UIR’) – The format of given data.
- **total\_users** (*int*, *optional*, *default*: *None*) – Total number of unique users in the data including train, val, and test sets.
- **total\_items** – Total number of unique items in the data including train, val, and test sets.
- **rating\_threshold** (*float*, *optional*, *default*: *1.0*) – The threshold to convert ratings into positive or negative feedback for ranking metrics.
- **exclude\_unknowns** (*bool*, *optional*, *default*: *False*) – Ignore unknown users and items (cold-start) during evaluation.
- **verbose** (*bool*, *optional*, *default*: *False*) – Output running log

**evaluate** (*model*, *metrics*, *user\_based*)

Evaluate given models according to given metrics

#### Parameters

- **model** (*cornac.models.Recommender*) – Recommender model to be evaluated.
- **metrics** (*iterable*) – List of metrics.

- **user\_based** (`bool`) – Evaluation mode. Whether results are averaging based on number of users or number of ratings.

```
classmethod from_splits(train_data, test_data, val_data=None, data_format='UIR', rating_threshold=1.0, exclude_unknowns=False, verbose=False)
```

Constructing evaluation method given data.

#### Parameters

- **train\_data** (`array-like`) – Training data
- **test\_data** (`array-like`) – Test data
- **val\_data** (`array-like`) – Validation data
- **data\_format** (`str`, `default: 'UIR'`) – The format of given data.
- **rating\_threshold** (`float`, `default: 1.0`) – Threshold to decide positive or negative preferences.
- **exclude\_unknowns** (`bool`, `default: False`) – Whether to exclude unknown users/items in evaluation.
- **verbose** (`bool`, `default: False`) – The verbosity flag.

**Returns** `method` – Evaluation method object.

**Return type** `<cornac.eval_methods.BaseMethod>`

## 6.2 Ratio Split

@author: Quoc-Tuan Truong <[tuantq.vnu@gmail.com](mailto:tuantq.vnu@gmail.com)>

```
class cornac.eval_methods.ratio_split.RatioSplit(data, fmt='UIR', test_size=0.2, val_size=0.0, rating_threshold=1.0, shuffle=True, seed=None, exclude_unknowns=False, verbose=False, **kwargs)
```

Train-Test Split Evaluation Method.

#### Parameters

- **data** (`...`, `required`) – The input data in the form of triplets (user, item, rating).
- **fmt** (`str`, `optional`, `default: "UIR"`) – The format of input data: - UIR: (user, item, rating) triplet data - UIRT: (user, item , rating, timestamp) quadruplet data
- **test\_size** (`float`, `optional`, `default: 0.2`) – The proportion of the test set, if > 1 then it is treated as the size of the test set.
- **val\_size** (`float`, `optional`, `default: 0.0`) – The proportion of the validation set, if > 1 then it is treated as the size of the validation set.
- **rating\_threshold** (`float`, `optional`, `default: 1.`) – The minimum value that is considered to be a good rating used for ranking, e.g, if the ratings are in {1, ..., 5}, then rating\_threshold = 4.
- **shuffle** (`bool`, `optional`, `default: True`) – Shuffle the data before splitting.
- **seed** (`bool`, `optional`, `default: None`) – Random seed.

- **exclude\_unknowns** (`bool`, *optional*, *default*: `False`) – Ignore unknown users and items (cold-start) during evaluation and testing
  - **verbose** (`bool`, *optional*, *default*: `False`) – Output running log
- evaluate** (*model*, *metrics*, *user\_based*)  
Evaluate given models according to given metrics

**Parameters**

- **model** (`cornac.models.Recommender`) – Recommender model to be evaluated.
- **metrics** (`iterable`) – List of metrics.
- **user\_based** (`bool`) – Evaluation mode. Whether results are averaging based on number of users or number of ratings.

## 6.3 Cross Validation

@author: Aghiles Salah

```
class cornac.eval_methods.cross_validation.CrossValidation(data,      fmt='UIR',
                                                          n_folds=5,       rating_threshold=1.0,
                                                          partition=None, exclude_unknowns=True,
                                                          verbose=False,
                                                          **kwargs)
```

Cross Validation Evaluation Method.

