
EvalML Documentation

Release 0.8.0

Feature Labs, Inc.

Dec 14, 2020

GETTING STARTED

1	Install	3
2	Quick Start	5
	Index	165



EvalML

EvalML is an AutoML library that builds, optimizes, and evaluates machine learning pipelines using domain-specific objective functions.

Combined with [Featuretools](#) and [Compose](#), EvalML can be used to create end-to-end machine learning solutions for classification and regression problems.

INSTALL

EvalML is available for Python 3.5+. It can be installed by running the following command:

```
pip install evalml --extra-index-url https://install.featurelabs.com/<license>/
```

Note for Windows users: The [XGBoost](#) library may not be pip-installable in some Windows environments. If you are encountering installation issues, please try installing XGBoost from [Github](#) before installing EvalML.

Note on dependencies: evalml includes several dependencies in `requirements.txt` by default: `xgboost` and `catboost` support pipelines built around those modeling libraries, and `plotly` and `ipywidgets` support plotting functionality in automl searches. These dependencies are not required in order to install and use evalml.

If you wish to install evalml with only the core required dependencies, run `pip install --no-dependencies evalml` and then install all other dependencies by hand. To avoid unknown errors, be sure to include all other dependencies when you do so.

QUICK START

```
[1]: import evalml
      from evalml import AutoClassificationSearch
```

2.1 Load Data

First, we load in the features and outcomes we want to use to train our model

```
[2]: X, y = evalml.demos.load_breast_cancer()
```

2.2 Configure search

EvalML has many options to configure the pipeline search. At the minimum, we need to define an objective function. For simplicity, we will use the F1 score in this example. However, the real power of EvalML is in using domain-specific *objective functions* or *building your own*.

Below EvalML utilizes Bayesian optimization (EvalML's default optimizer) to search and find the best pipeline defined by the given objective.

```
[3]: automl = AutoClassificationSearch(objective="f1",
                                     max_pipelines=5)
```

In order to validate the results of the pipeline creation and optimization process, we will save some of our data as a holdout set.

```
[4]: X_train, X_holdout, y_train, y_holdout = evalml.preprocessing.split_data(X, y, test_
      ↪size=.2)
```

When we call `.search()`, the search for the best pipeline will begin. There is no need to wrangle with missing data or categorical variables as EvalML includes various preprocessing steps (like imputation, one-hot encoding, feature selection) to ensure you're getting the best results. As long as your data is in a single table, EvalML can handle it. If not, you can reduce your data to a single table by utilizing [Featuretools](#) and its Entity Sets.

You can find more information on pipeline components and how to integrate your own custom pipelines into EvalML [here](#).

```
[5]: automl.search(X_train, y_train)
```

```
*****
* Beginning pipeline search *
*****
```

Optimizing for F1. Greater score is better.

Searching up to 5 pipelines.

```
FigureWidget({
  'data': [{'mode': 'lines+markers',
            'name': 'Best Score',
            'type'...
```

```
Cat Boost Classification Pipeline:      20%|          | Elapsed:00:07
Logistic Regression Pipeline:          40%|          | Elapsed:00:08
Logistic Regression Pipeline:          60%|          | Elapsed:00:08
XGBoost Classification Pipeline:        80%|          | Elapsed:00:08
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:
```

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

```
[20:07:13] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:
```

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

```
[20:07:15] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:
```

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

```
[20:07:18] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.
```

```
XGBoost Classification Pipeline:        80%|          | Elapsed:00:15
Cat Boost Classification Pipeline:      100%|| Elapsed:00:16
Optimization finished                   100%|| Elapsed:00:16
```

2.3 See Pipeline Rankings

After the search is finished we can view all of the pipelines searched, ranked by score. Internally, EvalML performs cross validation to score the pipelines. If it notices a high variance across cross validation folds, it will warn you. EvalML also provides additional *guardrails* to analyze your data to assist you in producing the best performing pipeline.

```
[6]: automl.rankings
```

	id	pipeline_name	score	high_variance_cv	\
0	2	Logistic Regression Pipeline	0.975559	False	
1	1	Logistic Regression Pipeline	0.972061	False	
2	4	Cat Boost Classification Pipeline	0.968803	False	
3	0	Cat Boost Classification Pipeline	0.963343	False	
4	3	XGBoost Classification Pipeline	0.939615	False	


```

                                parameters
0  {'impute_strategy': 'median', 'penalty': 'l2',...
1  {'impute_strategy': 'median', 'penalty': 'l2',...
2  {'impute_strategy': 'most_frequent', 'n_estima...
3  {'impute_strategy': 'most_frequent', 'n_estima...
4  {'impute_strategy': 'most_frequent', 'percent_...

```

2.4 Describe pipeline

If we are interested in see more details about the pipeline, we can describe it using the `id` from the rankings table:

```
[7]: automl.describe_pipeline(3)
```

```

*****
* XGBoost Classification Pipeline *
*****

Supported Problem Types: Binary Classification, Multiclass Classification
Model Family: XGBoost Classifier
Objective to Optimize: F1 (greater is better)
Number of features: 2

Pipeline Steps
=====
1. One Hot Encoder
   * top_n : 10
2. Simple Imputer
   * impute_strategy : most_frequent
   * fill_value : None
3. RF Classifier Select From Model
   * percent_features : 0.08032569761590808
   * threshold : mean
4. XGBoost Classifier
   * eta : 0.020218397440325723
   * max_depth : 20
   * min_child_weight : 9.614396430577418
   * n_estimators : 848

Training

```

(continues on next page)

(continued from previous page)

```
=====
Training for Binary Classification problems.
Total training time (including CV): 6.9 seconds

Cross Validation
-----
```

	F1	Precision	Recall	AUC	Log Loss	MCC	# Training	# Testing
0	0.952	0.967	0.937	0.976	0.172	0.876	303.000	152.000
1	0.940	0.977	0.905	0.973	0.215	0.853	303.000	152.000
2	0.927	0.918	0.937	0.970	0.208	0.800	304.000	151.000
mean	0.940	0.954	0.926	0.973	0.198	0.843	-	-
std	0.012	0.032	0.018	0.003	0.023	0.039	-	-
coef of var	0.013	0.034	0.020	0.003	0.117	0.046	-	-

2.5 Select Best pipeline

We can now select best pipeline and score it on our holdout data:

```
[8]: pipeline = automl.best_pipeline
      pipeline.score(X_holdout, y_holdout)

[8]: (0.993103448275862, {})
```

We can also visualize the structure of our pipeline:

```
[9]: pipeline.graph()

[9]:
```

2.6 Whats next?

Head into the more in-depth automated walkthrough [here](#) or any advanced topics below.

2.6.1 Objective Functions

The **objective function** is what EvalML maximizes (or minimizes) as it completes the pipeline search. As it gets feedback from building pipelines, it tunes the hyperparameters to build optimized models. Therefore, it is critical to have an objective function that captures the how the model's predictions will be used in a business setting.

List of Available Objective Functions

Most AutoML libraries optimize for generic machine learning objective functions. Frequently, the scores produced by the generic machine learning objective diverge from how the model will be evaluated in the real world.

In EvalML, we can train and optimize the model for a specific problem by optimizing a domain-specific objectives functions or by defining our own custom objective function.

Currently, EvalML has two domain specific objective functions with more being developed. For more information on these objective functions click on the links below.

- [Fraud Detection](#)
- [Lead Scoring](#)

Build your own objective Functions

Often times, the objective function is very specific to the use-case or business problem. To get the right objective to optimize requires thinking through the decisions or actions that will be taken using the model and assigning the cost/benefit to doing that correctly or incorrectly based on known outcomes in the training data.

Once you have determined the objective for your business, you can provide that to EvalML to optimize by defining a custom objective function. Read more [here](#).

2.6.2 Building a Fraud Prediction Model with EvalML

In this demo, we will build an optimized fraud prediction model using EvalML. To optimize the pipeline, we will set up an objective function to minimize the percentage of total transaction value lost to fraud. At the end of this demo, we also show you how introducing the right objective during the training is over 4x better than using a generic machine learning metric like AUC.

```
[1]: import evalml
      from evalml import AutoClassificationSearch
      from evalml.objectives import FraudCost
```

Configure “Cost of Fraud”

To optimize the pipelines toward the specific business needs of this model, you can set your own assumptions for the cost of fraud. These parameters are

- `retry_percentage` - what percentage of customers will retry a transaction if it is declined?
- `interchange_fee` - how much of each successful transaction do you collect?
- `fraud_payout_percentage` - the percentage of fraud will you be unable to collect
- `amount_col` - the column in the data the represents the transaction amount

Using these parameters, EvalML determines attempt to build a pipeline that will minimize the financial loss due to fraud.

```
[2]: fraud_objective = FraudCost(retry_percentage=.5,
                                interchange_fee=.02,
                                fraud_payout_percentage=.75,
                                amount_col='amount')
```

Search for best pipeline

In order to validate the results of the pipeline creation and optimization process, we will save some of our data as a holdout set

```
[3]: X, y = evalml.demos.load_fraud(n_rows=2500)
```

	Number of Features
Boolean	1
Categorical	6
Numeric	5
Number of training examples: 2500	
Labels	

(continues on next page)

(continued from previous page)

```
False      85.92%
True       14.08%
Name: fraud, dtype: object
```

EvalML natively supports one-hot encoding. Here we keep 1 out of the 6 categorical columns to decrease computation time.

```
[4]: X = X.drop(['datetime', 'expiration_date', 'country', 'region', 'provider'], axis=1)

X_train, X_holdout, y_train, y_holdout = evalml.preprocessing.split_data(X, y, test_
↳size=0.2, random_state=0)

print(X.dtypes)

card_id          int64
store_id         int64
amount          int64
currency         object
customer_present bool
lat             float64
lng             float64
dtype: object
```

