

---

**skrobot**

***Release 1.0.8***

**Medoid AI**

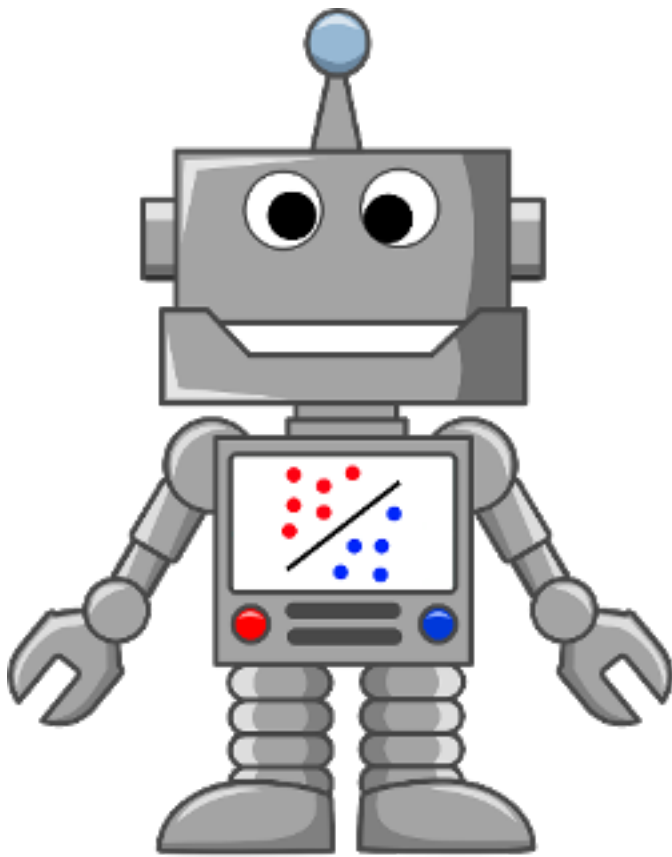
**Nov 29, 2020**



## CONTENTS

<b>1</b>	<b>API Reference</b>	<b>3</b>
<b>2</b>	<b>What is it about?</b>	<b>23</b>
<b>3</b>	<b>Why does it exists?</b>	<b>25</b>
<b>4</b>	<b>How do I install it?</b>	<b>27</b>
<b>5</b>	<b>Which are the components?</b>	<b>29</b>
<b>6</b>	<b>How do I use it?</b>	<b>31</b>
<b>7</b>	<b>Sample of generated results?</b>	<b>37</b>
<b>8</b>	<b>The people behind it?</b>	<b>51</b>
<b>9</b>	<b>Can I contribute?</b>	<b>53</b>
<b>10</b>	<b>What license do you use?</b>	<b>55</b>
	<b>Python Module Index</b>	<b>57</b>
	<b>Index</b>	<b>59</b>





**skrobot**



## API REFERENCE

### 1.1 skrobot package

#### 1.1.1 Subpackages

skrobot.core package

Submodules

skrobot.core.experiment module

**class** skrobot.core.experiment.**Experiment** (*experiments\_repository*)

Bases: object

The *Experiment* class can be used to build, track and run an experiment.

It can run *BaseTask* tasks in the context of an experiment.

When building an experiment and/or running tasks, various metadata as well as task-related files are stored for tracking experiments.

Lastly, an experiment can be configured to send notifications when running a task, which can be useful for teams who need to get notified for the progress of the experiment.

**\_\_init\_\_** (*experiments\_repository*)

This is the constructor method and can be used to create a new object instance of *Experiment* class.

**Parameters** **experiments\_repository** (*str*) – The root directory path under which a unique directory is created for the experiment.

**set\_notifier** (*notifier*: skrobot.notification.base\_notifier.BaseNotifier)

Optional method.

Set the experiment's notifier.

**Parameters** **notifier** (*BaseNotifier*) – The experiment's notifier.

**Returns** The object instance itself.

**Return type** *Experiment*

**set\_source\_code\_file\_path** (*source\_code\_file\_path*)

Optional method.

Set the experiment's source code file path.

**Parameters** **source\_code\_file\_path** (*str*) – The experiment's source code file path.

**Returns** The object instance itself.

**Return type** *Experiment*

**set\_experimenter** (*experimenter*)

Optional method.

Set the experimenter's name.

By default the experimenter's name is *anonymous*. However, if you want to override it you can pass a new name.

**Parameters** **experimenter** (*str*) – The experimenter's name.

**Returns** The object instance itself.

**Return type** *Experiment*

**build** ()

Build the *Experiment*.

When an experiment is built, it creates a unique directory under which it stores various experiment-related metadata and files for tracking reasons.

Specifically, under the experiment's directory an *experiment.log* JSON file is created, which contains a unique auto-generated experiment ID, the current date & time, and the experimenter's name.

Also, the experiment's directory name contains the experimenter's name as well as current date & time.

Lastly, in case *set\_source\_code\_file\_path()* is used, the experiment's source code file is copied also under the experiment's directory.

**Returns** The object instance itself.

**Return type** *Experiment*

**run** (*task*)

Run a *BaseTask* task.

When running a task, its recorded parameters (e.g., *train\_task.params*) and any other task-related generated files are stored under experiment's directory for tracking reasons.

The task's recorded parameters are in JSON format.

Also, in case *set\_notifier()* is used to set a notifier, a notification is sent for the success or failure (including the error message) of the task's execution.

Lastly, in case an exception occurs, a text file (e.g., *train\_task.errors*) is generated under experiment's directory containing the error message.

**Parameters** **task** (*BaseTask*) – The task to run.

**Returns** The task's result.

**Return type** Depends on the *task* parameter.



## skrobot.core.task\_runner module

**class** skrobot.core.task\_runner.**TaskRunner** (*output\_directory\_path*)

Bases: object

The *TaskRunner* class is a simplified version (in functionality) of the *Experiment* class.

It leaves out all the “experiment” stuff and is focused mostly in the execution and tracking of *BaseTask* tasks.

**\_\_init\_\_** (*output\_directory\_path*)

This is the constructor method and can be used to create a new object instance of *TaskRunner* class.

**Parameters** *output\_directory\_path* (*str*) – The output directory path under which task-related generated files are stored.

**run** (*task*)

Run a *BaseTask* task.

When running a task, its recorded parameters (e.g., *train\_task.params*) and any other task-related generated files are stored under output directory for tracking reasons.

The task’s recorded parameters are in JSON format.

Lastly, in case an exception occurs, a text file (e.g., *train\_task.errors*) is generated under output directory containing the error message.

**Parameters** *task* (*BaseTask*) – The task to run.

**Returns** The task’s result.

**Return type** Depends on the *task* parameter.

## skrobot.feature\_selection package

### Submodules

#### skrobot.feature\_selection.column\_selector module

**class** skrobot.feature\_selection.column\_selector.**ColumnSelector** (*cols*,  
*drop\_axis=False*)

Bases: sklearn.base.BaseEstimator

The *ColumnSelector* class is an implementation of a column selector for scikit-learn pipelines.

It can be used for manual feature selection to select specific columns from an input data set.

It can select columns either by integer indices or by names.

**\_\_init\_\_** (*cols*, *drop\_axis=False*)

This is the constructor method and can be used to create a new object instance of *ColumnSelector* class.

#### Parameters

- **cols** (*list*) – A non-empty list specifying the columns to be selected. For example, [1, 4, 5] to select the 2nd, 5th, and 6th columns, and ['A','C','D'] to select the columns A, C and D.
- **drop\_axis** (*bool*, *optional*) – Can be used to reshape the output data set from (n\_samples, 1) to (n\_samples) by dropping the last axis. It defaults to False.

**fit\_transform** (*X*, *y=None*)

Returns a slice of the input data set.

**Parameters**

- **X** ({NumPy array, pandas DataFrame, SciPy sparse matrix}) – Input vectors of shape (n\_samples, n\_features), where n\_samples is the number of samples and n\_features is the number of features.
- **y** (None) – Ignored.

**Returns** Subset of the input data set of shape (n\_samples, k\_features), where n\_samples is the number of samples and k\_features <= n\_features.

**Return type** {NumPy array, SciPy sparse matrix}

**t\_transform** (*X*, *y=None*)

Returns a slice of the input data set.

**Parameters**

- **X** ({NumPy array, pandas DataFrame, SciPy sparse matrix}) – Input vectors of shape (n\_samples, n\_features), where n\_samples is the number of samples and n\_features is the number of features.
- **y** (None) – Ignored.

**Returns** Subset of the input data set of shape (n\_samples, k\_features), where n\_samples is the number of samples and k\_features <= n\_features.

**Return type** {NumPy array, SciPy sparse matrix}

**fit** (*X*, *y=None*)

This is a mock method and does nothing.

**Parameters**

- **X** (None) – Ignored.
- **y** (None) – Ignored.

**Returns** The object instance itself.

**Return type** *ColumnSelector*

## skrobot.notification package

### Submodules

#### skrobot.notification.base\_notifier module

**class** skrobot.notification.base\_notifier.**BaseNotifier**

Bases: abc.ABC

The *BaseNotifier* is an abstract base class for implementing notifiers.

A notifier can be used to send notifications.

**abstract notify** (*message*)

An abstract method for sending the notification.

**Parameters** **message** (*str*) – The notification's message.

## skrobot.tasks package

### Submodules

#### skrobot.tasks.base\_cross\_validation\_task module

**class** skrobot.tasks.base\_cross\_validation\_task.**BaseCrossValidationTask** (*type\_name*, *args*)

Bases: *skrobot.tasks.base\_task.BaseTask*

The *BaseCrossValidationTask* is an abstract base class for implementing tasks that use cross-validation functionality.

It can support both stratified k-fold cross-validation as well as cross-validation with user-defined folds.

By default, stratified k-fold cross-validation is used with the default parameters of *stratified\_folds()* method.

**\_\_init\_\_** (*type\_name*, *args*)

This is the constructor method and can be used from child *BaseCrossValidationTask* implementations.

#### Parameters

- **type\_name** (*str*) – The task's type name. A common practice is to pass the name of the task's class.
- **args** (*dict*) – The task's parameters. A common practice is to pass the parameters at the time of task's object creation. It is a dictionary of key-value pairs, where the key is the parameter name and the value is the parameter value.

**custom\_folds** (*folds\_file\_path*, *fold\_column*='fold')

Optional method.

Use cross-validation with user-defined custom folds.

#### Parameters

- **folds\_file\_path** (*str*) – The path to the file containing the user-defined folds for the samples. The file needs to be formatted with the same separation delimiter (comma for CSV, tab for TSV, etc.) as the one used in the input data set files provided to the task. The file must contain two data columns and the first row must be the header. The first column is for the sample IDs and needs to be the same as the one used in the input data set files provided to the task. The second column is for the fold IDs (e.g., 1 through 5, A through D, etc.).
- **fold\_column** (*str*, *optional*) – The column name for the fold IDs. It defaults to 'fold'.

**Returns** The object instance itself.

**Return type** *BaseCrossValidationTask*

**stratified\_folds** (*total\_folds*=3, *shuffle*=False)

Optional method.

Use stratified k-fold cross-validation.

The folds are made by preserving the percentage of samples for each class.

#### Parameters

- **total\_folds** (*int*, *optional*) – Number of folds. Must be at least 2. It defaults to 3.
- **shuffle** (*bool*, *optional*) – Whether to shuffle each class's samples before splitting into batches. Note that the samples within each split will not be shuffled. It defaults to False.

**Returns** The object instance itself.

**Return type** *BaseCrossValidationTask*

**get\_configuration()**

Get the task's parameters.

**Returns** The task's parameters as a dictionary of key-value pairs, where the key is the parameter name and the value is the parameter value.

**Return type** dict

**get\_type()**

Get the task's type name.

**Returns** The task's type name.

**Return type** str

**abstract run** (*output\_directory*)

An abstract method for running the task.

**Parameters** **output\_directory** (*str*) – The output directory path under which task-related generated files are stored.

## skrobot.tasks.base\_task module

**class** skrobot.tasks.base\_task.**BaseTask** (*type\_name*, *args*)

Bases: abc.ABC

The *BaseTask* is an abstract base class for implementing tasks.

A task is a configurable and reproducible piece of code built on top of scikit-learn that can be used in machine learning pipelines.

**\_\_init\_\_** (*type\_name*, *args*)

This is the constructor method and can be used from child *BaseTask* implementations.

**Parameters**

- **type\_name** (*str*) – The task's type name. A common practice is to pass the name of the task's class.
- **args** (*dict*) – The task's parameters. A common practice is to pass the parameters at the time of task's object creation. It is a dictionary of key-value pairs, where the key is the parameter name and the value is the parameter value.

**get\_type()**

Get the task's type name.

**Returns** The task's type name.

**Return type** str

**get\_configuration()**

Get the task's parameters.

**Returns** The task’s parameters as a dictionary of key-value pairs, where the key is the parameter name and the value is the parameter value.

**Return type** dict

**abstract run** (*output\_directory*)

An abstract method for running the task.

**Parameters** **output\_directory** (*str*) – The output directory path under which task-related generated files are stored.

## skrobot.tasks.evaluation\_cross\_validation\_task module

```
class skrobot.tasks.evaluation_cross_validation_task.EvaluationCrossValidationTask (estimator,
train_data,
test_data,
es-
ti-
ma-
tor_params,
field_delim,
fea-
ture_columns,
id_columns,
la-
bel_columns,
ran-
dom_seed=
thresh-
old_selection,
met-
ric_greater,
thresh-
old_tuning,
1.0,
0.01),
ex-
port_classification,
ex-
port_confusion,
ex-
port_roc_curve,
ex-
port_pr_curve,
ex-
port_false_
ex-
port_false_
ex-
port_also_
fs-
core_beta=
```

Bases: *skrobot.tasks.base\_cross\_validation\_task.BaseCrossValidationTask*

The *EvaluationCrossValidationTask* class can be used to evaluate a scikit-learn estimator/pipeline

on some data.