**Parameters**

- **data** (.. , *required*) – Input data in the triplet format (user\_id, item\_id, rating\_val).
- **n\_folds** (`int`, *optional*, *default*: 5) – The number of folds for cross validation.
- **rating\_threshold** (`float`, *optional*, *default*: 1.) – The minimum value that is considered to be a good rating, e.g, if the ratings are in {1, ..., 5}, then rating\_threshold = 4.
- **partition** (`array-like, shape (n_observed_ratings,)`, *optional*, *default*: `None`) – The partition of ratings into n\_folds (fold label of each rating) If None, random partitioning is performed to assign each rating into a fold.
- **rating\_threshold** – The minimum value that is considered to be a good rating used for ranking, e.g, if the ratings are in {1, ..., 5}, then rating\_threshold = 4.
- **exclude\_unknowns** (`bool`, *optional*, *default*: `False`) – Ignore unknown users and items (cold-start) during evaluation and testing
- **verbose** (`bool`, *optional*, *default*: `False`) – Output running log

**evaluate** (*model*, *metrics*, *user\_based*)

Evaluate given models according to given metrics

**Parameters**

- **model** (`cornac.models.Recommender`) – Recommender model to be evaluated.
- **metrics** (`iterable`) – List of metrics.

- **user\_based** (*bool*) – Evaluation mode. Whether results are averaging based on number of users or number of ratings.

# CHAPTER 7

---

## Experiment

---

```
class cornac.experiment.Experiment(eval_method, models, metrics, user_based=True, verbose=False)
```

Experiment Class

### Parameters

- **eval\_method** (*BaseMethod object, required*) – The evaluation method (e.g., RatioSplit).
- **models** (*array of objects Recommender, required*) – A collection of recommender models to evaluate, e.g., [C2pf, Hpf, Pmf].
- **metrics** (*array of object metrics, required*) – A collection of metrics to use to evaluate the recommender models, e.g., [Ndcg, Mrr, Recall].
- **user\_based** (*bool, optional, default: True*) – Performance will be averaged based on number of users for rating metrics. If *False*, results will be averaged over number of ratings.
- **avg\_results** (*DataFrame, default: None*) – The average result per model.
- **user\_results** (*dictionary, default: {}*) – Results per user for each model. Result of user u, of metric m, of model d will be `user_results[d][m][u]`



# CHAPTER 8

---

## Built-in datasets

---

### 8.1 MovieLens

@author: Quoc-Tuan Truong <tuantq.vnu@gmail.com>

MovieLens: <https://grouplens.org/datasets/movielens/>

`cornac.datasets.movielens.load_100k(fmt='UIR')`

Load the MovieLens 100K dataset

**Parameters** `fmt` (`str`, `default: 'UIR'`) – Data format to be returned.

**Returns** `data` – Data in the form of a list of tuples depending on the given data format.

**Return type** array-like

`cornac.datasets.movielens.load_1m(fmt='UIR')`

Load the MovieLens 1M dataset

**Parameters** `fmt` (`str`, `default: 'UIR'`) – Data format to be returned.

**Returns** `data` – Data in the form of a list of tuples depending on the given data format.

**Return type** array-like

`cornac.datasets.movielens.load_plot()`

Load the plots of movies provided @ <http://dm.postech.ac.kr/~cartopy/ConvMF/>

**Returns** `movie_plots` – A dictionary with keys are movie ids and values are text plots.

**Return type** Dict

### 8.2 Netflix

@author: Quoc-Tuan Truong <tuantq.vnu@gmail.com>

Data: <https://www.kaggle.com/netflix-inc/netflix-prize-data/>

```
cornac.datasets.netflix.load_data(fmt='UIR')
```

Load the Netflix entire dataset - Number of ratings: 100,480,507 - Number of users: 480,189 - Number of items: 17,770

**Parameters** `fmt` (`str`, `default: 'UIR'`) – Data format to be returned.

**Returns** `data` – Data in the form of a list of tuples depending on the given data format.

**Return type** array-like

```
cornac.datasets.netflix.load_data_small(fmt='UIR')
```

Load a small subset of the Netflix dataset. We draw this subsample such that every user has at least 10 items and each item has at least 10 users. - Number of ratings: 607,803 - Number of users: 10,000 - Number of items: 5,000