Because the fraud labels are binary, we will use `AutoClassificationSearch`. When we call `.search()`, the search for the best pipeline will begin.

```
[5]: automl = AutoClassificationSearch(objective=fraud_objective,
                                     additional_objectives=['auc', 'recall', 'precision
↳'],
                                     max_pipelines=5)

automl.search(X_train, y_train)

*****
* Beginning pipeline search *
*****

Optimizing for Fraud Cost. Lower score is better.

Searching up to 5 pipelines.

FigureWidget({
  'data': [{'mode': 'lines+markers',
            'name': 'Best Score',
            'type': ...

Cat Boost Classification Pipeline: 20%|          | Elapsed:00:03
Logistic Regression Pipeline:      40%|          | Elapsed:00:04
Logistic Regression Pipeline:      60%|          | Elapsed:00:05
XGBoost Classification Pipeline:    80%|          | Elapsed:00:05

/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳future release. To remove this warning, do the following: 1) Pass option_
↳use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your_
↳labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].
```

```
[20:05:13] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↪evaluation metric used with the objective 'binary:logistic' was changed from
↪'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↪behavior.
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↪lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:
```

```
The use of label encoder in XGBClassifier is deprecated and will be removed in a
↪future release. To remove this warning, do the following: 1) Pass option
↪use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↪labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].
```

```
[20:05:17] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↪evaluation metric used with the objective 'binary:logistic' was changed from
↪'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↪behavior.
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↪lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:
```

```
The use of label encoder in XGBClassifier is deprecated and will be removed in a
↪future release. To remove this warning, do the following: 1) Pass option
↪use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↪labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].
```

```
[20:05:20] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↪evaluation metric used with the objective 'binary:logistic' was changed from
↪'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↪behavior.
```

```
XGBoost Classification Pipeline:      80%| | Elapsed:00:16
Cat Boost Classification Pipeline:    100%| Elapsed:00:16
Optimization finished                 100%| Elapsed:00:16
```

View rankings and select pipeline

Once the fitting process is done, we can see all of the pipelines that were searched, ranked by their score on the fraud detection objective we defined

```
[6]: automl.rankings
```

```
[6]:
```

	id	pipeline_name	score	high_variance_cv	\
0	4	Cat Boost Classification Pipeline	0.007838	False	
1	3	XGBoost Classification Pipeline	0.007838	False	
2	2	Logistic Regression Pipeline	0.008140	False	
3	1	Logistic Regression Pipeline	0.008538	False	
4	0	Cat Boost Classification Pipeline	0.008906	False	

```

                                parameters
0 {'impute_strategy': 'most_frequent', 'n_estima...
1 {'impute_strategy': 'most_frequent', 'percent...
2 {'impute_strategy': 'median', 'penalty': 'l2',...
3 {'impute_strategy': 'median', 'penalty': 'l2',...
4 {'impute_strategy': 'most_frequent', 'n_estima...
```

to select the best pipeline we can run

```
[7]: best_pipeline = automl.best_pipeline
```

Describe pipeline

You can get more details about any pipeline. Including how it performed on other objective functions.

```
[8]: automl.describe_pipeline(automl.rankings.iloc[0]["id"])

*****
* Cat Boost Classification Pipeline *
*****

Supported Problem Types: Binary Classification, Multiclass Classification
Model Family: CatBoost Classifier
Objective to Optimize: Fraud Cost (lower is better)
Number of features: 7

Pipeline Steps
=====
1. Simple Imputer
   * impute_strategy : most_frequent
   * fill_value : None
2. CatBoost Classifier
   * n_estimators : 109
   * eta : 0.7991585642167238
   * max_depth : 4

Training
=====
Training for Binary Classification problems.
Total training time (including CV): 0.7 seconds

Cross Validation
-----
```

	Fraud Cost	AUC	Recall	Precision	# Training	# Testing
0	0.008	0.860	1.000	0.141	1333.000	667.000
1	0.008	0.819	1.000	0.141	1333.000	667.000
2	0.008	0.815	1.000	0.141	1334.000	666.000
mean	0.008	0.831	1.000	0.141	-	-
std	0.000	0.025	0.000	0.000	-	-
coef of var	0.007	0.030	0.000	0.001	-	-

Evaluate on hold out

Finally, we retrain the best pipeline on all of the training data and evaluate on the holdout

```
[9]: best_pipeline.fit(X_train, y_train)

[9]: <evalml.pipelines.classification.catboost.CatBoostClassificationPipeline at_
     ↪ 0x7ff679773b10>
```

Now, we can score the pipeline on the hold out data using both the fraud cost score and the AUC.

```
[10]: best_pipeline.score(X_holdout, y_holdout, other_objectives=["auc", fraud_objective])
```



```
[10]: (0.007766323050145169,
OrderedDict([('AUC', 0.8461794019933554),
('Fraud Cost', 0.007766323050145169)]))
```

Why optimize for a problem-specific objective?

To demonstrate the importance of optimizing for the right objective, let's search for another pipeline using AUC, a common machine learning metric. After that, we will score the holdout data using the fraud cost objective to see how the best pipelines compare.

```
[11]: automl_auc = AutoClassificationSearch(objective='auc',
additional_objectives=['recall', 'precision'],
max_pipelines=5)

automl_auc.search(X_train, y_train)
```

```
*****
* Beginning pipeline search *
*****

Optimizing for AUC. Greater score is better.

Searching up to 5 pipelines.
```

```
FigureWidget({
  'data': [{'mode': 'lines+markers',
            'name': 'Best Score',
            'type': ...}]
})
```

Cat Boost Classification Pipeline:	20%	Elapsed:00:03
Logistic Regression Pipeline:	40%	Elapsed:00:03
Logistic Regression Pipeline:	60%	Elapsed:00:04
XGBoost Classification Pipeline:	80%	Elapsed:00:04

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

[20:05:29] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.

/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

[20:05:32] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.
```

(continues on next page)

(continued from previous page)

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳future release. To remove this warning, do the following: 1) Pass option
↳use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

[20:05:35] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳evaluation metric used with the objective 'binary:logistic' was changed from
↳'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳behavior.
XGBoost Classification Pipeline:      80%|  | Elapsed:00:13
Cat Boost Classification Pipeline:    100%|| Elapsed:00:13
Optimization finished                 100%|| Elapsed:00:13
```

like before, we can look at the rankings and pick the best pipeline

```
[12]: automl_auc.rankings
```

```
[12]:      id      pipeline_name      score  high_variance_cv  \
0    3  XGBoost Classification Pipeline  0.846968          False
1    4  Cat Boost Classification Pipeline  0.837858          False
2    0  Cat Boost Classification Pipeline  0.835164          False
3    2      Logistic Regression Pipeline  0.806344          False
4    1      Logistic Regression Pipeline  0.806338          False

                                parameters
0  {'impute_strategy': 'most_frequent', 'percent...
1  {'impute_strategy': 'most_frequent', 'n_estima...
2  {'impute_strategy': 'most_frequent', 'n_estima...
3  {'impute_strategy': 'median', 'penalty': 'l2',...
4  {'impute_strategy': 'median', 'penalty': 'l2',...
```

```
[13]: best_pipeline_auc = automl_auc.best_pipeline
```

```
# train on the full training data
best_pipeline_auc.fit(X_train, y_train)
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳future release. To remove this warning, do the following: 1) Pass option
↳use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

[20:05:39] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳evaluation metric used with the objective 'binary:logistic' was changed from
↳'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳behavior.
```

```
[13]: <evalml.pipelines.classification.xgboost.XGBoostPipeline at 0x7ff679f32310>
```

```
[14]: # get the fraud score on holdout data
best_pipeline_auc.score(X_holdout, y_holdout, other_objectives=["auc", fraud_
↪objective])

[14]: (0.8151162790697674,
OrderedDict([('AUC', 0.8151162790697674),
('Fraud Cost', 0.03655681280302016)]))

[15]: # fraud score on fraud optimized again
best_pipeline.score(X_holdout, y_holdout, other_objectives=["auc", fraud_objective])

[15]: (0.007766323050145169,
OrderedDict([('AUC', 0.8461794019933554),
('Fraud Cost', 0.007766323050145169)]))
```

When we optimize for AUC, we can see that the AUC score from this pipeline is better than the AUC score from the pipeline optimized for fraud cost. However, the losses due to fraud are over 3% of the total transaction amount when optimized for AUC and under 1% when optimized for fraud cost. As a result, we lose more than 2% of the total transaction amount by not optimizing for fraud cost specifically.

This happens because optimizing for AUC does not take into account the user-specified `retry_percentage`, `interchange_fee`, `fraud_payout_percentage` values. Thus, the best pipelines may produce the highest AUC but may not actually reduce the amount loss due to your specific type fraud.

This example highlights how performance in the real world can diverge greatly from machine learning metrics.