The following evaluation results can be generated on-demand for hold-out test data set as well as train/validation cross-validation folds:

- PR / ROC Curves
- Confusion Matrixes
- Classification Reports
- Performance Metrics
- False Positives
- False Negatives

It can support both stratified k-fold cross-validation as well as cross-validation with user-defined folds.

By default, stratified k-fold cross-validation is used with the default parameters of `stratified_folds()` method.

```
__init__(estimator, train_data_set_file_path, test_data_set_file_path=None, estimator_params=None, field_delimiter=',', feature_columns='all', id_column='id', label_column='label', random_seed=42, threshold_selection_by='f1', metric_greater_is_better=True, threshold_tuning_range=(0.01, 1.0, 0.01), export_classification_reports=False, export_confusion_matrixes=False, export_roc_curves=False, export_pr_curves=False, export_false_positives_reports=False, export_false_negatives_reports=False, export_also_for_train_folds=False, fscore_beta=1)
```

This is the constructor method and can be used to create a new object instance of `EvaluationCrossValidationTask` class.

#### Parameters

- **estimator** (*scikit-learn {estimator, pipeline}*) – It can be either an estimator (e.g., LogisticRegression) or a pipeline ending with an estimator. The estimator needs to be able to predict probabilities through a `predict_proba` method.
- **train\_data\_set\_file\_path** (*str*) – The file path of the input train data set. It can be either a URL or a disk file path.
- **test\_data\_set\_file\_path** (*str, optional*) – The file path of the input test data set. It can be either a URL or a disk file path. It defaults to `None`.
- **estimator\_params** (*dict, optional*) – The parameters to override in the provided estimator/pipeline. It defaults to `None`.
- **field\_delimiter** (*str, optional*) – The separation delimiter (comma for CSV, tab for TSV, etc.) used in the input train/test data set files. It defaults to `','`.
- **feature\_columns** (*{str, list}, optional*) – Either `'all'` to use from the input train/test data set files all the columns or a list of column names to select specific columns. It defaults to `'all'`.
- **id\_column** (*str, optional*) – The name of the column in the input train/test data set files containing the sample IDs. It defaults to `'id'`.
- **label\_column** (*str, optional*) – The name of the column in the input train/test data set files containing the ground truth labels. It defaults to `'label'`.
- **random\_seed** (*int, optional*) – The random seed used in the random number generator. It can be used to reproduce the output. It defaults to 42.

- **threshold\_selection\_by** (*{str, float}, optional*) – The evaluation results will be generated either for a specific provided threshold value (e.g., 0.49) or for the best threshold found from threshold tuning, based on a specific provided metric (e.g., 'f1', 'f0.55'). It defaults to 'f1'.
- **metric\_greater\_is\_better** (*bool, optional*) – This flag will control the direction of searching of the best threshold and it depends on the provided metric in `threshold_selection_by`. True, means that greater metric values is better and False means the opposite. It defaults to True.
- **threshold\_tuning\_range** (*tuple, optional*) – A range in form (start\_value, stop\_value, step\_size) for generating a sequence of threshold values in threshold tuning. It generates the sequence by incrementing the start value using the step size until it reaches the stop value. It defaults to (0.01, 1.0, 0.01).
- **export\_classification\_reports** (*bool, optional*) – If this task will export classification reports. It defaults to False.
- **export\_confusion\_matrixes** (*bool, optional*) – If this task will export confusion matrixes. It defaults to False.
- **export\_roc\_curves** (*bool, optional*) – If this task will export ROC curves. It defaults to False.
- **export\_pr\_curves** (*bool, optional*) – If this task will export PR curves. It defaults to False.
- **export\_false\_positives\_reports** (*bool, optional*) – If this task will export false positives reports. It defaults to False.
- **export\_false\_negatives\_reports** (*bool, optional*) – If this task will export false negatives reports. It defaults to False.
- **export\_also\_for\_train\_folds** (*bool, optional*) – If this task will export the evaluation results also for the train folds of cross-validation. It defaults to False.
- **fscore\_beta** (*float, optional*) – The beta parameter in F-measure. It determines the weight of recall in the score.  $\beta < 1$  lends more weight to precision, while  $\beta > 1$  favors recall ( $\beta \rightarrow 0$  considers only precision,  $\beta \rightarrow +\infty$  only recall). It defaults to 1.

**run** (*output\_directory*)

Run the task.

All of the evaluation results are stored as files under the output directory path.

**Parameters** **output\_directory** (*str*) – The output directory path under which task-related generated files are stored.

**Returns** The task's result. Specifically, the threshold used along with its related performance metrics and summary metrics from all cross-validation splits as well as hold-out test data set.

**Return type** dict

**custom\_folds** (*folds\_file\_path, fold\_column='fold'*)

Optional method.

Use cross-validation with user-defined custom folds.

**Parameters**

- **folds\_file\_path** (*str*) – The path to the file containing the user-defined folds for the samples. The file needs to be formatted with the same separation delimiter (comma for

CSV, tab for TSV, etc.) as the one used in the input data set files provided to the task. The file must contain two data columns and the first row must be the header. The first column is for the sample IDs and needs to be the same as the one used in the input data set files provided to the task. The second column is for the fold IDs (e.g., 1 through 5, A through D, etc.).

- **fold\_column** (*str*, *optional*) – The column name for the fold IDs. It defaults to 'fold'.

**Returns** The object instance itself.

**Return type** *BaseCrossValidationTask*

**get\_configuration()**

Get the task's parameters.

**Returns** The task's parameters as a dictionary of key-value pairs, where the key is the parameter name and the value is the parameter value.

**Return type** dict

**get\_type()**

Get the task's type name.

**Returns** The task's type name.

**Return type** str

**stratified\_folds** (*total\_folds=3*, *shuffle=False*)

Optional method.

Use stratified k-fold cross-validation.

The folds are made by preserving the percentage of samples for each class.

**Parameters**

- **total\_folds** (*int*, *optional*) – Number of folds. Must be at least 2. It defaults to 3.
- **shuffle** (*bool*, *optional*) – Whether to shuffle each class's samples before splitting into batches. Note that the samples within each split will not be shuffled. It defaults to False.

**Returns** The object instance itself.

**Return type** *BaseCrossValidationTask*



**skrobot.tasks.feature\_selection\_cross\_validation\_task module**

```
class skrobot.tasks.feature_selection_cross_validation_task.FeatureSelectionCrossValidationTask
```

Bases: *skrobot.tasks.base\_cross\_validation\_task.BaseCrossValidationTask*

The *FeatureSelectionCrossValidationTask* class can be used to perform feature selection with Recursive Feature Elimination using a scikit-learn estimator on some data.

A scikit-learn preprocessor can be used on the input train data set before feature selection runs.

It can support both stratified k-fold cross-validation as well as cross-validation with user-defined folds.

By default, stratified k-fold cross-validation is used with the default parameters of *stratified\_folds()* method.

```
__init__(estimator, train_data_set_file_path, estimator_params=None, field_delimiter=',', preprocessor=None, preprocessor_params=None, min_features_to_select=1, scoring='f1', feature_columns='all', id_column='id', label_column='label', random_seed=42, verbose=3, n_jobs=1)
```

This is the constructor method and can be used to create a new object instance of *FeatureSelectionCrossValidationTask* class.

**Parameters**

- **estimator** (*scikit-learn estimator*) – An estimator (e.g., LogisticRegression). It needs to provide feature importances through either a `coef_` or a `feature_importances_` attribute.
- **train\_data\_set\_file\_path** (*str*) – The file path of the input train data set. It can be either a URL or a disk file path.

- **estimator\_params** (*dict, optional*) – The parameters to override in the provided estimator. It defaults to None.
- **field\_delimiter** (*str, optional*) – The separation delimiter (comma for CSV, tab for TSV, etc.) used in the input train data set file. It defaults to ‘,’.
- **preprocessor** (*scikit-learn preprocessor, optional*) – The preprocessor you want to run on the input train data set before feature selection. You can set for example a scikit-learn ColumnTransformer, OneHotEncoder, etc. It defaults to None.
- **preprocessor\_params** (*dict, optional*) – The parameters to override in the provided preprocessor. It defaults to None.
- **min\_features\_to\_select** (*int, optional*) – The minimum number of features to be selected. This number of features will always be scored. It defaults to 1.
- **scoring** (*{str, callable}, optional*) – A single scikit-learn scorer string (e.g., ‘f1’) or a callable that is built with scikit-learn `make_scorer`. Note that when using custom scorers, each scorer should return a single value. It defaults to ‘f1’.
- **feature\_columns** (*{str, list}, optional*) – Either ‘all’ to use from the input train data set file all the columns or a list of column names to select specific columns. It defaults to ‘all’.
- **id\_column** (*str, optional*) – The name of the column in the input train data set file containing the sample IDs. It defaults to ‘id’.
- **label\_column** (*str, optional*) – The name of the column in the input train data set file containing the ground truth labels. It defaults to ‘label’.
- **random\_seed** (*int, optional*) – The random seed used in the random number generator. It can be used to reproduce the output. It defaults to 42.
- **verbose** (*int, optional*) – Controls the verbosity of output. The higher, the more messages. It defaults to 3.
- **n\_jobs** (*int, optional*) – Number of jobs to run in parallel. -1 means using all processors. It defaults to 1.

**run** (*output\_directory*)

Run the task.

The selected features are returned as a result and also stored in a *features\_selected.txt* text file under the output directory path.

**Parameters** **output\_directory** (*str*) – The output directory path under which task-related generated files are stored.

**Returns** The task’s result. Specifically, the selected features, which can be either column names from the input train data set or column indexes from the preprocessed data set, depending on whether a preprocessor was used or not.

**Return type** list

**custom\_folds** (*folds\_file\_path, fold\_column='fold'*)

Optional method.

Use cross-validation with user-defined custom folds.

**Parameters**

- **folds\_file\_path** (*str*) – The path to the file containing the user-defined folds for the samples. The file needs to be formatted with the same separation delimiter (comma for CSV, tab for TSV, etc.) as the one used in the input data set files provided to the task. The

file must contain two data columns and the first row must be the header. The first column is for the sample IDs and needs to be the same as the one used in the input data set files provided to the task. The second column is for the fold IDs (e.g., 1 through 5, A through D, etc.).

- **fold\_column** (*str*, *optional*) – The column name for the fold IDs. It defaults to 'fold'.

**Returns** The object instance itself.

**Return type** *BaseCrossValidationTask*

**get\_configuration()**

Get the task's parameters.

**Returns** The task's parameters as a dictionary of key-value pairs, where the key is the parameter name and the value is the parameter value.

**Return type** dict

**get\_type()**

Get the task's type name.

**Returns** The task's type name.

**Return type** str

**stratified\_folds** (*total\_folds=3*, *shuffle=False*)

Optional method.

Use stratified k-fold cross-validation.

The folds are made by preserving the percentage of samples for each class.

**Parameters**

- **total\_folds** (*int*, *optional*) – Number of folds. Must be at least 2. It defaults to 3.
- **shuffle** (*bool*, *optional*) – Whether to shuffle each class's samples before splitting into batches. Note that the samples within each split will not be shuffled. It defaults to False.

**Returns** The object instance itself.

**Return type** *BaseCrossValidationTask*

## skrobot.tasks.hyperparameters\_search\_cross\_validation\_task module

**class** skrobot.tasks.hyperparameters\_search\_cross\_validation\_task.**HyperParametersSearchCross**

Bases: *skrobot.tasks.base\_cross\_validation\_task.BaseCrossValidationTask*

The *HyperParametersSearchCrossValidationTask* class can be used to search the best hyperparameters of a scikit-learn estimator/pipeline on some data.

### Cross-Validation

It can support both stratified k-fold cross-validation as well as cross-validation with user-defined folds.

By default, stratified k-fold cross-validation is used with the default parameters of *stratified\_folds()* method.

### Search

It can support both grid search as well as random search.

By default, grid search is used.

```
__init__(estimator, search_params, train_data_set_file_path, estimator_params=None,
         field_delimiter=',', scorers=['roc_auc', 'average_precision', 'f1', 'precision', 'recall', 'accuracy'],
         feature_columns='all', id_column='id', label_column='label', objective_score='f1',
         random_seed=42, verbose=3, n_jobs=1, return_train_score=True)
```

This is the constructor method and can be used to create a new object instance of `HyperParametersSearchCrossValidationTask` class.

### Parameters

- **estimator** (*scikit-learn {estimator, pipeline}*) – It can be either an estimator (e.g., `LogisticRegression`) or a pipeline ending with an estimator.
- **search\_params** (*{dict, list of dictionaries}*) – Dictionary with hyperparameters names as keys and lists of hyperparameter settings to try as values, or a list of such dictionaries, in which case the grids spanned by each dictionary in the list are explored. This enables searching over any sequence of hyperparameter settings.
- **train\_data\_set\_file\_path** (*str*) – The file path of the input train data set. It can be either a URL or a disk file path.
- **estimator\_params** (*dict, optional*) – The parameters to override in the provided estimator/pipeline. It defaults to `None`.
- **field\_delimiter** (*str, optional*) – The separation delimiter (comma for CSV, tab for TSV, etc.) used in the input train data set file. It defaults to `','`.
- **scorers** (*{list, dict}, optional*) – Multiple metrics to evaluate the predictions on the hold-out data. Either give a list of (unique) strings or a dict with names as keys and callables as values. The callables should be scorers built using `scikit-learn make_scorer`. Note that when using custom scorers, each scorer should return a single value. It defaults to `['roc_auc', 'average_precision', 'f1', 'precision', 'recall', 'accuracy']`.
- **feature\_columns** (*{str, list}, optional*) – Either `'all'` to use from the input train data set file all the columns or a list of column names to select specific columns. It defaults to `'all'`.
- **id\_column** (*str, optional*) – The name of the column in the input train data set file containing the sample IDs. It defaults to `'id'`.
- **label\_column** (*str, optional*) – The name of the column in the input train data set file containing the ground truth labels. It defaults to `'label'`.
- **objective\_score** (*str, optional*) – The scorer that would be used to find the best hyperparameters for refitting the best estimator/pipeline at the end. It defaults to `'f1'`.
- **random\_seed** (*int, optional*) – The random seed used in the random number generator. It can be used to reproduce the output. It defaults to 42.
- **verbose** (*int, optional*) – Controls the verbosity of output. The higher, the more messages. It defaults to 3.
- **n\_jobs** (*int, optional*) – Number of jobs to run in parallel. -1 means using all processors. It defaults to 1.
- **return\_train\_score** (*bool, optional*) – If `False`, training scores will not be computed and returned. Computing training scores is used to get insights on how different parameter settings impact the overfitting/underfitting trade-off. It defaults to `True`.