**Parameters** `fmt` (`str`, `default: 'UIR'`) – Data format to be returned.

**Returns** `data` – Data in the form of a list of tuples depending on the given data format.

**Return type** array-like

## 8.3 Tradesy

@author: Quoc-Tuan Truong <tuantq.vnu@gmail.com>

Original data: <http://jmcauley.ucsd.edu/data/tradesy/> This data is used in the VBPR paper. After cleaning the data, we have: - Number of feedback: 394,421 (410,186 is reported but there are duplicates) - Number of users: 19,243 (19,823 is reported due to duplicates) - Number of items: 165,906 (166,521 is reported due to duplicates)

```
cornac.datasets.tradesy.load_data()
```

Load the feedback observations

**Returns** `data` – Data in the form of a list of tuples (user, item, 1).

**Return type** array-like

```
cornac.datasets.tradesy.load_feature()
```

Load the item visual feature

**Returns** `data` – Item-feature dictionary. Each feature vector is a Numpy array of size 4096.

**Return type** dict

## 8.4 Amazon Office

@author: Aghiles Salah <asalah@smu.edu.sg>

This data is built based on the Amazon datasets provided by Julian McAuley at: <http://jmcauley.ucsd.edu/data/amazon/>

```
cornac.datasets.amazon_office.load_context(data_format='UIR')
```

Load the item-item interactions

**Parameters** `data_format` (`str`, `default: 'UIR'`) – Data format to be returned.

**Returns** `data` – Data in the form of a list of tuples depending on the specified data format.

**Return type** array-like

```
cornac.datasets.amazon_office.load_rating(data_format='UIR')
```

Load the user-item ratings

**Parameters** `data_format` (`str`, `default: 'UIR'`) – Data format to be returned.

**Returns** `data` – Data in the form of a list of tuples depending on the specified data format.

**Return type** array-like



---

## Python Module Index

---

### C

cornac.data, 7  
cornac.data.graph, 16  
cornac.data.image, 19  
cornac.data.reader, 19  
cornac.data.testset, 15  
cornac.data.text, 16  
cornac.data.trainset, 12  
cornac.datasets, 49  
cornac.datasets.amazon\_office, 50  
cornac.datasets.movielens, 49  
cornac.datasets.netflix, 49  
cornac.datasets.tradesy, 50  
cornac.eval\_methods, 43  
cornac.eval\_methods.base\_method, 43  
cornac.eval\_methods.cross\_validation,  
    45  
cornac.eval\_methods.ratio\_split, 44  
cornac.experiment, 47  
cornac.metrics, 39  
cornac.models, 21  
cornac.models.c2pf.recom\_c2pf, 22  
cornac.models.cdl.recom\_cdl, 30  
cornac.models.coe.recom\_coe, 26  
cornac.models.conv\_mf.recom\_convvmf, 36  
cornac.models.hpf.recom\_hpf, 31  
cornac.models.ibpr.recom\_ibpr, 24  
cornac.models.mf.recom\_mf, 35  
cornac.models.online\_ibpr.recom\_online\_ibpr,  
    25  
cornac.models.pcrl.recom\_pcrl, 21  
cornac.models.pmf.recom\_pmf, 33  
cornac.models.skm.recom\_skmeans, 29  
cornac.models.vbpr.recom\_vbpr, 27