2.6.3 Building a Lead Scoring Model with EvalML

In this demo, we will build an optimized lead scoring model using EvalML. To optimize the pipeline, we will set up an objective function to maximize the revenue generated with true positives while taking into account the cost of false positives. At the end of this demo, we also show you how introducing the right objective during the training is over 6x better than using a generic machine learning metric like AUC.

```
[1]: import evalml
from evalml import AutoClassificationSearch
from evalml.objectives import LeadScoring
```

Configure LeadScoring

To optimize the pipelines toward the specific business needs of this model, you can set your own assumptions for how much value is gained through true positives and the cost associated with false positives. These parameters are

- `true_positive` - dollar amount to be gained with a successful lead
- `false_positive` - dollar amount to be lost with an unsuccessful lead

Using these parameters, EvalML builds a pipeline that will maximize the amount of revenue per lead generated.

```
[2]: lead_scoring_objective = LeadScoring(
    true_positives=1000,
    false_positives=-10
)
```

Dataset

We will be utilizing a dataset detailing a customer's job, country, state, zip, online action, the dollar amount of that action and whether they were a successful lead.

```
[3]: import pandas as pd

customers = pd.read_csv('s3://featurelabs-static/lead_scoring_ml_apps/customers.csv')
interactions = pd.read_csv('s3://featurelabs-static/lead_scoring_ml_apps/interactions.
↪ csv')
leads = pd.read_csv('s3://featurelabs-static/lead_scoring_ml_apps/previous_leads.csv')

X = customers.merge(interactions, on='customer_id').merge(leads, on='customer_id')
y = X['label']

X = X.drop(['customer_id', 'date_registered', 'birthday', 'phone', 'email',
            'owner', 'company', 'id', 'time_x',
            'session', 'referrer', 'time_y', 'label'], axis=1)

display(X.head())
```

	job	country	state	zip	action	amount
0	Engineer, mining	NaN	NY	60091.0	page_view	NaN
1	Psychologist, forensic	US	CA	NaN	purchase	135.23
2	Psychologist, forensic	US	CA	NaN	page_view	NaN
3	Air cabin crew	US	NaN	60091.0	download	NaN
4	Air cabin crew	US	NaN	60091.0	page_view	NaN

Search for best pipeline

In order to validate the results of the pipeline creation and optimization process, we will save some of our data as a holdout set

EvalML natively supports one-hot encoding and imputation so the above NaN and categorical values will be taken care of.

```
[4]: X_train, X_holdout, y_train, y_holdout = evalml.preprocessing.split_data(X, y, test_
↪ size=0.2, random_state=0)

print(X.dtypes)
```

job	object
country	object
state	object
zip	float64
action	object
amount	float64
dtype:	object

Because the lead scoring labels are binary, we will use `AutoClassificationSearch`. When we call `.search()`, the search for the best pipeline will begin.

```
[5]: automl = AutoClassificationSearch(objective=lead_scoring_objective,
                                     additional_objectives=['auc'],
                                     max_pipelines=5)

automl.search(X_train, y_train)
```

```
*****
* Beginning pipeline search *
*****
```

Optimizing for Lead Scoring. Greater score is better.

Searching up to 5 pipelines.

```
FigureWidget({
  'data': [{'mode': 'lines+markers',
            'name': 'Best Score',
            'type'...
```

```
Cat Boost Classification Pipeline:      20%|          | Elapsed:00:04
Logistic Regression Pipeline:          40%|          | Elapsed:00:07
Logistic Regression Pipeline:          60%|          | Elapsed:00:08
XGBoost Classification Pipeline:        80%|          | Elapsed:00:08
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:
```

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

```
[20:05:57] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:
```

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

```
[20:06:00] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:
```

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

```
[20:06:03] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.
```

```
XGBoost Classification Pipeline:        80%|          | Elapsed:00:19
Cat Boost Classification Pipeline:      100%|         | Elapsed:00:20
Optimization finished                   100%|         | Elapsed:00:20
```

View rankings and select pipeline

Once the fitting process is done, we can see all of the pipelines that were searched, ranked by their score on the lead scoring objective we defined

```
[6]: automl.rankings
```

```
[6]:
```

	id	pipeline_name	score	high_variance_cv	\
0	3	XGBoost Classification Pipeline	14.706382	False	
1	0	Cat Boost Classification Pipeline	13.224749	True	
2	1	Logistic Regression Pipeline	12.973101	True	
3	2	Logistic Regression Pipeline	12.973101	True	
4	4	Cat Boost Classification Pipeline	11.848536	True	

```

                                parameters
0  {'impute_strategy': 'most_frequent', 'percent_...
1  {'impute_strategy': 'most_frequent', 'n_estima...
2  {'impute_strategy': 'median', 'penalty': 'l2',...
3  {'impute_strategy': 'median', 'penalty': 'l2',...
4  {'impute_strategy': 'most_frequent', 'n_estima...
```

to select the best pipeline we can run

```
[7]: best_pipeline = automl.best_pipeline
```

Describe pipeline

You can get more details about any pipeline. Including how it performed on other objective functions.

```
[8]: automl.describe_pipeline(automl.rankings.iloc[0]["id"])
```

```

*****
* XGBoost Classification Pipeline *
*****

Supported Problem Types: Binary Classification, Multiclass Classification
Model Family: XGBoost Classifier
Objective to Optimize: Lead Scoring (greater is better)
Number of features: 1

Pipeline Steps
=====
1. One Hot Encoder
    * top_n : 10
2. Simple Imputer
    * impute_strategy : most_frequent
    * fill_value : None
3. RF Classifier Select From Model
    * percent_features : 0.08032569761590808
    * threshold : mean
4. XGBoost Classifier
    * eta : 0.020218397440325723
    * max_depth : 20
    * min_child_weight : 9.614396430577418
    * n_estimators : 848
```

(continues on next page)

(continued from previous page)

```

Training
=====
Training for Binary Classification problems.
Total training time (including CV): 10.4 seconds

Cross Validation
-----

```

	Lead Scoring	AUC	# Training	# Testing
0	15.619	0.531	3099.000	1550.000
1	13.826	0.492	3099.000	1550.000
2	14.674	0.502	3100.000	1549.000
mean	14.706	0.508	-	-
std	0.897	0.020	-	-
coef of var	0.061	0.040	-	-

Evaluate on hold out

Finally, we retrain the best pipeline on all of the training data and evaluate on the holdout

```

[9]: best_pipeline.fit(X_train, y_train)

/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

[20:06:08] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.

[9]: <evalml.pipelines.classification.xgboost.XGBoostPipeline at 0x7fa6b179ee90>

```

Now, we can score the pipeline on the hold out data using both the lead scoring score and the AUC.

```

[10]: best_pipeline.score(X_holdout, y_holdout, other_objectives=["auc", lead_scoring_
↳ objective])

[10]: (11.453138435081685,
OrderedDict([('AUC', 0.515593834995467),
            ('Lead Scoring', 11.453138435081685)]))

```

Why optimize for a problem-specific objective?

To demonstrate the importance of optimizing for the right objective, let's search for another pipeline using AUC, a common machine learning metric. After that, we will score the holdout data using the lead scoring objective to see how the best pipelines compare.

```

[11]: automl_auc = evalml.AutoClassificationSearch(objective='auc',
additional_objectives=[],
max_pipelines=5)

```

(continues on next page)

(continued from previous page)

```
automl_auc.search(X_train, y_train)
```

```
*****
* Beginning pipeline search *
*****
```

```
Optimizing for AUC. Greater score is better.
```

```
Searching up to 5 pipelines.
```

```
FigureWidget({
  'data': [{'mode': 'lines+markers',
            'name': 'Best Score',
            'type'...
```

```
Cat Boost Classification Pipeline:    20%|          | Elapsed:00:04
Logistic Regression Pipeline:         40%|          | Elapsed:00:04
Logistic Regression Pipeline:         60%|          | Elapsed:00:05
XGBoost Classification Pipeline:       80%|          | Elapsed:00:05
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:
```

```
The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].
```

```
[20:06:19] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:
```

```
The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].
```

```
[20:06:22] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.
```

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:
```

```
The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].
```



```
[20:06:25] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↪evaluation metric used with the objective 'binary:logistic' was changed from
↪'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↪behavior.
XGBoost Classification Pipeline:      80%| | Elapsed:00:15
Cat Boost Classification Pipeline:    100%| Elapsed:00:15
Optimization finished                 100%| Elapsed:00:15
```

like before, we can look at the rankings and pick the best pipeline

```
[12]: automl_auc.rankings
```

```
[12]:      id      pipeline_name      score  high_variance_cv  \
0    4  Cat Boost Classification Pipeline  0.928045          False
1    0  Cat Boost Classification Pipeline  0.891013          False
2    2      Logistic Regression Pipeline  0.695681          False
3    1      Logistic Regression Pipeline  0.695662          False
4    3  XGBoost Classification Pipeline  0.509056          False

      parameters
0  {'impute_strategy': 'most_frequent', 'n_estima...
1  {'impute_strategy': 'most_frequent', 'n_estima...
2  {'impute_strategy': 'median', 'penalty': 'l2',...
3  {'impute_strategy': 'median', 'penalty': 'l2',...
4  {'impute_strategy': 'most_frequent', 'percent_...
```

```
[13]: best_pipeline_auc = automl_auc.best_pipeline
```

```
# train on the full training data
best_pipeline_auc.fit(X_train, y_train)
```

```
[13]: <evalml.pipelines.classification.catboost.CatBoostClassificationPipeline at
↪0x7fa6b1737d90>
```

```
[14]: # get the auc and lead scoring score on holdout data
```

```
best_pipeline_auc.score(X_holdout, y_holdout, other_objectives=["auc", lead_scoring_
↪objective])
```

```
[14]: (0.9197793895436688,
OrderedDict([('AUC', 0.9197793895436688),
            ('Lead Scoring', -0.017196904557179708)]))
```

When we optimize for AUC, we can see that the AUC score from this pipeline is better than the AUC score from the pipeline optimized for lead scoring. However, the revenue per lead gained was only \$7 per lead when optimized for AUC and was \$45 when optimized for lead scoring. As a result, we would gain up to 6x the amount of revenue if we optimized for lead scoring.

This happens because optimizing for AUC does not take into account the user-specified `true_positive` (dollar amount to be gained with a successful lead) and `false_positive` (dollar amount to be lost with an unsuccessful lead) values. Thus, the best pipelines may produce the highest AUC but may not actually generate the most revenue through lead scoring.

This example highlights how performance in the real world can diverge greatly from machine learning metrics.

2.6.4 Custom Objective Functions

Often times, the objective function is very specific to the use-case or business problem. To get the right objective to optimize requires thinking through the decisions or actions that will be taken using the model and assigning a

cost/benefit to doing that correctly or incorrectly based on known outcomes in the training data.

Once you have determined the objective for your business, you can provide that to EvalML to optimize by defining a custom objective function.

How to Create a Objective Function

To create a custom objective function, we must define 2 functions

- The **“objective function”**: this function takes the predictions, true labels, and any other information about the future and returns a score of how well the model performed.
- The **“decision function”**: this function takes prediction probabilities that were output from the model and a threshold and returns a prediction.

To evaluate a particular model, EvalML automatically finds the best threshold to pass to the decision function to generate predictions and then scores the resulting predictions using the objective function. The score from the objective function determines which set of pipeline hyperparameters EvalML will try next.

To give a concrete example, let’s look at how the fraud detection objective function is built.