**grid\_search()**

Optional method.

Use the grid search method when searching the best hyperparameters.

**Returns** The object instance itself.

**Return type** *HyperParametersSearchCrossValidationTask*

**random\_search** (*n\_iters=200*)

Optional method.

Use the random search method when searching the best hyperparameters.

**Parameters** **n\_iters** (*int, optional*) – Number of hyperparameter settings that are sampled. *n\_iters* trades off runtime vs quality of the solution. It defaults to 200.

**Returns** The object instance itself.

**Return type** *HyperParametersSearchCrossValidationTask*

**run** (*output\_directory*)

Run the task.

The search results (*search\_results*) are stored also in a *search\_results.html* file as a static HTML table under the output directory path.

**Parameters** **output\_directory** (*str*) – The output directory path under which task-related generated files are stored.

**Returns** The task's result. Specifically, **1**) *best\_estimator*: The estimator/pipeline that was chosen by the search, i.e. estimator/pipeline which gave best score on the hold-out data. **2**) *best\_params*: The hyperparameters setting that gave the best results on the hold-out data. **3**) *best\_score*: Mean cross-validated score of the *best\_estimator*. **4**) *search\_results*: Metrics measured for each of the hyperparameters setting in the search. **5**) *best\_index*: The index (of the *search\_results*) which corresponds to the best candidate hyperparameters setting.

**Return type** dict

**custom\_folds** (*folds\_file\_path, fold\_column='fold'*)

Optional method.

Use cross-validation with user-defined custom folds.

**Parameters**

- **folds\_file\_path** (*str*) – The path to the file containing the user-defined folds for the samples. The file needs to be formatted with the same separation delimiter (comma for CSV, tab for TSV, etc.) as the one used in the input data set files provided to the task. The file must contain two data columns and the first row must be the header. The first column is for the sample IDs and needs to be the same as the one used in the input data set files provided to the task. The second column is for the fold IDs (e.g., 1 through 5, A through D, etc.).
- **fold\_column** (*str, optional*) – The column name for the fold IDs. It defaults to 'fold'.

**Returns** The object instance itself.

**Return type** *BaseCrossValidationTask*

**get\_configuration** ()

Get the task's parameters.

**Returns** The task's parameters as a dictionary of key-value pairs, where the key is the parameter name and the value is the parameter value.

**Return type** dict

**get\_type()**

Get the task's type name.

**Returns** The task's type name.

**Return type** str

**stratified\_folds** (*total\_folds=3, shuffle=False*)

Optional method.

Use stratified k-fold cross-validation.

The folds are made by preserving the percentage of samples for each class.

**Parameters**

- **total\_folds** (*int, optional*) – Number of folds. Must be at least 2. It defaults to 3.
- **shuffle** (*bool, optional*) – Whether to shuffle each class's samples before splitting into batches. Note that the samples within each split will not be shuffled. It defaults to False.

**Returns** The object instance itself.

**Return type** *BaseCrossValidationTask*

## skrobot.tasks.prediction\_task module

```
class skrobot.tasks.prediction_task.PredictionTask(estimator, data_set_file_path,  
                                                    field_delimiter=',', feature  
                                                    columns='all',  
                                                    id_column='id', predic  
                                                    tion_column='prediction', thresh  
                                                    old=0.5)
```

Bases: *skrobot.tasks.base\_task.BaseTask*

The *PredictionTask* class can be used to predict new data using a scikit-learn estimator/pipeline.

```
__init__(estimator, data_set_file_path, field_delimiter=',', feature_columns='all', id_column='id',  
          prediction_column='prediction', threshold=0.5)
```

This is the constructor method and can be used to create a new object instance of *PredictionTask* class.

**Parameters**

- **estimator** (*scikit-learn {estimator, pipeline}*) – It can be either an estimator (e.g., LogisticRegression) or a pipeline ending with an estimator. The estimator needs to be able to predict probabilities through a *predict\_proba* method.
- **data\_set\_file\_path** (*str*) – The file path of the input data set. It can be either a URL or a disk file path.
- **field\_delimiter** (*str, optional*) – The separation delimiter (comma for CSV, tab for TSV, etc.) used in the input data set file. It defaults to ','.
- **feature\_columns** (*{str, list}, optional*) – Either 'all' to use from the input data set file all the columns or a list of column names to select specific columns. It defaults to 'all'.
- **id\_column** (*str, optional*) – The name of the column in the input data set file containing the sample IDs. It defaults to 'id'.

- **prediction\_column** (*str*, *optional*) – The name of the column for the predicted binary class labels. It defaults to ‘prediction’.
- **threshold** (*float*, *optional*) – The threshold to use for converting the predicted probability into a binary class label. It defaults to 0.5.

**run** (*output\_directory*)

Run the task.

The predictions are returned as a result and also stored in a *predictions.csv* CSV file under the output directory path.

**Parameters** **output\_directory** (*str*) – The output directory path under which task-related generated files are stored.

**Returns** The task’s result. Specifically, the predictions for the input data set, containing the sample IDs, the predicted binary class labels, and the predicted probabilities for the positive class.

**Return type** pandas DataFrame

**get\_configuration** ()

Get the task’s parameters.

**Returns** The task’s parameters as a dictionary of key-value pairs, where the key is the parameter name and the value is the parameter value.

**Return type** dict

**get\_type** ()

Get the task’s type name.

**Returns** The task’s type name.

**Return type** str

## skrobot.tasks.train\_task module

```
class skrobot.tasks.train_task.TrainTask(estimator, train_data_set_file_path, estimator_params=None, field_delimiter=',', feature_columns='all', id_column='id', label_column='label', random_seed=42)
```

Bases: *skrobot.tasks.base\_task.BaseTask*

The *TrainTask* class can be used to fit a scikit-learn estimator/pipeline on train data.

```
__init__(estimator, train_data_set_file_path, estimator_params=None, field_delimiter=',', feature_columns='all', id_column='id', label_column='label', random_seed=42)
```

This is the constructor method and can be used to create a new object instance of *TrainTask* class.

### Parameters

- **estimator** (*scikit-learn {estimator, pipeline}*) – It can be either an estimator (e.g., LogisticRegression) or a pipeline ending with an estimator.
- **train\_data\_set\_file\_path** (*str*) – The file path of the input train data set. It can be either a URL or a disk file path.
- **estimator\_params** (*dict*, *optional*) – The parameters to override in the provided estimator/pipeline. It defaults to None.
- **field\_delimiter** (*str*, *optional*) – The separation delimiter (comma for CSV, tab for TSV, etc.) used in the input train data set file. It defaults to ‘,’.



- **feature\_columns** (*{str, list}, optional*) – Either ‘all’ to use from the input train data set file all the columns or a list of column names to select specific columns. It defaults to ‘all’.
- **id\_column** (*str, optional*) – The name of the column in the input train data set file containing the sample IDs. It defaults to ‘id’.
- **label\_column** (*str, optional*) – The name of the column in the input train data set file containing the ground truth labels. It defaults to ‘label’.
- **random\_seed** (*int, optional*) – The random seed used in the random number generator. It can be used to reproduce the output. It defaults to 42.

**run** (*output\_directory*)

Run the task.

The fitted estimator/pipeline is returned as a result and also stored in a *trained\_model.pkl* pickle file under the output directory path.

**Parameters** **output\_directory** (*str*) – The output directory path under which task-related generated files are stored.

**Returns** The task’s result. Specifically, the fitted estimator/pipeline.

**Return type** dict

**get\_configuration** ()

Get the task’s parameters.

**Returns** The task’s parameters as a dictionary of key-value pairs, where the key is the parameter name and the value is the parameter value.

**Return type** dict

**get\_type** ()

Get the task’s type name.

**Returns** The task’s type name.

**Return type** str



## WHAT IS IT ABOUT?

**skrobot** is a Python module for designing, running and tracking Machine Learning experiments / tasks. It is built on top of [scikit-learn](#) framework.



## WHY DOES IT EXISTS?

It can help Data Scientists and Machine Learning Engineers:

- to keep track of modelling experiments / tasks
- to automate the repetitive (and boring) stuff when designing modelling pipelines
- to spend more time on the things that truly matter when solving a problem



## HOW DO I INSTALL IT?

```
$ pip install skrobot
```





## WHICH ARE THE COMPONENTS?

**NOTE :** Currently, skrobot can be used only for binary classification problems.

### 5.1 For the module's users

Component	What is this?
Train Task	This task can be used to fit a scikit-learn estimator on some data.
Prediction Task	This task can be used to predict new data using a scikit-learn estimator.
Evaluation Cross Validation Task	This task can be used to evaluate a scikit-learn estimator on some data.
Feature Selection Cross Validation Task	This task can be used to perform feature selection with Recursive Feature Elimination using a scikit-learn estimator on some data.
Hyperparameters Search Cross Validation Task	This task can be used to search the best hyperparameters of a scikit-learn estimator on some data.
Experiment	This is used to build, track and run an experiment. It can run tasks in the context of an experiment.
Task Runner	This is a simplified version (in functionality) of the Experiment component. It leaves out all the “experiment” stuff and is focused mostly in the execution and tracking of tasks.

### 5.2 For the module's developers

Component	What is this?
Base Task	All tasks inherit from this component. A task is a configurable and reproducible piece of code built on top of scikit-learn that can be used in machine learning pipelines.
Base Cross Validation Task	All tasks that use cross validation functionality inherit from this component.
Base Notifier	All notifiers inherit from this component. A notifier can be used to send success / failure notifications for tasks execution.



## HOW DO I USE IT?

The following examples use many of skrobot's components to build a machine learning modelling pipeline. Please try them and we would love to have your feedback! Furthermore, many examples can be found in the project's [repository](#).

### 6.1 Example on Titanic Dataset

The following example has generated the following [results](#).

```
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.linear_model import LogisticRegression

from skrobot.core import Experiment
from skrobot.tasks import TrainTask
from skrobot.tasks import PredictionTask
from skrobot.tasks import FeatureSelectionCrossValidationTask
from skrobot.tasks import EvaluationCrossValidationTask
from skrobot.tasks import HyperParametersSearchCrossValidationTask
from skrobot.feature_selection import ColumnSelector
from skrobot.notification import BaseNotifier

##### Initialization Code

train_data_set_file_path = 'https://bit.ly/titanic-data-train'
test_data_set_file_path = 'https://bit.ly/titanic-data-test'
new_data_set_file_path = 'https://bit.ly/titanic-data-new'

random_seed = 42

id_column = 'PassengerId'

label_column = 'Survived'

numerical_features = ['Age', 'Fare', 'SibSp', 'Parch']

categorical_features = ['Embarked', 'Sex', 'Pclass']

numeric_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer()),
```

(continues on next page)

(continued from previous page)

```

        ('scaler', StandardScaler()))

categorical_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='most_frequent')),
    ('encoder', OneHotEncoder(handle_unknown='ignore'))])

preprocessor = ColumnTransformer(transformers=[
    ('numerical_transformer', numeric_transformer, numerical_features),
    ('categorical_transformer', categorical_transformer, categorical_features)])

classifier = LogisticRegression(solver='liblinear', random_state=random_seed)

search_params = {
    "classifier__C" : [ 1.e-01, 1.e+00, 1.e+01 ],
    "classifier__penalty" : [ "l1", "l2" ],
    "preprocessor__numerical_transformer__imputer__strategy" : [ "mean", "median" ]
}

##### skrobot Code

# Define a Notifier (This is optional and you can implement any notifier you want, e.
# g. for Slack / Trello / Discord)
class ConsoleNotifier(BaseNotifier):
    def notify (self, message):
        print(message)

# Build an Experiment
experiment = Experiment('experiments-output').set_source_code_file_path(__file__).set_
    ↪ experimenter('echatzikyriakidis').set_notifier(ConsoleNotifier()).build()

# Run Feature Selection Task
features_columns = experiment.run(FeatureSelectionCrossValidationTask,
    ↪ (estimator=classifier,
                                     train_data_set_
    ↪ file_path=train_data_set_file_path,
                                     ↪
    ↪ preprocessor=preprocessor,
                                     min_features_
    ↪ to_select=4,
                                     id_column=id_
    ↪ column,
                                     label_
    ↪ column=label_column,
                                     random_
    ↪ seed=random_seed).stratified_folds(total_folds=5, shuffle=True))

pipe = Pipeline(steps=[('preprocessor', preprocessor),
    ('selector', ColumnSelector(cols=features_columns)),
    ('classifier', classifier)])

# Run Hyperparameters Search Task
hyperparameters_search_results = experiment.
    ↪ run(HyperParametersSearchCrossValidationTask (estimator=pipe,
    ↪
    ↪ search_params=search_params,
    ↪
    ↪ train_data_set_file_path=train_data_set_file_path,

```

(continues on next page)