---

## Index

---

### A

AUC (*class in cornac.metrics*), 39

### B

BaseMethod                   (class                   in  
                                *cornac.eval\_methods.base\_method*), 43  
BaseTokenizer (*class in cornac.data.text*), 16  
batch() (*cornac.data.graph.GraphModule method*),  
      16  
batch() (*cornac.data.GraphModule method*), 8  
batch\_feature()           (*cornac.data.FeatureModule  
                                method*), 7  
batch\_image()              (*cornac.data.image.ImageModule  
                                method*), 19  
batch\_image() (*cornac.data.ImageModule method*),  
      8  
batch\_seq()                (*cornac.data.text.TextModule method*),  
      19  
batch\_seq() (*cornac.data.TextModule method*), 8  
batch\_tfidf()              (*cornac.data.text.TextModule  
                                method*), 19  
batch\_tfidf() (*cornac.data.TextModule method*), 8  
batch\_tokenize()            (*cornac.data.text.BaseTokenizer  
                                method*), 16  
batch\_tokenize()            (*cornac.data.text.Tokenizer  
                                method*), 16  
BPR (*class in cornac.models.bpr.recom\_bpr*), 32  
build() (*cornac.data.FeatureModule method*), 7  
build() (*cornac.data.graph.GraphModule method*),  
      16  
build() (*cornac.data.GraphModule method*), 8  
build() (*cornac.data.image.ImageModule method*), 19  
build() (*cornac.data.ImageModule method*), 8  
build() (*cornac.data.text.TextModule method*), 19  
build() (*cornac.data.TextModule method*), 8

### C

C2PF (*class in cornac.models.c2pf.recom\_c2pf*), 22  
CDL (*class in cornac.models.cdl.recom\_cdl*), 30

COE (*class in cornac.models.coe.recom\_coe*), 26  
ConvMF (*class in cornac.models.conv\_mf.recom\_convmf*),  
      36  
cornac.data (*module*), 7  
cornac.data.graph (*module*), 16  
cornac.data.image (*module*), 19  
cornac.data.reader (*module*), 19  
cornac.data.testset (*module*), 15  
cornac.data.text (*module*), 16  
cornac.data.trainset (*module*), 12  
cornac.datasets (*module*), 49  
cornac.datasets.amazon\_office (*module*), 50  
cornac.datasets.movielens (*module*), 49  
cornac.datasets.netflix (*module*), 49  
cornac.datasets.tradesy (*module*), 50  
cornac.eval\_methods (*module*), 43  
cornac.eval\_methods.base\_method (*module*),  
      43  
cornac.eval\_methods.cross\_validation  
                                (*module*), 45  
cornac.eval\_methods.ratio\_split (*module*),  
      44  
cornac.experiment (*module*), 47  
cornac.metrics (*module*), 39  
cornac.models (*module*), 21  
cornac.models.c2pf.recom\_c2pf (*module*), 22  
cornac.models.cdl.recom\_cdl (*module*), 30  
cornac.models.coe.recom\_coe (*module*), 26  
cornac.models.conv\_mf.recom\_convmf (*mod-  
ule*), 36  
cornac.models.hpf.recom\_hpf (*module*), 31  
cornac.models.ibpr.recom\_ibpr (*module*), 24  
cornac.models.mf.recom\_mf (*module*), 35  
cornac.models.online\_ibpr.recom\_online\_ibpr  
                                (*module*), 25  
cornac.models.pcrl.recom\_pcrl (*module*), 21  
cornac.models.pmf.recom\_pmf (*module*), 33  
cornac.models.skm.recom\_skmeans (*module*),  
      29  
cornac.models.vbpr.recom\_vbpr (*module*), 27

```

CountVectorizer (class in cornac.data.text), 17
CrossValidation (class in cornac.eval_methods.cross_validation), 45
from_uir() (cornac.data.trainset.MatrixTrainSet
            class method), 12

E
evaluate() (cornac.eval_methods.base_method.BaseMethod
            method), 43
evaluate() (cornac.eval_methods.cross_validation.CrossValidation