```
[1]: from evalml.objectives.objective_base import ObjectiveBase

class FraudCost(ObjectiveBase):
    """Score the percentage of money lost of the total transaction amount process due
    to fraud"""
    name = "Fraud Cost"
    needs_fitting = True
    greater_is_better = False
    uses_extra_columns = True
    score_needs_proba = False

    def __init__(self, retry_percentage=.5, interchange_fee=.02,
                  fraud_payout_percentage=1.0, amount_col='amount', verbose=False):
        """Create instance of FraudCost

        Args:
            retry_percentage (float): what percentage of customers will retry a
            transaction if it
            is declined? Between 0 and 1. Defaults to .5

            interchange_fee (float): how much of each successful transaction do you
            collect?
            Between 0 and 1. Defaults to .02

            fraud_payout_percentage (float): how percentage of fraud will you be
            unable to collect.
            Between 0 and 1. Defaults to 1.0

            amount_col (str): name of column in data that contains the amount.
            defaults to "amount"
        """
        self.retry_percentage = retry_percentage
        self.interchange_fee = interchange_fee
        self.fraud_payout_percentage = fraud_payout_percentage
        self.amount_col = amount_col
        super().__init__(verbose=verbose)
```

(continues on next page)

(continued from previous page)

```

def decision_function(self, y_predicted, extra_cols, threshold):
    """Determine if transaction is fraud given predicted probabilities,
       dataframe with transaction amount, and threshold"""

    transformed_probs = (y_predicted * extra_cols[self.amount_col])
    return transformed_probs > threshold

def objective_function(self, y_predicted, y_true, extra_cols):
    """Calculate amount lost to fraud given predictions, true values, and
    ↪dataframe
       with transaction amount"""

    # extract transaction using the amount columns in users data
    transaction_amount = extra_cols[self.amount_col]

    # amount paid if transaction is fraud
    fraud_cost = transaction_amount * self.fraud_payout_percentage

    # money made from interchange fees on transaction
    interchange_cost = transaction_amount * (1 - self.retry_percentage) * self.
    ↪interchange_fee

    # calculate cost of missing fraudulent transactions
    false_negatives = (y_true & ~y_predicted) * fraud_cost

    # calculate money lost from fees
    false_positives = (~y_true & y_predicted) * interchange_cost

    loss = false_negatives.sum() + false_positives.sum()

    loss_per_total_processed = loss / transaction_amount.sum()

    return loss_per_total_processed

```

2.6.5 Setting up pipeline search

Designing the right machine learning pipeline and picking the best parameters is a time-consuming process that relies on a mix of data science intuition as well as trial and error. EvalML streamlines the process of selecting the best modeling algorithms and parameters, so data scientists can focus their energy where it is most needed.

How it works

EvalML selects and tunes machine learning pipelines built of numerous steps. This includes encoding categorical data, missing value imputation, feature selection, feature scaling, and finally machine learning. As EvalML tunes pipelines, it uses the objective function selected and configured by the user to guide its search.

At each iteration, EvalML uses cross-validation to generate an estimate of the pipeline's performances. If a pipeline has high variance across cross-validation folds, it will provide a warning. In this case, the pipeline may not perform reliably in the future.

EvalML is designed to work well out of the box. However, it provides numerous methods for you to control the search described below.

Selecting problem type

EvalML supports both classification and regression problems. You select your problem type by importing the appropriate class.

```
[1]: import evalml
    from evalml import AutoClassificationSearch, AutoRegressionSearch

[2]: AutoClassificationSearch()

[2]: <evalml.automl.auto_classification_search.AutoClassificationSearch at 0x7fef5c9a1390>

[3]: AutoRegressionSearch()

[3]: <evalml.automl.auto_regression_search.AutoRegressionSearch at 0x7fef1dc67550>
```

Setting the Objective Function

The only required parameter to start searching for pipelines is the objective function. Most domain-specific objective functions require you to specify parameters based on your business assumptions. You can do this before you initialize your pipeline search. For example

```
[4]: from evalml.objectives import FraudCost

    fraud_objective = FraudCost(
        retry_percentage=.5,
        interchange_fee=.02,
        fraud_payout_percentage=.75,
        amount_col='amount'
    )

    AutoClassificationSearch(objective=fraud_objective)

[4]: <evalml.automl.auto_classification_search.AutoClassificationSearch at 0x7fef5c9c9c50>
```

Evaluate on Additional Objectives

Additional objectives can be scored on during the evaluation process. To add another objective, use the `additional_objectives` parameter in `AutoClassificationSearch` or `AutoRegressionSearch`. The results of these additional objectives will then appear in the results of `describe_pipeline`.

```
[5]: from evalml.objectives import FraudCost

    fraud_objective = FraudCost(
        retry_percentage=.5,
        interchange_fee=.02,
        fraud_payout_percentage=.75,
        amount_col='amount'
    )

    AutoClassificationSearch(objective='AUC', additional_objectives=[fraud_objective])

[5]: <evalml.automl.auto_classification_search.AutoClassificationSearch at 0x7feef7cd1e10>
```

Selecting Model Types

By default, all model types are considered. You can control which model types to search with the `allowed_model_families` parameters

```
[6]: automl = AutoClassificationSearch(objective="f1",
                                     allowed_model_families=["random_forest"])
```

you can see the possible pipelines that will be searched after initialization

```
[7]: automl.possible_pipelines
[7]: [evalml.pipelines.classification.random_forest.RFClassificationPipeline]
```

you can see a list of all supported models like this

```
[8]: evalml.list_model_families("binary") # `binary` for binary classification and
      ↪ `multiclass` for multiclass classification
[8]: [<ModelFamily.RANDOM_FOREST: 'random_forest'>,
      <ModelFamily.CATBOOST: 'catboost'>,
      <ModelFamily.LINEAR_MODEL: 'linear_model'>,
      <ModelFamily.XGBOOST: 'xgboost'>]
```

```
[9]: evalml.list_model_families("regression")
[9]: [<ModelFamily.LINEAR_MODEL: 'linear_model'>,
      <ModelFamily.CATBOOST: 'catboost'>,
      <ModelFamily.RANDOM_FOREST: 'random_forest'>]
```

Limiting Search Time

You can limit the search time by specifying a maximum number of pipelines and/or a maximum amount of time. EvalML won't build new pipelines after the maximum time has passed or the maximum number of pipelines have been built. If a limit is not set, then a maximum of 5 pipelines will be built.

The maximum search time can be specified as a integer in seconds or as a string in seconds, minutes, or hours.

```
[10]: AutoClassificationSearch(objective="f1",
                               max_pipelines=5,
                               max_time=60)

AutoClassificationSearch(objective="f1",
                          max_time="1 minute")
[10]: <evalml.automl.auto_classification_search.AutoClassificationSearch at 0x7feef7c7c290>
```

To start, EvalML samples 10 sets of hyperparameters chosen randomly for each possible pipeline. Therefore, we recommend setting `max_pipelines` at least 10 times the number of possible pipelines.

```
[11]: n_possible_pipelines = len(AutoClassificationSearch(objective="f1").possible_
      ↪pipelines)

[12]: AutoClassificationSearch(objective="f1",
                              max_time=60)
[12]: <evalml.automl.auto_classification_search.AutoClassificationSearch at 0x7feef7c7c7d0>
```

Early Stopping

You can also limit search time by providing a patience value for early stopping. With a patience value, EvalML will stop searching when the best objective score has not been improved upon for n iterations. The patience value must be a positive integer. You can also provide a tolerance value where EvalML will only consider a score as an improvement over the best score if the difference was greater than the tolerance percentage.

```
[13]: from evalml.demos import load_diabetes

X, y = load_diabetes()
automl = AutoRegressionSearch(objective="MSE", patience=2, tolerance=0.01, max_
↳ pipelines=10)
automl.search(X, y)
```

```
*****
* Beginning pipeline search *
*****

Optimizing for MSE. Lower score is better.

Searching up to 10 pipelines.

FigureWidget({
  'data': [{'mode': 'lines+markers',
            'name': 'Best Score',
            'type': ...

Linear Regression Pipeline:          10%|          | Elapsed:00:00
Random Forest Regression Pipeline:  20%|          | Elapsed:00:05
Random Forest Regression Pipeline:  30%|          | Elapsed:00:15
Random Forest Regression Pipeline:  40%|          | Elapsed:00:24
Linear Regression Pipeline:          50%|          | Elapsed:00:24

2 iterations without improvement. Stopping search early...
Optimization finished                50%|          | Elapsed:00:24
```

```
[14]: automl.rankings
```

```
[14]:
```

	id	pipeline_name	score	high_variance_cv	\
0	2	Random Forest Regression Pipeline	3614.391928	False	
1	1	Random Forest Regression Pipeline	3760.685056	False	
2	3	Random Forest Regression Pipeline	3769.284869	False	
3	0	Linear Regression Pipeline	26225.781788	False	
4	4	Linear Regression Pipeline	26225.781788	False	

```

parameters
0 {'impute_strategy': 'most_frequent', 'percent_...
1 {'impute_strategy': 'median', 'percent_feature...
2 {'impute_strategy': 'median', 'percent_feature...
3 {'impute_strategy': 'median', 'fit_intercept':...
4 {'impute_strategy': 'most_frequent', 'fit_inte...
```

Control Cross Validation

EvalML cross-validates each model it tests during its search. By default it uses 3-fold cross-validation. You can optionally provide your own cross-validation method.

```
[15]: from sklearn.model_selection import StratifiedKFold

automl = AutoClassificationSearch(objective="f1",
                                cv=StratifiedKFold(5))
```

2.6.6 Exploring search results

After finishing a pipeline search, we can inspect the results. First, let's build a search of 10 different pipelines to explore.

```
[1]: import evalml
from evalml import AutoClassificationSearch

X, y = evalml.demos.load_breast_cancer()

automl = AutoClassificationSearch(objective="f1",
                                max_pipelines=5)

automl.search(X, y)
```

```
*****
* Beginning pipeline search *
*****

Optimizing for F1. Greater score is better.