(continued from previous page)

```

↳ id_column=id_column,
↳ label_column=label_column,
↳ random_seed=random_seed).random_search(n_iters=100).stratified_folds(total_
↳ folds=5, shuffle=True))

# Run Evaluation Task
evaluation_results = experiment.run(EvaluationCrossValidationTask(estimator=pipe,
                                                                    estimator_
↳ params=hyperparameters_search_results['best_params'],
                                                                    train_data_set_file_
↳ path=train_data_set_file_path,
                                                                    test_data_set_file_
↳ path=test_data_set_file_path,
                                                                    id_column=id_column,
                                                                    label_column=label_
↳ column,
                                                                    random_seed=random_
↳ seed,
                                                                    export_
↳ classification_reports=True,
                                                                    export_confusion_
↳ matrixes=True,
                                                                    export_pr_
↳ curves=True,
                                                                    export_roc_
↳ curves=True,
                                                                    export_false_
↳ positives_reports=True,
                                                                    export_false_
↳ negatives_reports=True,
                                                                    export_also_for_
↳ train_folds=True).stratified_folds(total_folds=5, shuffle=True))

# Run Train Task
train_results = experiment.run(TrainTask(estimator=pipe,
                                                                    estimator_params=hyperparameters_search_
↳ results['best_params'],
                                                                    train_data_set_file_path=train_data_set_file_
↳ path,
                                                                    id_column=id_column,
                                                                    label_column=label_column,
                                                                    random_seed=random_seed))

# Run Prediction Task
predictions = experiment.run(PredictionTask(estimator=train_results['estimator'],
                                                                    data_set_file_path=new_data_set_file_path,
                                                                    id_column=id_column,
                                                                    prediction_column=label_column,
                                                                    threshold=evaluation_results['threshold
↳ ']))

# Print in-memory results
print(features_columns)

```

(continues on next page)

(continued from previous page)

```

print(hyperparameters_search_results['best_params'])
print(hyperparameters_search_results['best_index'])
print(hyperparameters_search_results['best_estimator'])
print(hyperparameters_search_results['best_score'])
print(hyperparameters_search_results['search_results'])

print(evaluation_results['threshold'])
print(evaluation_results['cv_threshold_metrics'])
print(evaluation_results['cv_splits_threshold_metrics'])
print(evaluation_results['cv_splits_threshold_metrics_summary'])
print(evaluation_results['test_threshold_metrics'])

print(train_results['estimator'])

print(predictions)

```

## 6.2 Example on SMS Spam Collection Dataset

The following example has generated the following results.

```

from sklearn.pipeline import Pipeline
from sklearn.feature_extraction.text import CountVectorizer, TfidfTransformer
from sklearn.feature_selection import SelectPercentile, chi2
from sklearn.linear_model import SGDClassifier

from skrobot.core import Experiment
from skrobot.tasks import TrainTask
from skrobot.tasks import PredictionTask
from skrobot.tasks import EvaluationCrossValidationTask
from skrobot.tasks import HyperParametersSearchCrossValidationTask
from skrobot.feature_selection import ColumnSelector

##### Initialization Code

train_data_set_file_path = 'https://bit.ly/sms-spam-ham-data-train'

test_data_set_file_path = 'https://bit.ly/sms-spam-ham-data-test'

new_data_set_file_path = 'https://bit.ly/sms-spam-ham-data-new'

field_delimiter = '\t'

random_seed = 42

pipe = Pipeline(steps=[
    ('column_selection', ColumnSelector(cols=['message'], drop_axis=True)),
    ('vectorizer', CountVectorizer()),
    ('tfidf', TfidfTransformer()),
    ('feature_selection', SelectPercentile(chi2)),
    ('classifier', SGDClassifier(loss='log'))])

search_params = {
    'classifier__max_iter': [ 20, 50, 80 ],
    'classifier__alpha': [ 0.00001, 0.000001 ],

```

(continues on next page)

(continued from previous page)

```

    'classifier_penalty': [ 'l2', 'elasticnet' ],
    "vectorizer_stop_words" : [ "english", None ],
    "vectorizer_ngram_range" : [ (1, 1), (1, 2) ],
    "vectorizer_max_df": [ 0.5, 0.75, 1.0 ],
    "tfidf_use_idf" : [ True, False ],
    "tfidf_norm" : [ 'l1', 'l2' ],
    "feature_selection_percentile" : [ 70, 60, 50 ]
}

##### skrobot Code

# Build an Experiment
experiment = Experiment('experiments-output').set_source_code_file_path(__file__).set_
↳ experimenter('echatzikyriakidis').build()

# Run Hyperparameters Search Task
hyperparameters_search_results = experiment.
↳ run(HyperParametersSearchCrossValidationTask (estimator=pipe,

↳ search_params=search_params,

↳ train_data_set_file_path=train_data_set_file_path,

↳ field_delimiter=field_delimiter,

↳ random_seed=random_seed).random_search().stratified_folds(total_folds=5,
↳ shuffle=True))

# Run Evaluation Task
evaluation_results = experiment.run(EvaluationCrossValidationTask(estimator=pipe,
                                                                    estimator_
↳ params=hyperparameters_search_results['best_params'],
                                                                    train_data_set_file_
↳ path=train_data_set_file_path,
                                                                    test_data_set_file_
↳ path=test_data_set_file_path,
                                                                    field_
↳ delimiter=field_delimiter,
                                                                    random_seed=random_
↳ seed,
                                                                    export_
↳ classification_reports=True,
                                                                    export_confusion_
↳ matrixes=True,
                                                                    export_pr_
↳ curves=True,
                                                                    export_roc_
↳ curves=True,
                                                                    export_false_
↳ positives_reports=True,
                                                                    export_false_
↳ negatives_reports=True,
                                                                    export_also_for_
↳ train_folds=True).stratified_folds(total_folds=5, shuffle=True))

# Run Train Task
train_results = experiment.run(TrainTask(estimator=pipe,

```

(continues on next page)

(continued from previous page)

```
↪results['best_params'],
↪path,
                                estimator_params=hyperparameters_search_
                                train_data_set_file_path=train_data_set_file_
                                field_delimiter=field_delimiter,
                                random_seed=random_seed))

# Run Prediction Task
predictions = experiment.run(PredictionTask(estimator=train_results['estimator'],
                                data_set_file_path=new_data_set_file_path,
                                field_delimiter=field_delimiter,
                                threshold=evaluation_results['threshold
↪']))

# Print in-memory results
print(hyperparameters_search_results['best_params'])
print(hyperparameters_search_results['best_index'])
print(hyperparameters_search_results['best_estimator'])
print(hyperparameters_search_results['best_score'])
print(hyperparameters_search_results['search_results'])

print(evaluation_results['threshold'])
print(evaluation_results['cv_threshold_metrics'])
print(evaluation_results['cv_splits_threshold_metrics'])
print(evaluation_results['cv_splits_threshold_metrics_summary'])
print(evaluation_results['test_threshold_metrics'])

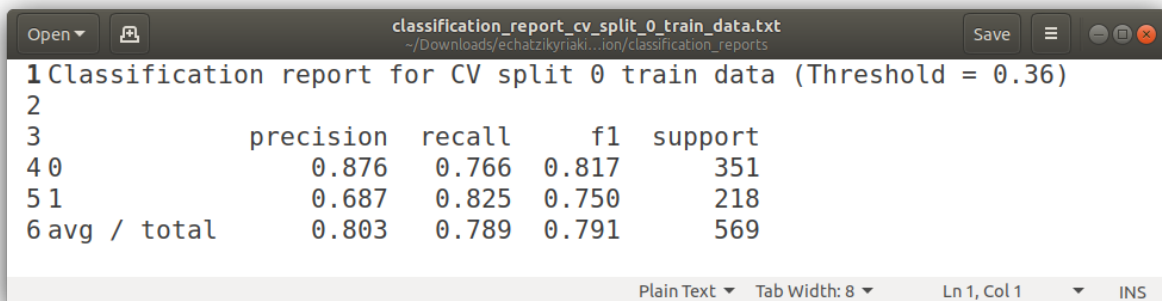
print(train_results['estimator'])

print(predictions)
```



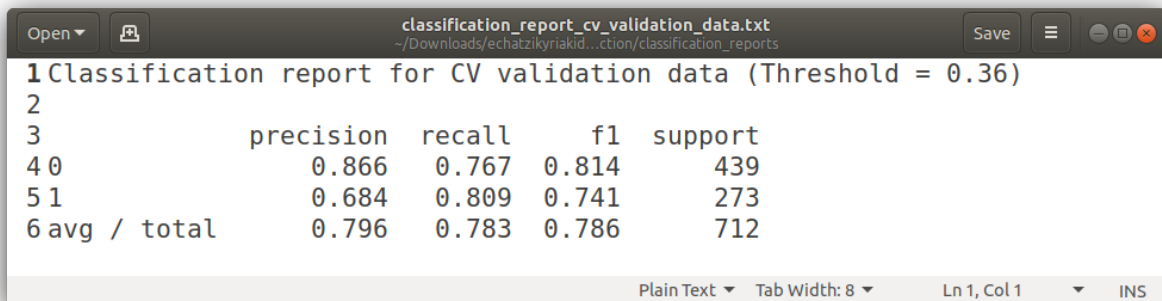
## SAMPLE OF GENERATED RESULTS?

### 7.1 Classification Reports



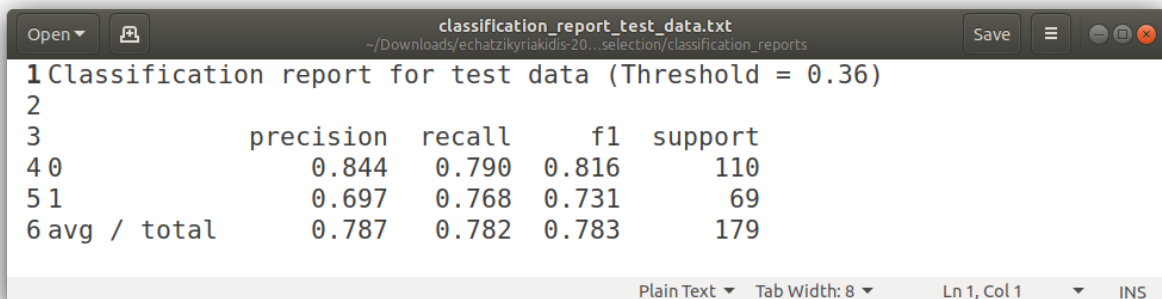
A screenshot of a text editor window titled "classification\_report\_cv\_split\_0\_train\_data.txt". The window shows a classification report for training data with a threshold of 0.36. The report includes precision, recall, f1 score, and support for classes 0 and 1, along with an average for the total. The editor interface includes an "Open" button, a file icon, a "Save" button, and window control buttons. The status bar at the bottom indicates "Plain Text", "Tab Width: 8", "Ln 1, Col 1", and "INS".

```
1 Classification report for CV split 0 train data (Threshold = 0.36)
2
3           precision    recall  f1   support
4 0           0.876      0.766  0.817       351
5 1           0.687      0.825  0.750       218
6 avg / total    0.803      0.789  0.791       569
```



A screenshot of a text editor window titled "classification\_report\_cv\_validation\_data.txt". The window shows a classification report for validation data with a threshold of 0.36. The report includes precision, recall, f1 score, and support for classes 0 and 1, along with an average for the total. The editor interface includes an "Open" button, a file icon, a "Save" button, and window control buttons. The status bar at the bottom indicates "Plain Text", "Tab Width: 8", "Ln 1, Col 1", and "INS".

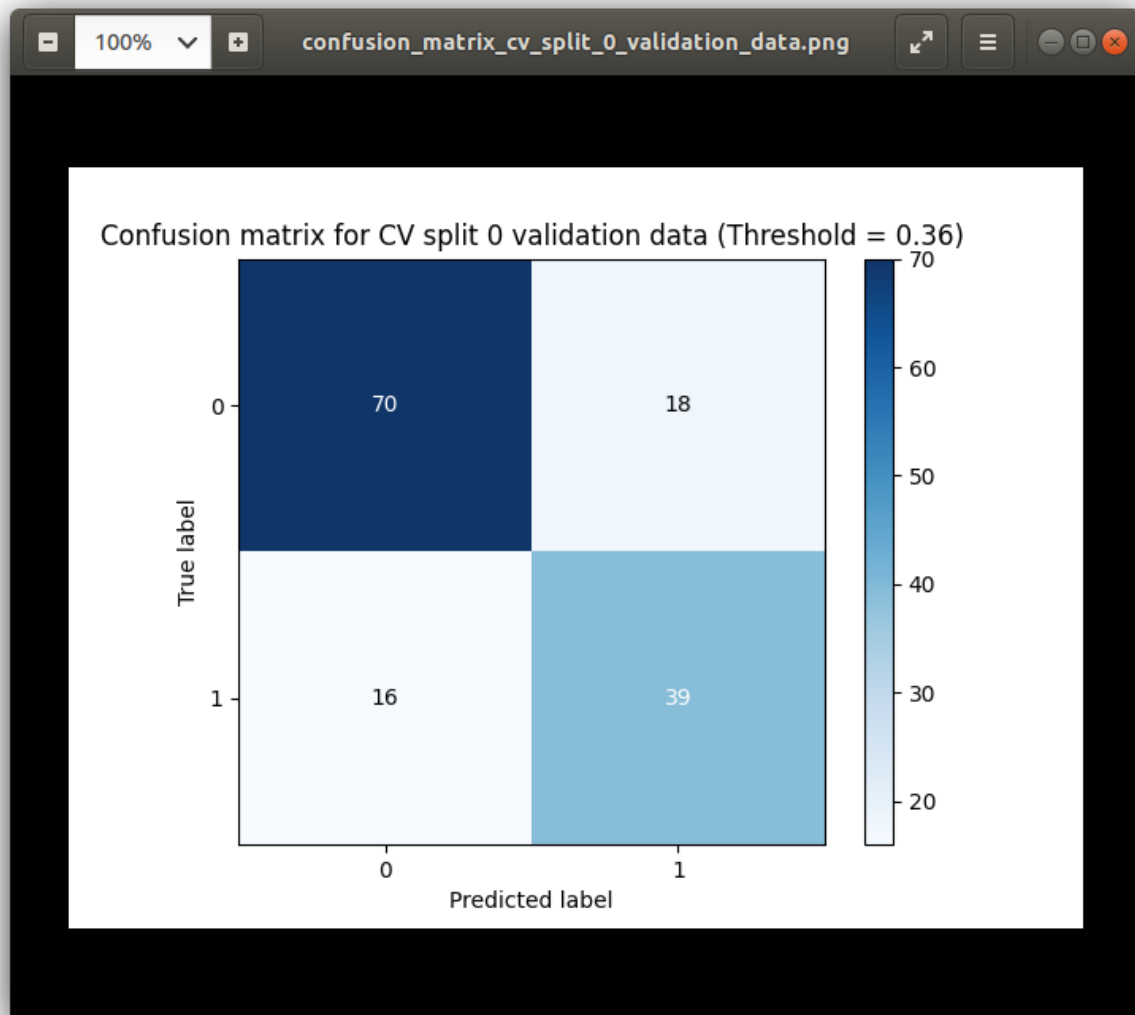
```
1 Classification report for CV validation data (Threshold = 0.36)
2
3           precision    recall  f1   support
4 0           0.866      0.767  0.814       439
5 1           0.684      0.809  0.741       273
6 avg / total    0.796      0.783  0.786       712
```