            method), 45
evaluate() (cornac.eval_methods.ratio_split.RatioSplit
            method), 45
Experiment (class in cornac.experiment), 47

F
feature_dim (cornac.data.FeatureModule attribute),
              7
FeatureModule (class in cornac.data), 7
features (cornac.data.FeatureModule attribute), 7
fit (cornac.models.bpr.recom_bpr.BPR attribute), 33
fit (cornac.models.mf.recom_mf.MF attribute), 35
fit () (cornac.data.text.CountVectorizer method), 18
fit () (cornac.models.c2pf.recom_c2pf.C2PF method),
       23
fit () (cornac.models.cdl.recom_cdl.CDL method), 31
fit () (cornac.models.coe.recom_coe.COE method), 27
fit () (cornac.models.conv_mf.recom_convmf.ConvMF
       method), 37
fit () (cornac.models.hpf.recom_hpf.HPF method), 32
fit () (cornac.models.ibpr.recom_ibpr.IBPR method),
       24
fit () (cornac.models.online_ibpr.recom_online_ibpr.OnlineIBPR
       method), 26
fit () (cornac.models.pcrl.recom_pcrl.PCRL method),
       22
fit () (cornac.models.pmf.recom_pmf.PMF method),
       34
fit () (cornac.models.skm.recom_skmeans.SKMeans
       method), 29
fit () (cornac.models.vbpr.recom_vbpr.VBPR method),
       28
fit_transform() (cornac.data.text.CountVectorizer
                 method), 18
FMeasure (class in cornac.metrics), 40
from_sequences() (cornac.data.text.Vocabulary
                  class method), 17
from_splits() (cornac.eval_methods.base_method.BaseMethod
                  class method), 44
from_tokens() (cornac.data.text.Vocabulary class
                  method), 17
from_uir() (cornac.data.MatrixTrainSet class
                  method), 9
from_uir() (cornac.data.TestSet class method), 11
from_uir() (cornac.data.testset.TestSet class
                  method), 15
from_uir() (cornac.data.trainset.TrainSet
                  class method), 8
get_iid() (cornac.data.TestSet method), 12
get_iid() (cornac.data.testset.TestSet method), 16
get_iid() (cornac.data.TrainSet method), 8
get_ratings() (cornac.data.TestSet method), 12
get_ratings() (cornac.data.testset.TestSet method),
               16
get_train_triplet()
    (cornac.data.graph.GraphModule method), 16
get_train_triplet() (cornac.data.GraphModule
                     method), 8
get_uid() (cornac.data.TestSet method), 12
get_uid() (cornac.data.testset.TestSet method), 16
get_uid() (cornac.data.TrainSet method), 9
get_uid() (cornac.data.trainset.TrainSet method), 14
GraphModule (class in cornac.data), 8
GraphModule (class in cornac.data.graph), 16

H
HPF (class in cornac.models.hpf.recom_hpf), 31

I
IBPR (class in cornac.models.ibpr.recom_ibpr), 24
idx_iter() (cornac.data.TrainSet static method), 9
idx_iter() (cornac.data.trainset.TrainSet static
            method), 14
idx_iter (cornac.data.TrainSet attribute), 9
iid_list (cornac.data.trainset.TrainSet attribute), 14
ImageModule (class in cornac.data), 8
ImageModule (class in cornac.data.image), 19
is_unk_item() (cornac.data.TrainSet method), 9
is_unk_item() (cornac.data.trainset.TrainSet
               method), 14
is_unk_user() (cornac.data.TrainSet method), 9
is_unk_user() (cornac.data.trainset.TrainSet
               method), 14
item_iter() (cornac.data.MatrixTrainSet method),
             10
item_iter() (cornac.data.trainset.MatrixTrainSet
             method), 13
load() (cornac.data.text.Vocabulary class method), 17
load_100k() (in module cornac.datasets.movielens),
             49
load_1m() (in module cornac.datasets.movielens), 49
load_context() (in module cornac.datasets.amazon_office),
                50
load_data() (in module cornac.datasets.netflix), 49