Searching up to 5 pipelines.
```

```
FigureWidget({
  'data': [{'mode': 'lines+markers',
            'name': 'Best Score',
            'type': ...
```

Cat Boost Classification Pipeline:	20%	Elapsed:00:07
Logistic Regression Pipeline:	40%	Elapsed:00:08
Logistic Regression Pipeline:	60%	Elapsed:00:08
XGBoost Classification Pipeline:	80%	Elapsed:00:08

```
/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

[20:04:56] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↳ evaluation metric used with the objective 'binary:logistic' was changed from
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↳ behavior.

/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↳ future release. To remove this warning, do the following: 1) Pass option
↳ use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↳ labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].
```

(continues on next page)

(continued from previous page)

```
[20:04:58] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↪evaluation metric used with the objective 'binary:logistic' was changed from
↪'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↪behavior.

/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↪lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning:

The use of label encoder in XGBClassifier is deprecated and will be removed in a
↪future release. To remove this warning, do the following: 1) Pass option
↪use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your
↪labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].

[20:05:01] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default
↪evaluation metric used with the objective 'binary:logistic' was changed from
↪'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old
↪behavior.
XGBoost Classification Pipeline:          80%| | Elapsed:00:15
Cat Boost Classification Pipeline:       100%| Elapsed:00:16
Optimization finished                    100%| Elapsed:00:16
```

View Rankings

A summary of all the pipelines built can be returned as a pandas DataFrame. It is sorted by score. EvalML knows based on our objective function whether higher or lower is better.

```
[2]: automl.rankings

[2]:   id      pipeline_name      score  high_variance_cv  \
0    2  Logistic Regression Pipeline  0.976371          False
1    1  Logistic Regression Pipeline  0.974941          False
2    0  Cat Boost Classification Pipeline  0.973464          False
3    4  Cat Boost Classification Pipeline  0.966453          False
4    3  XGBoost Classification Pipeline  0.935899          False

      parameters
0  {'impute_strategy': 'median', 'penalty': 'l2',...
1  {'impute_strategy': 'median', 'penalty': 'l2',...
2  {'impute_strategy': 'most_frequent', 'n_estima...
3  {'impute_strategy': 'most_frequent', 'n_estima...
4  {'impute_strategy': 'most_frequent', 'percent...
```

Describe Pipeline

Each pipeline is given an id. We can get more information about any particular pipeline using that id. Here, we will get more information about the pipeline with id = 0.

```
[3]: automl.describe_pipeline(0)

*****
* Cat Boost Classification Pipeline *
*****
```

(continues on next page)

(continued from previous page)

```
Supported Problem Types: Binary Classification, Multiclass Classification
Model Family: CatBoost Classifier
Objective to Optimize: F1 (greater is better)
Number of features: 30

Pipeline Steps
=====
1. Simple Imputer
    * impute_strategy : most_frequent
    * fill_value : None
2. CatBoost Classifier
    * n_estimators : 369
    * eta : 0.4236547993389048
    * max_depth : 6

Training
=====
Training for Binary Classification problems.
Total training time (including CV): 7.2 seconds

Cross Validation
-----
```

	F1	Precision	Recall	AUC	Log Loss	MCC	# Training	# Testing
0	0.954	0.958	0.950	0.985	0.163	0.877	379.000	190.000
1	0.983	0.967	1.000	0.996	0.092	0.955	379.000	190.000
2	0.983	0.975	0.992	0.996	0.076	0.955	380.000	189.000
mean	0.973	0.967	0.980	0.993	0.110	0.929	-	-
std	0.017	0.009	0.027	0.006	0.046	0.045	-	-
coef of var	0.018	0.009	0.028	0.007	0.418	0.049	-	-

Get Pipeline

We can get the object of any pipeline via their id as well:

```
[4]: automl.get_pipeline(0)
[4]: <evalml.pipelines.classification.catboost.CatBoostClassificationPipeline at 0x7f7aea94c650>
```

Get best pipeline

If we specifically want to get the best pipeline, there is a convenient access

```
[5]: automl.best_pipeline
[5]: <evalml.pipelines.classification.logistic_regression.LogisticRegressionPipeline at 0x7f7aeac74e90>
```

Feature Importances

We can get the feature importances of the resulting pipeline

```
[6]: pipeline = automl.get_pipeline(0)
      pipeline.feature_importances
```

```
[6]:
```

	feature	importance
0	worst radius	26.873775
1	worst texture	7.612433
2	mean concave points	7.145519
3	mean texture	7.024352
4	worst symmetry	6.361393
5	worst concave points	4.964684
6	worst perimeter	3.983425
7	mean fractal dimension	3.922057
8	radius error	3.447547
9	worst smoothness	2.880494
10	worst area	2.583201
11	worst concavity	2.413365
12	mean concavity	2.407919
13	mean radius	2.284723
14	mean compactness	2.012972
15	perimeter error	1.918985
16	mean symmetry	1.740673
17	worst compactness	1.655886
18	area error	1.606211
19	mean smoothness	1.443077
20	compactness error	1.435200
21	concave points error	1.040189
22	mean area	0.965964
23	smoothness error	0.834013
24	worst fractal dimension	0.436113
25	fractal dimension error	0.401008
26	mean perimeter	0.284681
27	texture error	0.182944
28	symmetry error	0.118948
29	concavity error	0.018248

We can also create a bar plot of the feature importances

```
[7]: pipeline.feature_importance_graph(pipeline)
```

Data type cannot be displayed: application/vnd.plotly.v1+json, text/html

Plot ROC

For binary classification tasks, we can also plot the ROC plot of a specific pipeline:

```
[8]: automl.plot.generate_roc_plot(0)
```

Data type cannot be displayed: application/vnd.plotly.v1+json, text/html

Access raw results

You can also get access to all the underlying data like this

```
[9]: automl.results
[9]: {'pipeline_results': {0: {'id': 0,
    'pipeline_name': 'Cat Boost Classification Pipeline',
    'pipeline_summary': 'CatBoost Classifier w/ Simple Imputer',
    'parameters': {'impute_strategy': 'most_frequent',
    'n_estimators': 369,
    'eta': 0.4236547993389048,
    'max_depth': 6},
    'score': 0.9734636352012647,
    'high_variance_cv': False,
    'training_time': 7.151692628860474,
    'cv_data': [{'all_objective_scores': OrderedDict([('F1',
    0.9535864978902954),
    ('Precision', 0.9576271186440678),
    ('Recall', 0.9495798319327731),
    ('AUC', 0.9850869925434962),
    ('Log Loss', 0.16289034062986435),
    ('MCC', 0.8767221685540966),
    ('ROC',
    (array([0.          , 0.          , 0.          , 0.01408451, 0.01408451,
    0.02816901, 0.02816901, 0.05633803, 0.05633803, 0.07042254,
    0.07042254, 0.08450704, 0.08450704, 0.14084507, 0.14084507,
    1.          ]),
    array([0.          , 0.00840336, 0.34453782, 0.34453782, 0.84033613,
    0.84033613, 0.93277311, 0.93277311, 0.94117647, 0.94117647,
    0.98319328, 0.98319328, 0.99159664, 0.99159664, 1.          ,
    1.          ])),
    array([1.99997424e+00, 9.99974244e-01, 9.99756520e-01, 9.
    ↪99745844e-01,
    9.88277062e-01, 9.81932629e-01, 7.35284284e-01, 6.
    ↪10443075e-01,
    5.94587484e-01, 5.82814719e-01, 9.45094859e-02, 4.
    ↪30750748e-02,
    3.99827348e-02, 1.48833097e-02, 1.48214210e-02, 3.
    ↪31091266e-05]))),
    ('Confusion Matrix',
    0    1
    0   66    5
    1    6  113),
    ('# Training', 379),
    ('# Testing', 190)]),
    'score': 0.9535864978902954},
    {'all_objective_scores': OrderedDict([('F1', 0.9834710743801653),
    ('Precision', 0.967479674796748),
    ('Recall', 1.0),
    ('AUC', 0.996330926736892),
    ('Log Loss', 0.09239450377828266),
    ('MCC', 0.9554966130892879),
    ('ROC',
    (array([0.          , 0.          , 0.          , 0.01408451, 0.01408451,
    0.02816901, 0.02816901, 0.04225352, 0.04225352, 1.          ↪
    ↪])),
    array([0.          , 0.00840336, 0.80672269, 0.80672269, 0.94957983,
    0.94957983, 0.98319328, 0.98319328, 1.          , 1.          ↪
    ↪])),
    array([1.99998792e+00, 9.99987915e-01, 9.98961984e-01, 9.
    ↪98941711e-01,
```

(continues on next page)

(continues on next page)

(continued from previous page)

```

1.          ]),
array([0.          , 0.00840336, 0.59663866, 0.59663866, 0.86554622,
0.86554622, 0.91596639, 0.91596639, 0.94117647, 0.94117647,
0.97478992, 0.97478992, 0.99159664, 0.99159664, 1.          ,
1.          ]),
array([2.00000000e+00, 1.00000000e+00, 9.99832384e-01, 9.
↪99827305e-01,
9.78642320e-01, 9.78473934e-01, 8.85901577e-01, 8.
↪46173258e-01,
8.11073923e-01, 7.85780332e-01, 5.55932686e-01, 5.
↪12973949e-01,
3.58737371e-01, 3.80356180e-04, 2.16404977e-04, 3.
↪27364828e-24]))),
('Confusion Matrix',
0  1
0 66  5
1  3 116),
('# Training', 379),
('# Testing', 190))),
'score': 0.9666666666666667},
{'all_objective_scores': OrderedDict([('F1', 0.979253112033195),
('Precision', 0.9672131147540983),
('Recall', 0.9915966386554622),
('AUC', 0.9983429991714996),
('Log Loss', 0.054377948871000094),
('MCC', 0.943843520216036),
('ROC',
array([0.          , 0.          , 0.          , 0.01408451, 0.01408451,
0.02816901, 0.02816901, 0.05633803, 0.05633803, 0.09859155,
0.09859155, 1.          ]),
array([0.          , 0.00840336, 0.96638655, 0.96638655, 0.97478992,
0.97478992, 0.98319328, 0.98319328, 0.99159664, 0.99159664,
1.          , 1.          ]),
array([2.00000000e+00, 1.00000000e+00, 9.02530651e-01, 8.
↪56837110e-01,
8.32273897e-01, 7.83350982e-01, 7.29357685e-01, 5.
↪16534638e-01,
5.11585048e-01, 2.30456831e-01, 2.07185770e-01, 4.
↪44775362e-44]))),
('Confusion Matrix',
0  1
0 67  4
1  1 118),
('# Training', 379),
('# Testing', 190))),
'score': 0.979253112033195},
{'all_objective_scores': OrderedDict([('F1', 0.9789029535864979),
('Precision', 0.9830508474576272),
('Recall', 0.9747899159663865),
('AUC', 0.9960384153661465),
('Log Loss', 0.07207639820328837),
('MCC', 0.9435040132749904),
('ROC',
array([0.          , 0.          , 0.          , 0.01428571, 0.01428571,
0.02857143, 0.02857143, 0.04285714, 0.04285714, 1.          ↪
↪])),
array([0.          , 0.00840336, 0.77310924, 0.77310924, 0.96638655,

```

(continues on next page)

(continued from previous page)