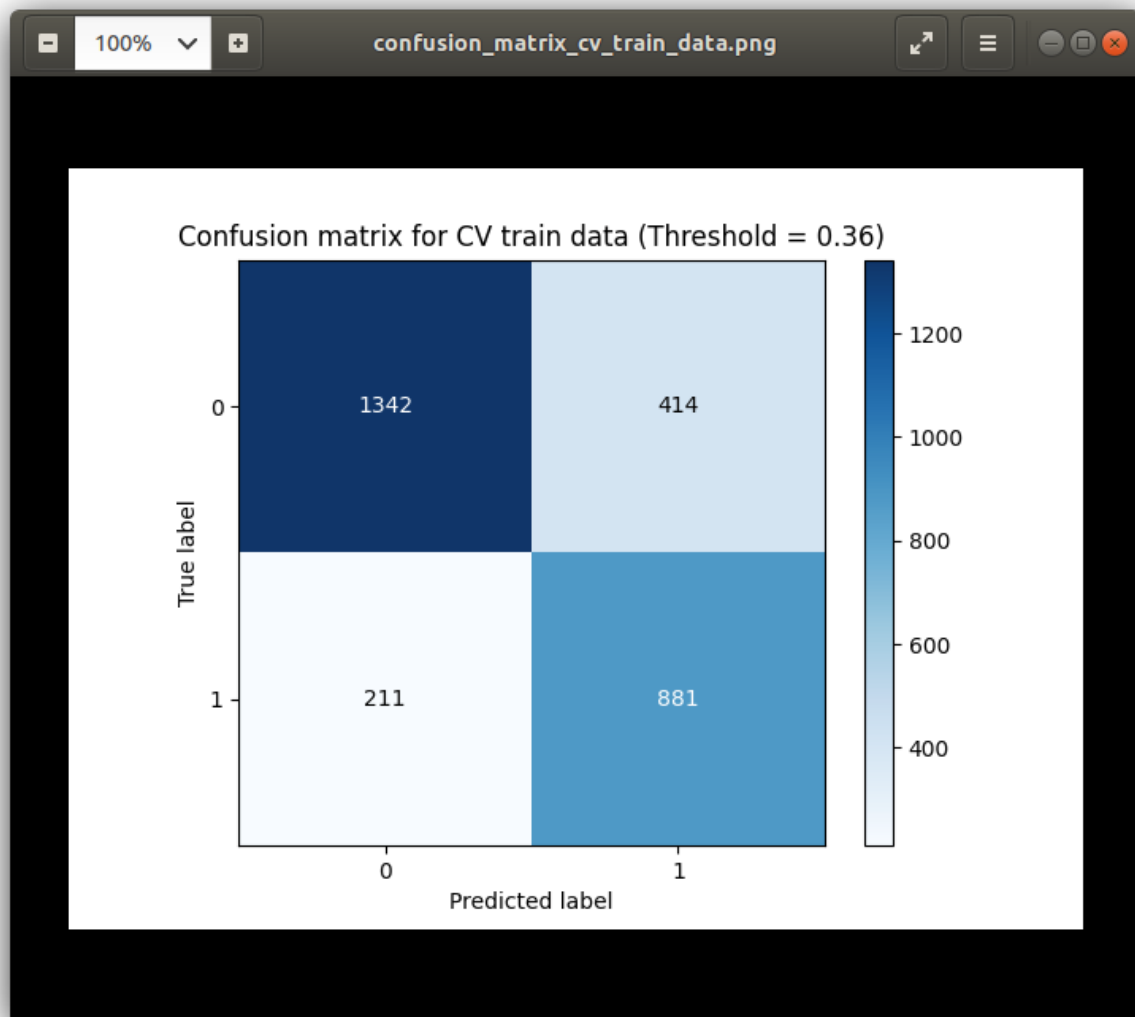


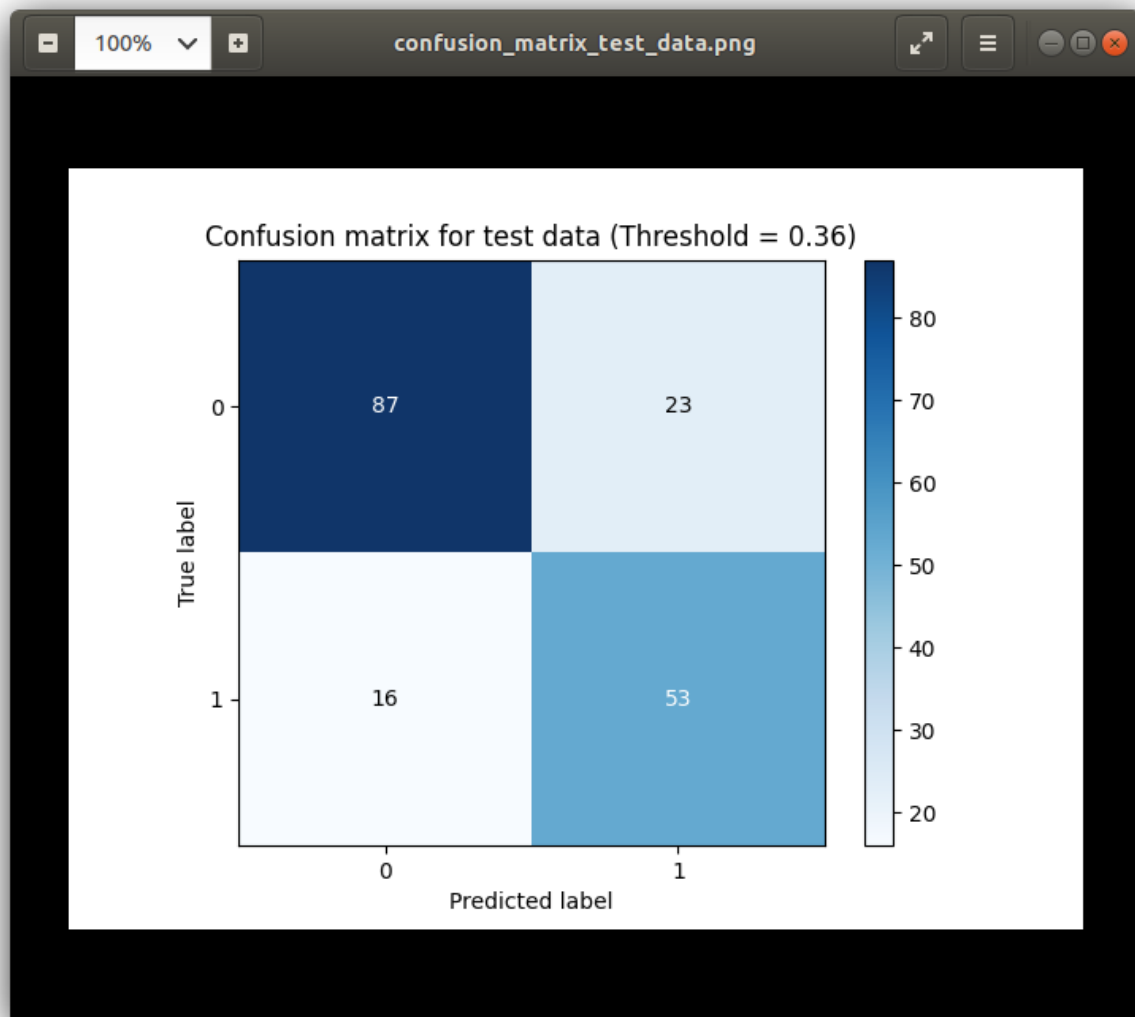
A screenshot of a text editor window titled "classification\_report\_test\_data.txt". The window shows a classification report for test data with a threshold of 0.36. The report includes precision, recall, f1 score, and support for classes 0 and 1, along with an average for the total. The editor interface includes an "Open" button, a file icon, a "Save" button, and window control buttons. The status bar at the bottom indicates "Plain Text", "Tab Width: 8", "Ln 1, Col 1", and "INS".

```
1 Classification report for test data (Threshold = 0.36)
2
3           precision    recall  f1   support
4 0           0.844      0.790  0.816       110
5 1           0.697      0.768  0.731        69
6 avg / total    0.787      0.782  0.783       179
```

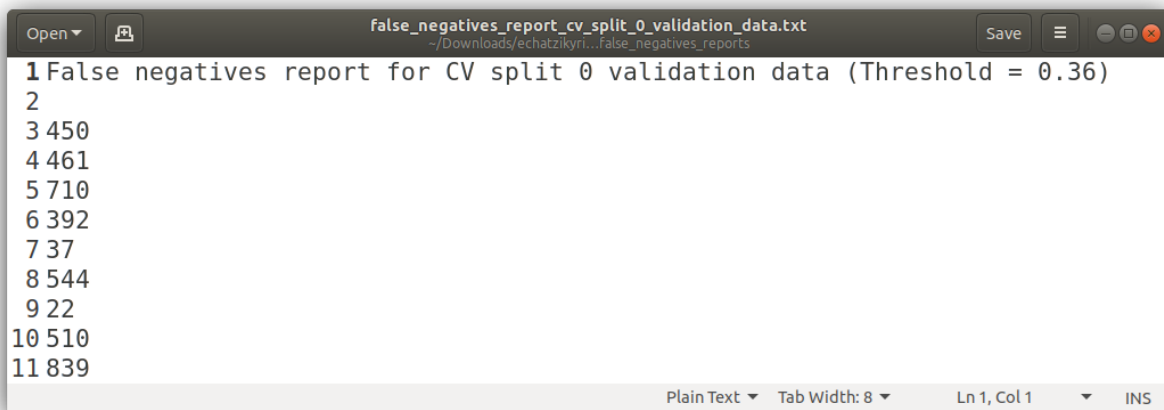
## 7.2 Confusion Matrixes







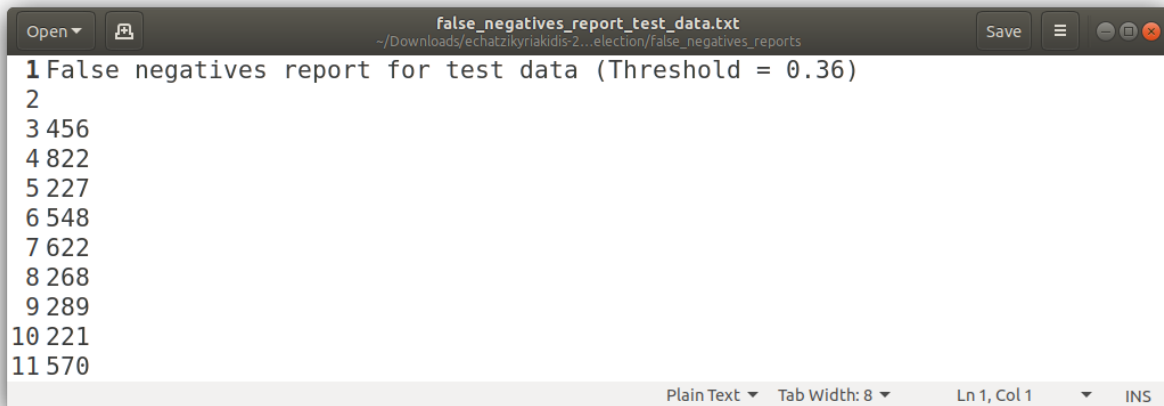
## 7.3 False Negatives



A screenshot of a text editor window titled "false\_negatives\_report\_cv\_split\_0\_validation\_data.txt". The window shows a list of 11 items, each consisting of a line number followed by a text description. The text is as follows:

```
1 False negatives report for CV split 0 validation data (Threshold = 0.36)
2
3 450
4 461
5 710
6 392
7 37
8 544
9 22
10 510
11 839
```

The editor interface includes a menu bar with "Open", "Save", and other icons. The status bar at the bottom indicates "Plain Text", "Tab Width: 8", "Ln 1, Col 1", and "INS".