```

load\_data () (in module cornac.datasets.tradesy), 50  
 load\_data\_small () (in module cornac.datasets.netflix), 50  
 load\_feature () (in module cornac.datasets.tradesy), 50  
 load\_plot () (in module cornac.datasets.movieLens), 49  
 load\_rating () (in module cornac.datasets.amazon\_office), 50

**M**

MAE (class in cornac.metrics), 40  
 MatrixTrainSet (class in cornac.data), 9  
 MatrixTrainSet (class in cornac.data.trainset), 12  
 MF (class in cornac.models.mf.recom\_mf), 35  
 MRR (class in cornac.metrics), 40  
 MultimodalTestSet (class in cornac.data), 12  
 MultimodalTestSet (class in cornac.data.testset), 15  
 MultimodalTrainSet (class in cornac.data), 11  
 MultimodalTrainSet (class in cornac.data.trainset), 14

**N**

name (cornac.metrics.MAE attribute), 40  
 name (cornac.metrics.RMSE attribute), 41  
 NCRR (class in cornac.metrics), 39  
 NDCG (class in cornac.metrics), 39  
 num\_items (cornac.data.TrainSet attribute), 9  
 num\_items (cornac.data.trainset.TrainSet attribute), 14  
 num\_users (cornac.data.TrainSet attribute), 9  
 num\_users (cornac.data.trainset.TrainSet attribute), 14

**O**

OnlineIBPR (class in cornac.models.online\_ibpr.recom\_online\_ibpr), 25

**P**

PCRL (class in cornac.models.pcrl.recom\_pcrl), 21  
 PMF (class in cornac.models.pmf.recom\_pmf), 33  
 Precision (class in cornac.metrics), 40

**R**

RatioSplit (class in cornac.eval\_methods.ratio\_split), 44  
 raw\_iid\_list (cornac.data.TrainSet attribute), 9  
 raw\_iid\_list (cornac.data.trainset.TrainSet attribute), 15  
 raw\_uid\_list (cornac.data.TrainSet attribute), 9  
 raw\_uid\_list (cornac.data.trainset.TrainSet attribute), 15

read\_ui () (in module cornac.data.reader), 19  
 read\_uir () (in module cornac.data.reader), 19  
 Recall (class in cornac.metrics), 40  
 RMSE (class in cornac.metrics), 41

**S**

save () (cornac.data.text.Vocabulary method), 17  
 score (cornac.models.bpr.recom\_bpr.BPR attribute), 33  
 score (cornac.models.mf.recom\_mf.MF attribute), 35  
 score () (cornac.models.c2pf.recom\_c2pf.C2PF method), 23  
 score () (cornac.models.cdl.recom\_cdl.CDL method), 31  
 score () (cornac.models.coe.recom\_coe.COE method), 27  
 score () (cornac.models.conv\_mf.recom\_convMF.ConvMF method), 37  
 score () (cornac.models.hpf.recom\_hpf.HPF method), 32  
 score () (cornac.models.ibpr.recom\_ibpr.IBPR method), 25  
 score () (cornac.models.online\_ibpr.recom\_online\_ibpr.OnlineIBPR method), 26  
 score () (cornac.models.pcrl.recom\_pcrl.PCRL method), 22  
 score () (cornac.models.pmf.recom\_pmf.PMF method), 34  
 score () (cornac.models.skm.recom\_skmmeans.SKMeans method), 29  
 score () (cornac.models.vbpr.recom\_vbpr.VBPR method), 28  
 SKMeans (class in cornac.models.skm.recom\_skmmeans), 29

**T**

TestSet (class in cornac.data), 11  
 TestSet (class in cornac.data.testset), 15  
 TextModule (class in cornac.data), 7  
 TextModule (class in cornac.data.text), 18  
 to\_idx () (cornac.data.text.Vocabulary method), 17  
 to\_text () (cornac.data.text.Vocabulary method), 17  
 tokenize () (cornac.data.text.BaseTokenizer method), 17  
 tokenize () (cornac.data.text.Tokenizer method), 16  
 Tokenizer (class in cornac.data.text), 16  
 TrainSet (class in cornac.data), 8  
 TrainSet (class in cornac.data.trainset), 14  
 transform () (cornac.data.text.CountVectorizer method), 18

**U**

uid\_list (cornac.data.TrainSet attribute), 9  
 uid\_list (cornac.data.trainset.TrainSet attribute), 15  
 uij\_iter () (cornac.data.MatrixTrainSet method), 10

uij\_iter() (*cornac.data.trainset.MatrixTrainSet method*), 13  
uir\_iter() (*cornac.data.MatrixTrainSet method*), 10  
uir\_iter() (*cornac.data.trainset.MatrixTrainSet method*), 13  
user\_iter() (*cornac.data.MatrixTrainSet method*),  
10  
user\_iter() (*cornac.data.trainset.MatrixTrainSet method*), 13  
users (*cornac.data.TestSet attribute*), 12  
users (*cornac.data.testset.TestSet attribute*), 16

## V

VBPR (*class in cornac.models.vbpr.recom\_vbpr*), 27  
Vocabulary (*class in cornac.data.text*), 17