```

0.96638655, 0.98319328, 0.98319328, 1.          , 1.
↪)],
        array([2.00000000e+00, 9.99999998e-01, 9.91123955e-01, 9.
↪89983111e-01,
        5.77651350e-01, 5.45079286e-01, 4.73073055e-01, 4.
↪09422136e-01,
        3.23433980e-01, 1.26109142e-17]))),
        ('Confusion Matrix',
         0    1
         0 68    2
         1    3 116),
        ('# Training', 380),
        ('# Testing', 189)]),
        'score': 0.9789029535864979}}],
2: {'id': 2,
    'pipeline_name': 'Logistic Regression Pipeline',
    'pipeline_summary': 'Logistic Regression Classifier w/ One Hot Encoder + Simple_
↪Imputer + Standard Scaler',
    'parameters': {'impute_strategy': 'median',
                    'penalty': 'l2',
                    'C': 4.804971952026824},
    'score': 0.976371018670262,
    'high_variance_cv': False,
    'training_time': 0.35906362533569336,
    'cv_data': [{'all_objective_scores': OrderedDict([('F1',
0.9666666666666667),
('Precision', 0.9586776859504132),
('Recall', 0.9747899159663865),
('AUC', 0.9899396378269617),
('Log Loss', 0.13000491300733263),
('MCC', 0.9097672817424011),
('ROC',
(array([0.          , 0.          , 0.          , 0.01408451, 0.01408451,
0.02816901, 0.02816901, 0.04225352, 0.04225352, 0.07042254,
0.07042254, 0.12676056, 0.12676056, 1.          ]),
array([0.          , 0.00840336, 0.59663866, 0.59663866, 0.85714286,
0.85714286, 0.93277311, 0.93277311, 0.96638655, 0.96638655,
0.99159664, 0.99159664, 1.          , 1.          ])),
array([2.00000000e+00, 1.00000000e+00, 9.99513300e-01, 9.
↪99411733e-01,
        9.72040206e-01, 9.64245540e-01, 8.24677844e-01, 8.
↪13137106e-01,
        6.23606303e-01, 5.74962709e-01, 3.88150049e-01, 1.
↪05841228e-01,
        8.88048039e-03, 1.95380932e-20]))),
        ('Confusion Matrix',
         0    1
         0 66    5
         1    3 116),
        ('# Training', 379),
        ('# Testing', 190)]),
        'score': 0.9666666666666667},
{'all_objective_scores': OrderedDict([('F1', 0.979253112033195),
('Precision', 0.9672131147540983),
('Recall', 0.9915966386554622),
('AUC', 0.9986980707776069),
('Log Loss', 0.051792807657331755),

```

(continues on next page)

(continued from previous page)

```

('MCC', 0.943843520216036),
('ROC',
 (array([0.          , 0.          , 0.          , 0.01408451, 0.01408451,
         0.04225352, 0.04225352, 0.08450704, 0.08450704, 1.          ]
↪))),
        array([0.          , 0.00840336, 0.96638655, 0.96638655, 0.98319328,
         0.98319328, 0.99159664, 0.99159664, 1.          , 1.          ]
↪))),
        array([2.00000000e+00, 9.99999999e-01, 9.15699240e-01, 8.
↪64894910e-01,
              7.51407817e-01, 6.63047926e-01, 6.26895191e-01, 4.
↪06750268e-01,
              3.09309176e-01, 1.59809498e-36]))),
('Confusion Matrix',
 0   1
0  67   4
1   1 118),
('# Training', 379),
('# Testing', 190)]),
'score': 0.979253112033195},
{'all_objective_scores': OrderedDict([('F1', 0.9831932773109243),
('Precision', 0.9831932773109243),
('Recall', 0.9831932773109243),
('AUC', 0.9965186074429773),
('Log Loss', 0.06435139219399516),
('MCC', 0.9546218487394958),
('ROC',
 (array([0.          , 0.          , 0.          , 0.01428571, 0.01428571,
         0.02857143, 0.02857143, 0.04285714, 0.04285714, 1.          ]
↪))),
        array([0.          , 0.00840336, 0.78151261, 0.78151261, 0.98319328,
         0.98319328, 0.99159664, 0.99159664, 1.          , 1.          ]
↪))),
        array([1.99999996e+00, 9.99999962e-01, 9.84583454e-01, 9.
↪76666534e-01,
              5.25351686e-01, 5.03605398e-01, 4.55872870e-01, 4.
↪36470183e-01,
              4.28464345e-01, 2.38837351e-14]))),
('Confusion Matrix',
 0   1
0  68   2
1   2 117),
('# Training', 380),
('# Testing', 189)]),
'score': 0.9831932773109243}}],
3: {'id': 3,
'pipeline_name': 'XGBoost Classification Pipeline',
'pipeline_summary': 'XGBoost Classifier w/ One Hot Encoder + Simple Imputer + RF
↪Classifier Select From Model',
'parameters': {'impute_strategy': 'most_frequent',
'percent_features': 0.08032569761590808,
'threshold': 'mean',
'eta': 0.020218397440325723,
'max_depth': 20,
'min_child_weight': 9.614396430577418,
'n_estimators': 848},
'score': 0.9358991763960539,

```

(continues on next page)

(continued from previous page)

```

'high_variance_cv': False,
'training_time': 6.933280944824219,
'cv_data': [{'all_objective_scores': OrderedDict([('F1',
0.9224137931034483),
('Precision', 0.9469026548672567),
('Recall', 0.8991596638655462),
('AUC', 0.9666232690259202),
('Log Loss', 0.2562118422632155),
('MCC', 0.8027688833516422),
('ROC',
(array([0.          , 0.01408451, 0.01408451, 0.01408451, 0.01408451,
0.01408451, 0.01408451, 0.01408451, 0.04225352, 0.04225352,
0.07042254, 0.07042254, 0.12676056, 0.12676056, 0.12676056,
0.15492958, 0.16901408, 0.18309859, 0.21126761, 0.25352113,
0.28169014, 0.36619718, 0.3943662 , 0.3943662 , 0.42253521,
0.45070423, 1.          ]),
array([0.          , 0.54621849, 0.62184874, 0.63865546, 0.65546218,
0.66386555, 0.71428571, 0.7394958 , 0.7394958 , 0.8487395 ,
0.86554622, 0.89915966, 0.89915966, 0.90756303, 0.94117647,
0.94117647, 0.96638655, 0.96638655, 0.96638655, 0.98319328,
0.99159664, 0.99159664, 0.99159664, 1.          , 1.          ,
1.          , 1.          ]),
array([1.9729209 , 0.97292095, 0.9717207 , 0.9710014 , 0.970229 ,
0.9697926 , 0.9423327 , 0.9413457 , 0.93782085, 0.92524594,
0.6996592 , 0.5642555 , 0.40121877, 0.28772253, 0.27541634,
0.16674423, 0.15323758, 0.06644133, 0.06502564, 0.04827413,
0.04722596, 0.04633046, 0.04199466, 0.03978153, 0.03914706,
0.03638416, 0.03603544], dtype=float32)))],
('Confusion Matrix',
0 1
0 65 6
1 12 107),
('# Training', 379),
('# Testing', 190)]),
'score': 0.9224137931034483},
{'all_objective_scores': OrderedDict([('F1', 0.937007874015748),
('Precision', 0.8814814814814815),
('Recall', 1.0),
('AUC', 0.9770978814060836),
('Log Loss', 0.23696691198764663),
('MCC', 0.826339982903411),
('ROC',
(array([0.          , 0.          , 0.          , 0.05633803, 0.05633803,
0.05633803, 0.08450704, 0.08450704, 0.14084507, 0.15492958,
0.15492958, 0.16901408, 0.16901408, 0.23943662, 0.26760563,
0.29577465, 0.32394366, 0.38028169, 0.43661972, 0.46478873,
0.49295775, 1.          ]),
array([0.          , 0.49579832, 0.55462185, 0.85714286, 0.89915966,
0.90756303, 0.90756303, 0.92436975, 0.96638655, 0.96638655,
0.98319328, 0.99159664, 1.          , 1.          , 1.          ,
1.          , 1.          , 1.          , 1.          , 1.          ,
1.          , 1.          , 1.          ]),
array([1.9701126 , 0.97011256, 0.96989095, 0.9684937 , 0.9554615 ,
0.92179286, 0.9034773 , 0.83730716, 0.8372069 , 0.75368625,
0.63067174, 0.62702805, 0.61390054, 0.49871758, 0.48841092,
0.42385587, 0.41924828, 0.08368076, 0.04792023, 0.0469924 ,
0.04657101, 0.04570549], dtype=float32)))],

```

(continues on next page)

(continued from previous page)