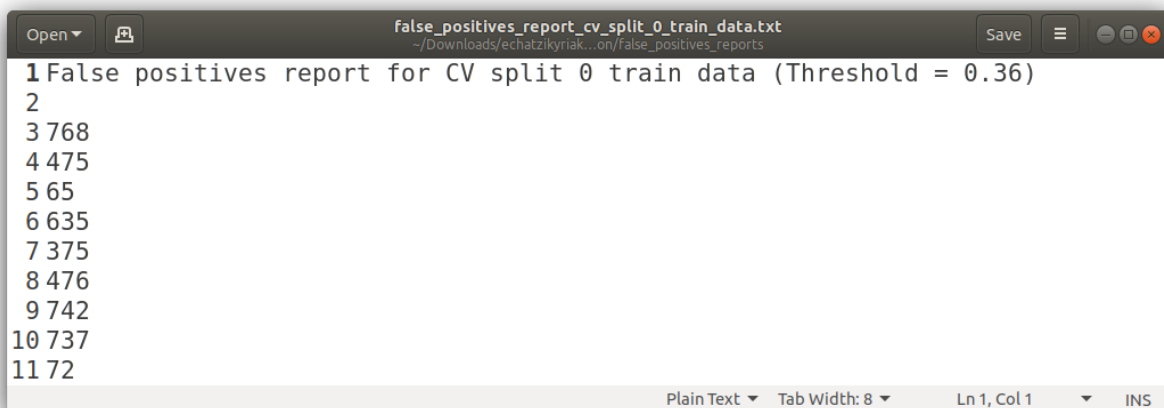


A screenshot of a text editor window titled "false\_negatives\_report\_test\_data.txt". The window shows a list of 11 items, each consisting of a line number followed by a text description. The text is as follows:

```
1 False negatives report for test data (Threshold = 0.36)
2
3 456
4 822
5 227
6 548
7 622
8 268
9 289
10 221
11 570
```

The editor interface includes a menu bar with "Open", "Save", and other icons. The status bar at the bottom indicates "Plain Text", "Tab Width: 8", "Ln 1, Col 1", and "INS".

## 7.4 False Positives



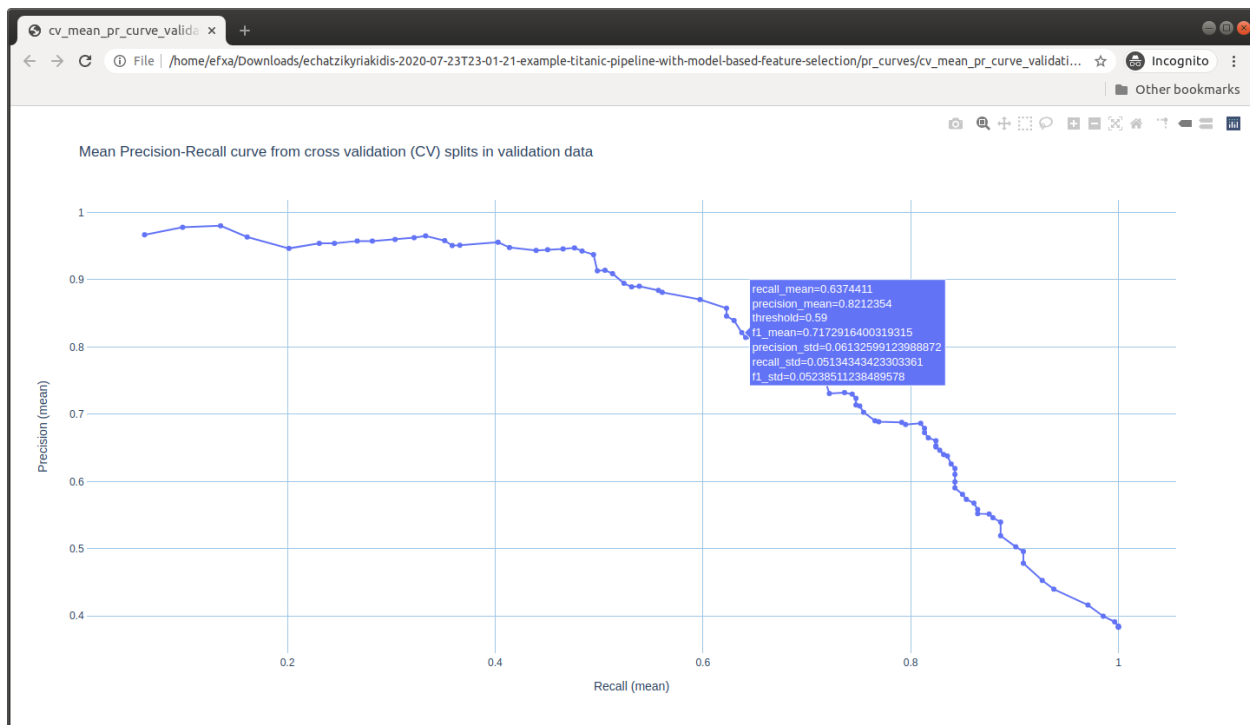
A screenshot of a text editor window titled "false\_positives\_report\_cv\_split\_0\_train\_data.txt". The window shows a list of 11 items, each consisting of a line number followed by a text description. The text is as follows:

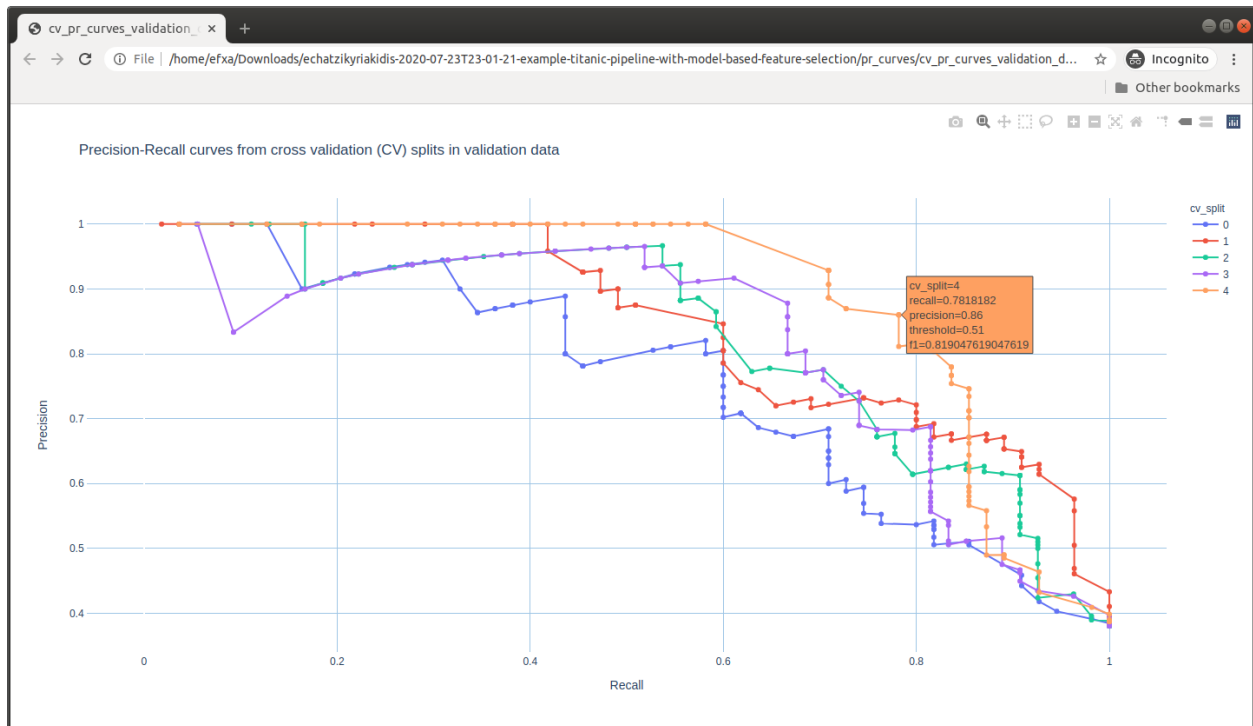
```
1 False positives report for CV split 0 train data (Threshold = 0.36)
2
3 768
4 475
5 65
6 635
7 375
8 476
9 742
10 737
11 72
```

The editor interface includes a menu bar with "Open", "Save", and other icons. The status bar at the bottom indicates "Plain Text", "Tab Width: 8", "Ln 1, Col 1", and "INS".

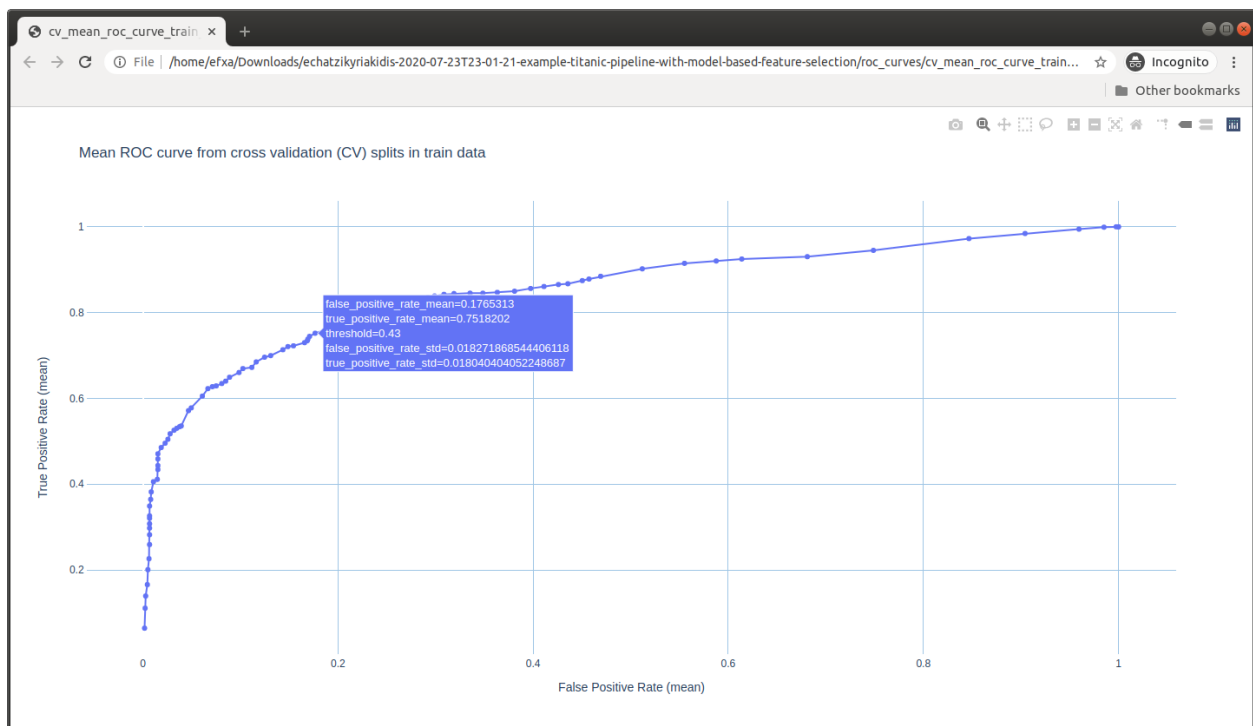
```
Open ▾  false_positives_report_test_data.txt  Save  ▮  -/Downloads/echatzikyriakidis-2...election/false_positives_reports
1 False positives report for test data (Threshold = 0.36)
2
3 358
4 246
5 333
6 528
7 594
8 542
9 410
10 618
11 453
Plain Text ▾  Tab Width: 8 ▾  Ln 1, Col 1 ▾  INS
```

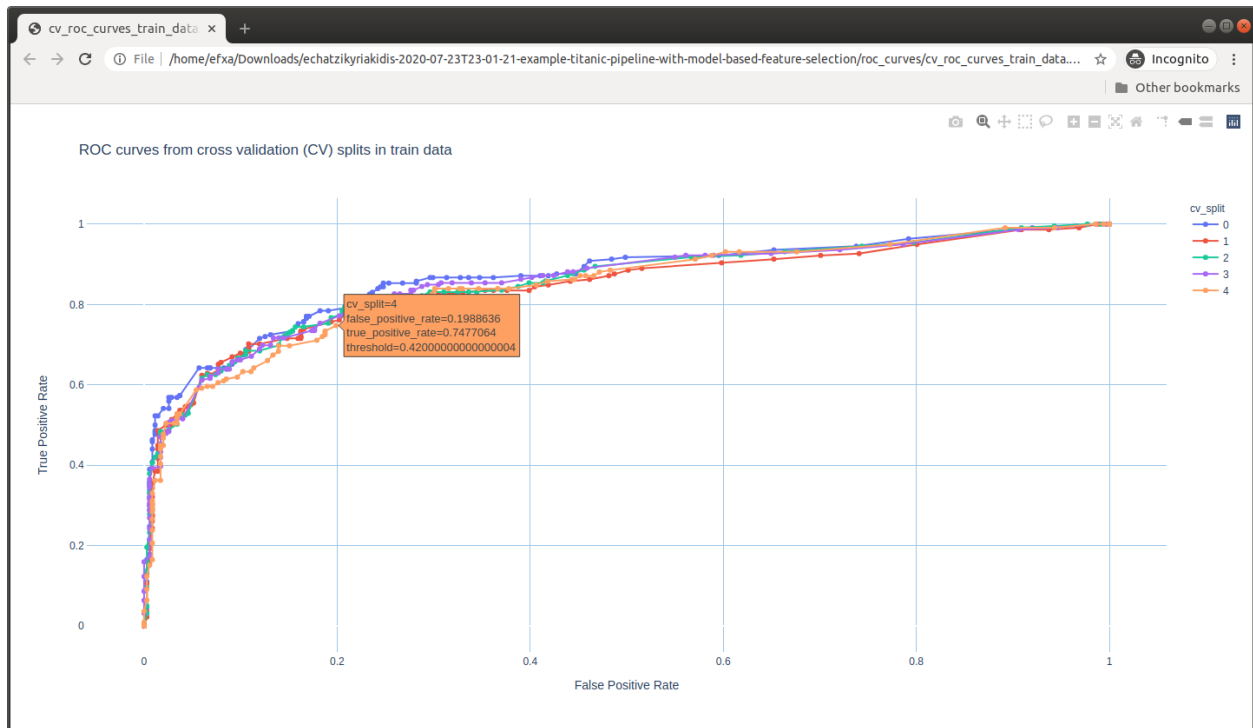
## 7.5 PR Curves





## 7.6 ROC Curves





## 7.7 Performance Metrics

On train / validation CV folds:

threshold	cv_split	validation_sensitivity	validation_specificity	validation_accuracy	validation_true_positive_rate	validation_false_positive_rate	validation_geometric_mean	validation_precision	validation
0.01	0	1.000000	0.000000	0.384615	1.000000	1.000000	0.000000	0.384615	1.000000
0.02	0	1.000000	0.000000	0.384615	1.000000	1.000000	0.000000	0.384615	1.000000
0.03	0	1.000000	0.000000	0.384615	1.000000	1.000000	0.000000	0.384615	1.000000
0.04	0	1.000000	0.000000	0.384615	1.000000	1.000000	0.000000	0.384615	1.000000
0.05	0	0.981818	0.045455	0.405594	0.981818	0.954545	0.211254	0.391304	0.981818
0.06	0	0.945455	0.125000	0.440559	0.945455	0.875000	0.343776	0.403101	0.945455
0.07	0	0.927273	0.193182	0.475524	0.927273	0.806818	0.423240	0.418033	0.927273
0.08	0	0.909091	0.284091	0.524476	0.909091	0.715909	0.508197	0.442478	0.909091
0.09	0	0.909091	0.329545	0.552448	0.909091	0.670455	0.547345	0.458716	0.909091
0.10	0	0.854545	0.477273	0.622378	0.854545	0.522727	0.638632	0.505376	0.854545
0.11	0	0.854545	0.488636	0.629371	0.854545	0.511364	0.646190	0.510870	0.854545
0.12	0	0.818182	0.500000	0.622378	0.818182	0.500000	0.639602	0.505618	0.818182
0.13	0	0.818182	0.522727	0.636364	0.818182	0.477273	0.653977	0.517241	0.818182
0.14	0	0.818182	0.545455	0.650350	0.818182	0.454545	0.668043	0.529412	0.818182
0.15	0	0.818182	0.545455	0.650350	0.818182	0.454545	0.668043	0.529412	0.818182
0.16	0	0.818182	0.556818	0.657343	0.818182	0.443182	0.674966	0.535714	0.818182
0.17	0	0.818182	0.556818	0.657343	0.818182	0.443182	0.674966	0.535714	0.818182
0.18	0	0.818182	0.568182	0.664336	0.818182	0.431818	0.681818	0.542169	0.818182
0.19	0	0.800000	0.568182	0.657343	0.800000	0.431818	0.674200	0.536585	0.800000
0.20	0	0.763636	0.590909	0.657343	0.763636	0.409091	0.671744	0.538462	0.763636
0.21	0	0.763636	0.613636	0.671329	0.763636	0.386364	0.684540	0.552632	0.763636
0.22	0	0.745455	0.625000	0.671329	0.745455	0.375000	0.682575	0.554054	0.745455
0.23	0	0.745455	0.647727	0.685315	0.745455	0.352273	0.694875	0.569444	0.745455
0.24	0	0.745455	0.681818	0.706294	0.745455	0.318182	0.712927	0.594203	0.745455
0.25	0	0.745455	0.681818	0.706294	0.745455	0.318182	0.712927	0.594203	0.745455
0.26	0	0.727273	0.681818	0.699301	0.727273	0.318182	0.704179	0.588235	0.727273
0.27	0	0.727273	0.704545	0.713287	0.727273	0.295455	0.715819	0.606061	0.727273
0.28	0	0.709091	0.704545	0.706294	0.709091	0.295455	0.706815	0.600000	0.709091

On hold-out test set:



threshold	test_sensitivity	test_specificity	test_accuracy	test_true_positive_rate	test_false_positive_rate	test_geometric_mean	test_precision	test_recall	test_f1
0.01	1.000000	0.000000	0.385475	1.000000	1.000000	0.000000	0.385475	1.000000	0.556452
0.02	1.000000	0.000000	0.385475	1.000000	1.000000	0.000000	0.385475	1.000000	0.556452
0.03	1.000000	0.009091	0.391061	1.000000	0.990909	0.095346	0.387640	1.000000	0.558704
0.04	1.000000	0.027273	0.402235	1.000000	0.972727	0.165145	0.392045	1.000000	0.563265
0.05	1.000000	0.027273	0.402235	1.000000	0.972727	0.165145	0.392045	1.000000	0.563265
0.06	1.000000	0.036364	0.407821	1.000000	0.963636	0.190693	0.394286	1.000000	0.565574
0.07	0.985507	0.090909	0.435754	0.985507	0.909091	0.299319	0.404762	0.985507	0.573840
0.08	0.956522	0.281818	0.541899	0.956522	0.718182	0.519197	0.455172	0.956522	0.616822
0.09	0.942029	0.327273	0.564246	0.942029	0.672727	0.555248	0.467626	0.942029	0.625000
0.10	0.927536	0.445455	0.631285	0.927536	0.554545	0.642787	0.512000	0.927536	0.659794
0.11	0.913043	0.472727	0.642458	0.913043	0.527273	0.656978	0.520661	0.913043	0.663158
0.12	0.913043	0.490909	0.653631	0.913043	0.509091	0.669493	0.529412	0.913043	0.670213
0.13	0.884058	0.536364	0.670391	0.884058	0.463636	0.688605	0.544643	0.884058	0.674033
0.14	0.869565	0.563636	0.681564	0.869565	0.436364	0.700085	0.555556	0.869565	0.677966
0.15	0.855072	0.572727	0.681564	0.855072	0.427273	0.699802	0.556604	0.855072	0.674286
0.16	0.840580	0.609091	0.698324	0.840580	0.390909	0.715534	0.574257	0.840580	0.682353
0.17	0.840580	0.609091	0.698324	0.840580	0.390909	0.715534	0.574257	0.840580	0.682353
0.18	0.840580	0.609091	0.698324	0.840580	0.390909	0.715534	0.574257	0.840580	0.682353
0.19	0.840580	0.618182	0.703911	0.840580	0.381818	0.720854	0.580000	0.840580	0.686391
0.20	0.840580	0.636364	0.715084	0.840580	0.363636	0.731378	0.591837	0.840580	0.694611
0.21	0.840580	0.636364	0.715084	0.840580	0.363636	0.731378	0.591837	0.840580	0.694611
0.22	0.840580	0.672727	0.737430	0.840580	0.327273	0.751985	0.617021	0.840580	0.711656
0.23	0.840580	0.718182	0.765363	0.840580	0.281818	0.776974	0.651685	0.840580	0.734177
0.24	0.840580	0.718182	0.765363	0.840580	0.281818	0.776974	0.651685	0.840580	0.734177
0.25	0.840580	0.736364	0.776536	0.840580	0.263636	0.786748	0.666667	0.840580	0.743590
0.26	0.840580	0.745455	0.782123	0.840580	0.254545	0.791590	0.674419	0.840580	0.748387
0.27	0.840580	0.745455	0.782123	0.840580	0.254545	0.791590	0.674419	0.840580	0.748387
0.28	0.826087	0.763636	0.787709	0.826087	0.236364	0.794248	0.686747	0.826087	0.750000
0.29	0.811504	0.772727	0.787709	0.811504	0.232727	0.791022	0.691358	0.811504	0.746667

## 7.8 Hyperparameters Search Results

index	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_preprocessor__numerical_transformer__imputer__strategy	param_classifier__penalty	param_classifier__C	
0	0.016693	0.000718	0.018089	0.000548	mean	l1	0.1	{'preprocessor__nume mean', 'classifier__per
1	0.017744	0.000700	0.019145	0.000597	median	l1	0.1	{'preprocessor__nume median', 'classifier__p
2	0.017219	0.000747	0.019203	0.000779	mean	l2	0.1	{'preprocessor__nume mean', 'classifier__per
3	0.017486	0.000722	0.019177	0.000649	median	l2	0.1	{'preprocessor__nume median', 'classifier__p
4	0.017819	0.000417	0.018958	0.000365	mean	l1	1	{'preprocessor__nume mean', 'classifier__per
5	0.019109	0.000578	0.019138	0.000399	median	l1	1	{'preprocessor__nume median', 'classifier__p
6	0.018570	0.001124	0.020416	0.001766	mean	l2	1	{'preprocessor__nume mean', 'classifier__per
7	0.017627	0.000371	0.019003	0.000563	median	l2	1	{'preprocessor__nume median', 'classifier__p
8	0.021885	0.000673	0.018920	0.000322	mean	l1	10	{'preprocessor__nume mean', 'classifier__per
9	0.022998	0.000909	0.019196	0.000564	median	l1	10	{'preprocessor__nume median', 'classifier__p
10	0.017580	0.001034	0.019417	0.001051	mean	l2	10	{'preprocessor__nume mean', 'classifier__per
11	0.019342	0.003049	0.019181	0.000388	median	l2	10	{'preprocessor__nume median', 'classifier__p

## 7.9 Task Parameters Logging

```

1 {
2   "type": "evaluation_cross_validation_task",
3   "estimator": "Pipeline(steps=[('preprocessor',\n                                ColumnTransformer(transformers=[('numerical_transformer',\n                                \"train_data_set_file_path\": \"https://bit.ly/titanic-data-train\",
4   \"test_data_set_file_path\": \"https://bit.ly/titanic-data-test\",
5   \"estimator_params\": {
6     \"preprocessor__numerical_transformer__imputer__strategy\": \"mean\",
7     \"classifier__penalty\": \"l1\",
8     \"classifier__C\": 10.0
9   },
10 },
11 \"field delimiter\": \",\",
12 \"feature_columns\": \"all\",
13 \"id_column\": \"PassengerId\",
14 \"label_column\": \"Survived\",
15 \"random_seed\": 42,
16 \"threshold_selection_by\": \"f1\",
17 \"metric_greater_is_better\": true,
18 \"threshold_tuning_range\": [
19   0.01,
20   1.0,
21   0.01
22 ],
23 \"export_classification_reports\": true,
24 \"export_confusion_matrixes\": true,
25 \"export_roc_curves\": true,
26 \"export_pr_curves\": true,
27 \"export_false_positives_reports\": true,
28 \"export_false_negatives_reports\": true,
29 \"export_also_for_train_folds\": true,
30 \"fscore_beta\": 1,
31 \"fold_options\": {
32   \"total_folds\": 5,
33   \"shuffle\": true
34 },
35 \"fold_method\": \"stratified\"
36 }

```

```

1 {
2   "type": "feature_selection_cross_validation_task",
3   "estimator": "LogisticRegression(random_state=42, solver='liblinear')",
4   "train_data_set_file_path": "https://bit.ly/titanic-data-train",
5   "estimator_params": null,
6   "field delimiter": \",\",
7   "preprocessor": "ColumnTransformer(transformers=[('numerical_transformer',\n                                Pipeline(steps=[('imputer', Si
8   "preprocessor_params": null,
9   "min_features_to_select": 4,
10  "scoring": "f1",
11  "feature_columns": "all",
12  "id_column": "PassengerId",
13  "label_column": "Survived",
14  "random_seed": 42,
15  "verbose": 3,
16  "n_jobs": 1,
17  "fold_options": {
18    "total_folds": 5,
19    "shuffle": true
20  },
21  "fold_method": "stratified"
22 }

```

```

hyper_parameters_search_cross_validation_task.params
~/Downloads/skrobot/...-based-feature-selection
1{
2  "type": "hyper_parameters_search_cross_validation_task",
3  "estimator": "Pipeline(steps=[('preprocessor',\n                                ColumnTransformer(transformers=[('numerical_transformer',\n0
4  "search_params": {
5    "classifier__C": [
6      0.1,
7      1.0,
8      10.0
9    ],
10   "classifier__penalty": [
11     "l1",
12     "l2"
13   ],
14   "preprocessor__numerical_transformer__imputer__strategy": [
15     "mean",
16     "median"
17   ]
18 },
19 "train_data_set_file_path": "https://bit.ly/titanic-data-train",
20 "estimator_params": null,
21 "field_delimiter": ",",
22 "scorers": [
23   "roc_auc",
24   "average_precision",
25   "f1",
26   "precision",
27   "recall",
28   "accuracy"
29 ],
30 "feature_columns": "all",
31 "id_column": "PassengerId",
32 "label_column": "Survived",
33 "objective_score": "f1",
34 "random_seed": 42,
35 "verbose": 3,
36 "n_jobs": 1,
37 "return_train_score": true,
38 "fold_options": {
train_task.params
~/Downloads/skrobot/...-based-feature-selection
1{
2  "type": "train_task",
3  "estimator": "Pipeline(steps=[('preprocessor',\n                                ColumnTransformer(transformers=[('numerical_transformer',\n0
4  "train_data_set_file_path": "https://bit.ly/titanic-data-train",
5  "estimator_params": {
6    "preprocessor__numerical_transformer__imputer__strategy": "mean",
7    "classifier__penalty": "l1",
8    "classifier__C": 10.0
9  },
10 "field_delimiter": ",",
11 "feature_columns": "all",
12 "id_column": "PassengerId",
13 "label_column": "Survived",
14 "random_seed": 42
15}

```

```
1 {
2   "type": "prediction task",
3   "estimator": "Pipeline(steps=[('preprocessor',\n                                ColumnTransformer(transformers=[('numerical_transformer',\n                                "data_set_file_path": "https://bit.ly/titanic-data-new",
4   "field_delimiter": ", ",
5   "feature_columns": "all",
6   "id_column": "PassengerId",
7   "prediction_column": "Survived",
8   "threshold": 0.36
9 }
10 }
```