```

        ('Confusion Matrix',
         0    1
         0 55   16
         1    0 119),
        ('# Training', 379),
        ('# Testing', 190)]),
    'score': 0.937007874015748},
{'all_objective_scores': OrderedDict([('F1', 0.9482758620689655),
    ('Precision', 0.9734513274336283),
    ('Recall', 0.9243697478991597),
    ('AUC', 0.9885954381752702),
    ('Log Loss', 0.1391644817456682),
    ('MCC', 0.8681704699715063),
    ('ROC',
     (array([0.          , 0.          , 0.          , 0.          , 0.          ,
              0.          , 0.          , 0.          , 0.          , 0.          ,
              0.          , 0.          , 0.          , 0.          , 0.          ,
              0.01428571, 0.01428571, 0.04285714, 0.04285714, 0.07142857,
              0.07142857, 0.08571429, 0.1          , 0.1          , 0.12857143,
              0.12857143, 0.17142857, 0.2          , 0.22857143, 0.22857143,
              0.28571429, 0.34285714, 0.38571429, 0.41428571, 0.42857143,
              0.47142857, 0.55714286, 0.61428571, 1.          ]),
      array([0.          , 0.28571429, 0.42016807, 0.44537815, 0.48739496,
              0.50420168, 0.55462185, 0.59663866, 0.60504202, 0.67226891,
              0.71428571, 0.74789916, 0.83193277, 0.85714286, 0.87394958,
              0.87394958, 0.90756303, 0.90756303, 0.93277311, 0.93277311,
              0.94117647, 0.94117647, 0.95798319, 0.96638655, 0.96638655,
              0.97478992, 0.97478992, 0.97478992, 0.97478992, 1.          ,
              1.          , 1.          , 1.          , 1.          ,
              1.          , 1.          , 1.          , 1.          ]),
      array([1.9841778 , 0.9841778 , 0.9838165 , 0.98342854, 0.9833071 ,
              0.98211664, 0.9810844 , 0.9787828 , 0.97762203, 0.9773312 ,
              0.975222 , 0.97440153, 0.89735514, 0.8955815 , 0.86498946,
              0.84353524, 0.810082 , 0.67681974, 0.49621513, 0.4496832 ,
              0.44711152, 0.4437709 , 0.43844405, 0.4295119 , 0.39135826,
              0.33074102, 0.33018234, 0.31763443, 0.29066208, 0.18481019,
              0.03418758, 0.03268759, 0.0274495 , 0.0270752 , 0.0269705 ,
              0.02627881, 0.02563108, 0.02260989, 0.0213987 ]),
      dtype=float32))),
    ('Confusion Matrix',
     0    1
     0 67   3
     1    9 110),
    ('# Training', 380),
    ('# Testing', 189)]),
    'score': 0.9482758620689655}}],
4: {'id': 4,
    'pipeline_name': 'Cat Boost Classification Pipeline',
    'pipeline_summary': 'CatBoost Classifier w/ Simple Imputer',
    'parameters': {'impute_strategy': 'most_frequent',
    'n_estimators': 109,
    'eta': 0.7991585642167238,
    'max_depth': 4},
    'score': 0.9664525764397752,
    'high_variance_cv': False,
    'training_time': 0.9719579219818115,
    'cv_data': [{'all_objective_scores': OrderedDict([('F1',

```

(continues on next page)

(continued from previous page)

```

0.9617021276595743),
('Precision', 0.9741379310344828),
('Recall', 0.9495798319327731),
('AUC', 0.9865072789679252),
('Log Loss', 0.16313398857144826),
('MCC', 0.9001633057441626),
('ROC',
 (array([0.          , 0.          , 0.          , 0.01408451, 0.01408451,
          0.02816901, 0.02816901, 0.04225352, 0.04225352, 0.05633803,
          0.05633803, 0.07042254, 0.07042254, 0.08450704, 0.08450704,
          1.          ]),
  array([0.          , 0.00840336, 0.34453782, 0.34453782, 0.84033613,
          0.84033613, 0.93277311, 0.93277311, 0.96638655, 0.96638655,
          0.97478992, 0.97478992, 0.98319328, 0.98319328, 1.          ,
          1.          ]),
  array([1.99999633e+00, 9.99996325e-01, 9.99942123e-01, 9.
↪99941223e-01,
          9.83889653e-01, 9.82884213e-01, 7.60889858e-01, 7.
↪00312819e-01,
          4.44799151e-01, 4.02629215e-01, 3.25498190e-01, 2.
↪40859456e-01,
          7.31464857e-02, 4.01578493e-02, 1.91156164e-02, 3.
↪89213440e-06]))),
('Confusion Matrix',
  0    1
  0  68    3
  1    6 113),
('# Training', 379),
('# Testing', 190))),
'score': 0.9617021276595743},
{'all_objective_scores': OrderedDict([('F1', 0.9626556016597511),
('Precision', 0.9508196721311475),
('Recall', 0.9747899159663865),
('AUC', 0.9874541365842111),
('Log Loss', 0.14974588678331058),
('MCC', 0.8984549429340701),
('ROC',
 (array([0.          , 0.          , 0.          , 0.01408451, 0.01408451,
          0.02816901, 0.02816901, 0.04225352, 0.04225352, 0.07042254,
          0.07042254, 0.08450704, 0.08450704, 1.          ]),
  array([0.          , 0.00840336, 0.41176471, 0.41176471, 0.83193277,
          0.83193277, 0.95798319, 0.95798319, 0.96638655, 0.96638655,
          0.97478992, 0.97478992, 1.          , 1.          ]),
  array([1.99999714e+00, 9.99997136e-01, 9.99900758e-01, 9.
↪99895250e-01,
          9.97706259e-01, 9.97275707e-01, 9.75388340e-01, 9.
↪39646167e-01,
          9.37469604e-01, 8.54371454e-01, 8.53389442e-01, 8.
↪42251927e-01,
          3.54062858e-01, 2.61300949e-06]))),
('Confusion Matrix',
  0    1
  0  65    6
  1    3 116),
('# Training', 379),
('# Testing', 190))),
'score': 0.9626556016597511},

```

(continues on next page)

(continued from previous page)

```
{'all_objective_scores': OrderedDict([('F1', 0.975),
                                     ('Precision', 0.9669421487603306),
                                     ('Recall', 0.9831932773109243),
                                     ('AUC', 0.9980792316926771),
                                     ('Log Loss', 0.06313235445917639),
                                     ('MCC', 0.9317727223276882),
                                     ('ROC',
                                      (array([0.          , 0.          , 0.          , 0.01428571, 0.01428571,
                                                0.02857143, 0.02857143, 0.05714286, 0.05714286, 0.07142857,
                                                0.07142857, 1.          ]),
                                       array([0.          , 0.00840336, 0.95798319, 0.95798319, 0.96638655,
                                                0.96638655, 0.97478992, 0.97478992, 0.99159664, 0.99159664,
                                                1.          , 1.          ]),
                                       array([1.99999867e+00, 9.99998668e-01, 9.22815771e-01, 8.
↪86401336e-01,
                                                8.72320864e-01, 8.49344800e-01, 7.87185615e-01, 7.
↪06786188e-01,
                                                4.73751084e-01, 3.98242251e-01, 3.72969449e-01, 1.
↪06636001e-06])))),
                                     ('Confusion Matrix',
                                      (0      1
                                       0 66    4
                                       1  2 117),
                                      ('# Training', 380),
                                      ('# Testing', 189))),
                                     'score': 0.975}}],
'search_order': [0, 1, 2, 3, 4])
```

2.6.7 Regression Example

```
[1]: import evalml
      from evalml import AutoRegressionSearch
      from evalml.demos import load_diabetes
      from evalml.pipelines import PipelineBase, get_pipelines

X, y = evalml.demos.load_diabetes()

automl = AutoRegressionSearch(objective="R2", max_pipelines=5)

automl.search(X, y)

*****
* Beginning pipeline search *
*****

Optimizing for R2. Greater score is better.

Searching up to 5 pipelines.

FigureWidget({
  'data': [{'mode': 'lines+markers',
            'name': 'Best Score',
            'type': ...
```

```

Linear Regression Pipeline:      20%|          | Elapsed:00:00
Random Forest Regression Pipeline: 40%|          | Elapsed:00:06
Random Forest Regression Pipeline: 60%|          | Elapsed:00:15
Random Forest Regression Pipeline: 80%|          | Elapsed:00:24
Linear Regression Pipeline:      100%| Elapsed:00:24
Optimization finished           100%| Elapsed:00:24

```

[2]: automl.rankings

```

[2]:   id      pipeline_name      score  high_variance_cv  \
0    2  Random Forest Regression Pipeline  0.388582          False
1    1  Random Forest Regression Pipeline  0.362982          False
2    3  Random Forest Regression Pipeline  0.361598          False
3    0      Linear Regression Pipeline -3.441755          False
4    4      Linear Regression Pipeline -3.441755          False

                                parameters
0  {'impute_strategy': 'most_frequent', 'percent_...
1  {'impute_strategy': 'median', 'percent_feature...
2  {'impute_strategy': 'median', 'percent_feature...
3  {'impute_strategy': 'median', 'fit_intercept':...
4  {'impute_strategy': 'most_frequent', 'fit_inte...

```

[3]: automl.best_pipeline

[3]: <evalml.pipelines.regression.random_forest.RFRegressionPipeline at 0x7fcdd9fe7410>

[4]: automl.get_pipeline(0)

[4]: <evalml.pipelines.regression.linear_regression.LinearRegressionPipeline at 0x7fcde387f450>

[5]: automl.describe_pipeline(0)

```

*****
* Linear Regression Pipeline *
*****

Supported Problem Types: Regression
Model Family: Linear Model
Objective to Optimize: R2 (greater is better)
Number of features: 10

Pipeline Steps
=====
1. One Hot Encoder
   * top_n : 10
2. Simple Imputer
   * impute_strategy : median
   * fill_value : None
3. Standard Scaler
4. Linear Regressor
   * fit_intercept : False
   * normalize : True