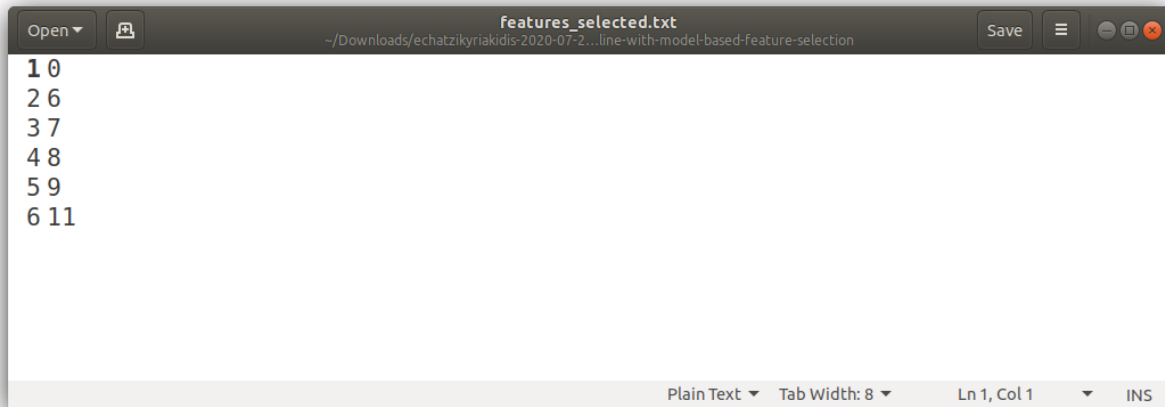
## 7.10 Experiment Logging

```
experiment.log
~/Downloads/echatzikyriakidis-2020-07-23T2...pipeline-with-model-based-feature-selection
Save

1 {
2   "datetime": "2020-07-23T23-01-21",
3   "experimenter": "echatzikyriakidis",
4   "experiment_id": "dc1095b704f546debd615928bcf7a994"
5 }
```

## 7.11 Features Selected

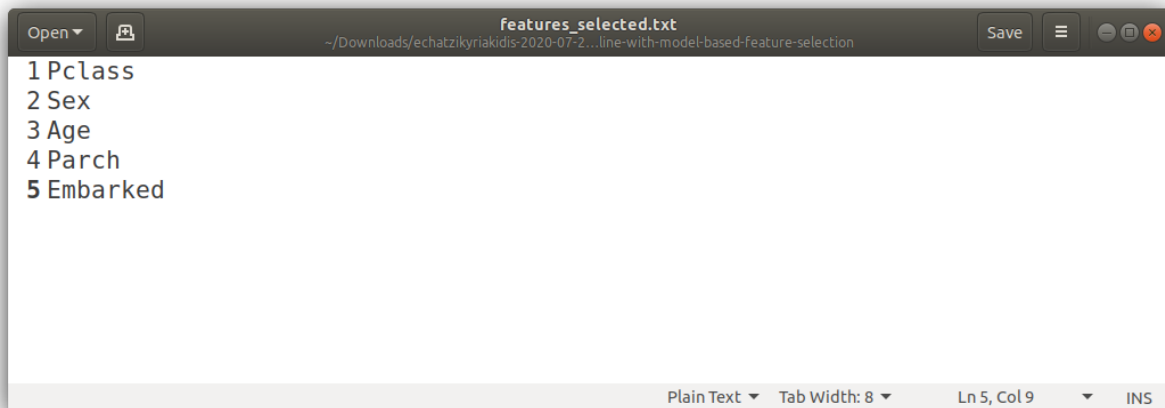
*The selected column indexes from the transformed features (this is generated when a preprocessor is used):*



A screenshot of a text editor window titled "features\_selected.txt". The window shows a list of six lines, each containing a line number followed by a space and a column index: "1 0", "2 6", "3 7", "4 8", "5 9", and "6 11". The editor interface includes a menu bar with "Open" and "Save" options, a status bar at the bottom indicating "Plain Text", "Tab Width: 8", and the current cursor position "Ln 1, Col 1".

```
1 0
2 6
3 7
4 8
5 9
6 11
```

*The selected column names from the original features (this is generated when no preprocessor is used):*



A screenshot of a text editor window titled "features\_selected.txt". The window shows a list of five lines, each containing a line number followed by a space and a column name: "1 Pclass", "2 Sex", "3 Age", "4 Parch", and "5 Embarked". The editor interface includes a menu bar with "Open" and "Save" options, a status bar at the bottom indicating "Plain Text", "Tab Width: 8", and the current cursor position "Ln 5, Col 9".

```
1 Pclass
2 Sex
3 Age
4 Parch
5 Embarked
```

## 7.12 Experiment Source Code

```

1 from sklearn.compose import ColumnTransformer
2 from sklearn.pipeline import Pipeline
3 from sklearn.impute import SimpleImputer
4 from sklearn.preprocessing import StandardScaler, OneHotEncoder
5 from sklearn.linear_model import LogisticRegression
6
7 from skrobot.core import Experiment
8 from skrobot.tasks import TrainTask
9 from skrobot.tasks import PredictionTask
10 from skrobot.tasks import FeatureSelectionCrossValidationTask
11 from skrobot.tasks import EvaluationCrossValidationTask
12 from skrobot.tasks import HyperParametersSearchCrossValidationTask
13 from skrobot.feature_selection import ColumnSelector
14 from skrobot.notification import BaseNotifier
15
16 ##### Initialization Code
17
18 train_data_set_file_path = 'https://bit.ly/titanic-data-train'
19
20 test_data_set_file_path = 'https://bit.ly/titanic-data-test'
21
22 new_data_set_file_path = 'https://bit.ly/titanic-data-new'
23
24 random_seed = 42
25
26 id_column = 'PassengerId'
27
28 label_column = 'Survived'
29
30 numerical_features = ['Age', 'Fare', 'SibSp', 'Parch']
31
32 categorical_features = ['Embarked', 'Sex', 'Pclass']
33
34 numeric_transformer = Pipeline(steps=[
35     ('imputer', SimpleImputer()),
36     ('scaler', StandardScaler())])
37
38 categorical_transformer = Pipeline(steps=[

```

## 7.13 Predictions

```

1 PassengerId,Survived,probability
2 530,0,0.24575978907773116
3 760,1,0.8367357017875539
4 532,0,0.13299699284500366
5 227,0,0.244909831257034
6 884,0,0.244909831257034
7 637,0,0.1333657497174666
8 718,1,0.7307918212381642
9 278,0,0.2361056710832674

```

## THE PEOPLE BEHIND IT?

Development:

- [Efstathios Chatzikyriakidis](#)

Support, testing and features recommendation:

- [Lefteris Kouloubri](#)
- [Antonis Markou](#)
- [Christina Chrysouli](#)
- [Michalis Chaviaras](#)

And last but not least, all the open-source contributors whose work went into [RELEASES](#).





## **CAN I CONTRIBUTE?**

Of course, the project is [Free Software](#) and you can contribute to it!



## WHAT LICENSE DO YOU USE?

See our [LICENSE](#) for more details.



## PYTHON MODULE INDEX

### S

`skrobot.core.experiment`, 3  
`skrobot.core.task_runner`, 5  
`skrobot.feature_selection.column_selector`,  
5  
`skrobot.notification.base_notifier`, 6  
`skrobot.tasks.base_cross_validation_task`,  
7  
`skrobot.tasks.base_task`, 8  
`skrobot.tasks.evaluation_cross_validation_task`,  
9  
`skrobot.tasks.feature_selection_cross_validation_task`,  
13  
`skrobot.tasks.hyperparameters_search_cross_validation_task`,  
16  
`skrobot.tasks.prediction_task`, 19  
`skrobot.tasks.train_task`, 20



## Symbols

`__init__()` (`skrobot.core.experiment.Experiment` method), 3  
`__init__()` (`skrobot.core.task_runner.TaskRunner` method), 5  
`__init__()` (`skrobot.feature_selection.column_selector.ColumnSelector` method), 5  
`__init__()` (`skrobot.tasks.base_cross_validation_task.BaseCrossValidationTask` method), 7  
`__init__()` (`skrobot.tasks.base_task.BaseTask` method), 8  
`__init__()` (`skrobot.tasks.evaluation_cross_validation_task.EvaluationCrossValidationTask` method), 10  
`__init__()` (`skrobot.tasks.feature_selection_cross_validation_task.FeatureSelectionCrossValidationTask` method), 13  
`__init__()` (`skrobot.tasks.hyperparameters_search_cross_validation_task.HyperparametersSearchCrossValidationTask` method), 16  
`__init__()` (`skrobot.tasks.prediction_task.PredictionTask` method), 19  
`__init__()` (`skrobot.tasks.train_task.TrainTask` method), 20

**B**  
`BaseCrossValidationTask` (class in `skrobot.tasks.base_cross_validation_task`), 7  
`BaseNotifier` (class in `skrobot.notification.base_notifier`), 6  
`BaseTask` (class in `skrobot.tasks.base_task`), 8  
`build()` (`skrobot.core.experiment.Experiment` method), 4

**C**  
`ColumnSelector` (class in `skrobot.feature_selection.column_selector`), 5  
`custom_folds()` (`skrobot.tasks.base_cross_validation_task.BaseCrossValidationTask` method), 7  
`custom_folds()` (`skrobot.tasks.evaluation_cross_validation_task.EvaluationCrossValidationTask` method), 11  
`custom_folds()` (`skrobot.tasks.feature_selection_cross_validation_task.FeatureSelectionCrossValidationTask` method), 14  
`custom_folds()` (`skrobot.tasks.hyperparameters_search_cross_validation_task.HyperparametersSearchCrossValidationTask` method), 18

**E**  
`EvaluationCrossValidationTask` (class in `skrobot.tasks.evaluation_cross_validation_task`), 9

**F**  
`FeatureSelectionCrossValidationTask` (class in `skrobot.tasks.feature_selection_cross_validation_task`), 13

**G**  
`get_configuration()` (`skrobot.tasks.base_cross_validation_task.BaseCrossValidationTask` method), 8  
`get_configuration()` (`skrobot.tasks.base_task.BaseTask` method), 8  
`get_configuration()` (`skrobot.tasks.evaluation_cross_validation_task.EvaluationCrossValidationTask` method), 12  
`get_configuration()` (`skrobot.tasks.feature_selection_cross_validation_task.FeatureSelectionCrossValidationTask` method), 15  
`get_configuration()` (`skrobot.tasks.hyperparameters_search_cross_validation_task.HyperparametersSearchCrossValidationTask` method), 18  
`get_configuration()` (`skrobot.tasks.prediction_task.PredictionTask` method), 19  
`get_configuration()` (`skrobot.tasks.train_task.TrainTask` method), 20

```

get_type() (skrobot.tasks.base_task.BaseTask method), 8
get_type() (skrobot.tasks.evaluation_cross_validation_task.EvaluationCrossValidationTask method), 12
get_type() (skrobot.tasks.feature_selection_cross_validation_task.FeatureSelectionCrossValidationTask method), 15
get_type() (skrobot.tasks.hyperparameters_search_cross_validation_task.HyperParametersSearchCrossValidationTask method), 18
get_type() (skrobot.tasks.prediction_task.PredictionTask method), 20
get_type() (skrobot.tasks.train_task.TrainTask method), 21
grid_search() (skrobot.tasks.hyperparameters_search_cross_validation_task.HyperParametersSearchCrossValidationTask method), 17

H
HyperParametersSearchCrossValidationTask (class in skrobot.tasks.hyperparameters_search_cross_validation_task), 16

M
module
    skrobot.core.experiment, 3
    skrobot.core.task_runner, 5
    skrobot.feature_selection.column_selector, 5
    skrobot.notification.base_notifier, 6
    skrobot.tasks.base_cross_validation_task, 7
    skrobot.tasks.base_task, 8
    skrobot.tasks.evaluation_cross_validation_task, 9
    skrobot.tasks.feature_selection_cross_validation_task, 13
    skrobot.tasks.hyperparameters_search_cross_validation_task, 16
    skrobot.tasks.prediction_task, 19
    skrobot.tasks.train_task, 20

N
notify() (skrobot.notification.base_notifier.BaseNotifier method), 6

P
PredictionTask (class in skrobot.tasks.prediction_task), 19

R
random_search() (skrobot.tasks.hyperparameters_search_cross_validation_task.HyperParametersSearchCrossValidationTask method), 18
run() (skrobot.core.experiment.Experiment method), 4
run() (skrobot.core.task_runner.TaskRunner method), 5
run() (skrobot.tasks.base_cross_validation_task.BaseCrossValidationTask method), 8
run() (skrobot.tasks.evaluation_cross_validation_task.EvaluationCrossValidationTask method), 9
run() (skrobot.tasks.feature_selection_cross_validation_task.FeatureSelectionCrossValidationTask method), 15
run() (skrobot.tasks.hyperparameters_search_cross_validation_task.HyperParametersSearchCrossValidationTask method), 18
run() (skrobot.tasks.prediction_task.PredictionTask method), 20
run() (skrobot.tasks.train_task.TrainTask method), 21
set_experimenter() (skrobot.core.experiment.Experiment method), 4
set_source_code_file_path() (skrobot.core.experiment.Experiment method), 3
skrobot.core.experiment, 3
skrobot.core.task_runner, 5
skrobot.feature_selection.column_selector, 5
skrobot.notification.base_notifier, 6
skrobot.tasks.base_cross_validation_task, 7
skrobot.tasks.base_task, 8
skrobot.tasks.evaluation_cross_validation_task, 9
skrobot.tasks.feature_selection_cross_validation_task, 13
skrobot.tasks.hyperparameters_search_cross_validation_task, 16
skrobot.tasks.prediction_task, 19
skrobot.tasks.train_task, 20
stratified_folds() (skrobot.tasks.base_cross_validation_task.BaseCrossValidationTask method), 7
stratified_folds() (skrobot.tasks.evaluation_cross_validation_task.EvaluationCrossValidationTask method), 12
stratified_folds() (skrobot.tasks.feature_selection_cross_validation_task.FeatureSelectionCrossValidationTask method), 15

```



`stratified_folds()`  
(*skrobot.tasks.hyperparameters\_search\_cross\_validation\_task.HyperParametersSearchCrossValidationTask*  
*method*), [19](#)

## T

`TaskRunner` (*class in skrobot.core.task\_runner*), [5](#)  
`TrainTask` (*class in skrobot.tasks.train\_task*), [20](#)  
`transform()` (*skrobot.feature\_selection.column\_selector.ColumnSelector*  
*method*), [6](#)