Training
=====
Training for Regression problems.

```

(continues on next page)

(continued from previous page)

Total training time (including CV): 0.1 seconds

Cross Validation

	R2	MAE	MSE	MedianAE	MaxError	ExpVariance	# Training	#
↪Testing								
0	-3.872	157.576	27852.855	156.885	312.634	0.471	294.000	
↪148.000								
1	-3.114	151.160	26101.724	147.264	271.512	0.487	295.000	
↪147.000								
2	-3.340	148.098	24722.766	149.071	294.176	0.510	295.000	
↪147.000								
mean	-3.442	152.278	26225.782	151.073	292.774	0.490	-	
↪-								
std	0.389	4.837	1568.727	5.113	20.597	0.020	-	
↪-								
coef of var	-0.113	0.032	0.060	0.034	0.070	0.040	-	
↪-								

2.6.8 EvalML Components and Pipelines

EvalML searches and trains multiple machine learning **pipelines** in order to find the best one for your data. Each pipeline is made up of various **components** that can learn from the data, transform the data and ultimately predict labels given new data. Below we'll show an example of an EvalML pipeline. You can find a more in-depth look into [components](#) or learn how you can construct and use your own [pipelines](#).

XGBoost Pipeline

The EvalML XGBoost Pipeline is made up of four different components: a one-hot encoder, a missing value imputer, a feature selector and an XGBoost estimator. We can see them here by calling `.plot()`:

```
[1]: from evalml.demos import load_breast_cancer
from evalml.pipelines import XGBoostPipeline

X, y = load_breast_cancer()

objective='recall'
parameters = {
    'Simple Imputer': {
        'impute_strategy': 'mean'
    },
    'RF Classifier Select From Model': {
        "percent_features": 0.5,
        "number_features": X.shape[1],
        "n_estimators": 20,
        "max_depth": 5
    },
    'XGBoost Classifier': {
        "n_estimators": 20,
        "eta": 0.5,
        "min_child_weight": 5,
        "max_depth": 10,
    }
}
```

(continues on next page)

(continued from previous page)

```

    }

xgp = XGBoostPipeline(objective='recall', parameters=parameters, random_state=5)
xgp.graph()

```

[1]:

From the above graph we can see each component and its parameters. Each component takes in data and feeds it to the next. You can see more detailed information by calling `.describe()`:

[2]: `xgp.describe()`

```

*****
* XGBoost Classification Pipeline *
*****

Supported Problem Types: Binary Classification, Multiclass Classification
Model Family: XGBoost Classifier
Objective to Optimize: Recall (greater is better)

Pipeline Steps
=====
1. One Hot Encoder
   * top_n : 10
2. Simple Imputer
   * impute_strategy : mean
   * fill_value : None
3. RF Classifier Select From Model
   * percent_features : 0.5
   * threshold : -inf
4. XGBoost Classifier
   * eta : 0.5
   * max_depth : 10
   * min_child_weight : 5
   * n_estimators : 20

```

You can then fit and score an individual pipeline:

[3]: `xgp.fit(X, y)`
`xgp.score(X, y)`

```

/home/docs/checkouts/readthedocs.org/user_builds/feature-labs-inc-evalml/envs/v0.8.0/
↳ lib/python3.7/site-packages/xgboost/sklearn.py:888: UserWarning: The use of label_
↳ encoder in XGBClassifier is deprecated and will be removed in a future release. To
↳ remove this warning, do the following: 1) Pass option use_label_encoder=False when
↳ constructing XGBClassifier object; and 2) Encode your labels (y) as integers
↳ starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].
   warnings.warn(label_encoder_deprecation_msg, UserWarning)

```

```

[20:07:31] WARNING: ../src/learner.cc:1061: Starting in XGBoost 1.3.0, the default_
↳ evaluation metric used with the objective 'binary:logistic' was changed from_
↳ 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old_
↳ behavior.

```

[3]: (0.9971988795518207, {})

2.6.9 EvalML Components

From the [overview](#), we see how each machine learning pipeline consists of individual components that process data before the data is ultimately sent to an estimator. Below we will describe each type of component in an EvalML pipeline.

Component Classes

Components can be split into two distinct classes: **transformers** and **estimators**.

```
[1]: import numpy as np
import pandas as pd
from evalml.pipelines.components import SimpleImputer

X = pd.DataFrame([[1, 2, 3], [1, np.nan, 3]])
display(X)
```

	0	1	2
0	1	2	3
1	1	NaN	3

Transformers take in data as input and output altered data. For example, an *imputer* takes in data and outputs filled in missing data with the mean, median, or most frequent value of each column.

A transformer can fit on data and then transform it in two steps by calling `.fit()` and `.transform()` or in one step by calling `fit_transform()`.

```
[2]: imp = SimpleImputer(impute_strategy="mean")
X = imp.fit_transform(X)

display(X)
```

	0	1	2
0	1	2.0	3
1	1	2.0	3

On the other hand, an estimator fits on data (X) and labels (y) in order to take in new data as input and return the predicted label as output. Therefore, an estimator can fit on data and labels by calling `.fit()` and then predict by calling `.predict()` on new data. An example of this would be the *LogisticRegressionClassifier*. We can now see how a transformer alters data to make it easier for an estimator to learn and predict.

```
[3]: from evalml.pipelines.components import LogisticRegressionClassifier

clf = LogisticRegressionClassifier()

X = X
y = [1, 0]

clf.fit(X, y)
clf.predict(X)
```

```
[3]: array([0, 0])
```

Component Types

Components can further separate into different types that serve different functionality. Below we will go over the different types of transformers and estimators.

Transformer Types

- Imputer: fills missing data
 - Ex: *SimpleImputer*
- Scaler: alters numerical data into different scales
 - Ex: *StandardScaler*
- Encoder: translates different data types
 - Ex: *OneHotEncoder*
- Feature Selection: selects most useful columns of data
 - Ex: *RFCClassifierSelectFromModel*

Estimator Types

- Regressor: predicts numerical or continuous labels
 - Ex: *LinearRegressor*
- Classifier: predicts categorical or discrete labels
 - Ex: *XGBoostClassifier*

2.6.10 Custom Pipelines in EvalML

EvalML pipelines consist of modular components combining any number of transformers and an estimator. This allows you to create pipelines that fit the needs of your data to achieve the best results. You can create your own pipeline like this:

```
[1]: from evalml.pipelines import PipelineBase
from evalml.pipelines.components import StandardScaler, SimpleImputer
from evalml.pipelines.components.estimators import LogisticRegressionClassifier

# objectives can be either a str or the evalml objective object
objective = 'Precision_Macro'

# pipeline needs to be a subclass of `PipelineBase`
class CustomPipeline(PipelineBase):
    # component_graph and problem_types are required class variables

    # components can be passed in as objects or as component name strings
    component_graph = ['Simple Imputer', StandardScaler(), 'Logistic Regression_
↳Classifier']
    supported_problem_types = ['binary', 'multiclass']

    # you can override component hyperparameter_ranges like so
    # ranges must adhere to skopt tuner
    custom_hyperparameters = {
        "impute_strategy": ["most_frequent"]
    }
```

(continues on next page)

(continued from previous page)

```
# a parameters dictionary is necessary to instantiate pipelines
parameters = {
    'Simple Imputer':{
        'impute_strategy':"most_frequent"
    },
    'Logistic Regression Classifier':{
        'penalty':'l2',
        'C':5,
    }
}

pipeline = CustomPipeline(parameters={}, objective=objective, random_state=3)
```

```
[2]: from evalml.demos import load_wine
```

```
X, y = load_wine()

pipeline.fit(X, y)
pipeline.score(X, y)
```

```
[2]: (1.0, {})
```

2.6.11 Guardrails

EvalML provides guardrails to help guide you in achieving the highest performing model. These utility functions help deal with overfitting, abnormal data, and missing data. These guardrails can be found under `evalml/guardrails/` `utils`. Below we will cover abnormal and missing data guardrails. You can find an in-depth look into overfitting guardrails [here](#).

Missing Data

Missing data or rows with NaN values provide many challenges for machine learning pipelines. In the worst case, many algorithms simply will not run with missing data! EvalML pipelines contain imputation *components* to ensure that doesn't happen. Imputation works by approximating missing values with existing values. However, if a column contains a high number of missing values a large percentage of the column would be approximated by a small percentage. This could potentially create a column without useful information for machine learning pipelines. By running the `detect_highly_null()` guardrail, EvalML will alert you to this potential problem by returning the columns that pass the missing values threshold.

```
[1]: import numpy as np
import pandas as pd

from evalml.guardrails.utils import detect_highly_null

X = pd.DataFrame(
    [
        [1, 2, 3],
        [0, 4, np.nan],
        [1, 4, np.nan],
        [9, 4, np.nan],
        [8, 6, np.nan]
    ]
)
```

(continues on next page)

(continued from previous page)

```
detect_highly_null(X, percent_threshold=0.8)
```

```
[1]: {2: 0.8}
```

Abnormal Data

EvalML provides two utility functions to check for abnormal data: `detect_outliers()` and `detect_id_columns()`.

ID Columns

ID columns in your dataset provide little to no benefit to a machine learning pipeline as the pipeline cannot extrapolate useful information from unique identifiers. Thus, `detect_id_columns()` reminds you if these columns exists.

```
[2]: from evalml.guardrails.utils import detect_id_columns

X = pd.DataFrame([[0, 53, 6325, 5], [1, 90, 6325, 10], [2, 90, 18, 20]], columns=['user_
↪number', 'cost', 'revenue', 'id'])

display(X)
print(detect_id_columns(X, threshold=0.95))
```

	user_number	cost	revenue	id
0	0	53	6325	5
1	1	90	6325	10
2	2	90	18	20

```
{'id': 1.0, 'user_number': 0.95}
```

Outliers

Outliers are observations that differ significantly from other observations in the same sample. Many machine learning pipelines suffer in performance if outliers are not dropped from the training set as they are not representative of the data. `detect_outliers()` uses Isolation Forests to notify you if a sample can be considered an outlier.

Below we generate a random dataset with some outliers.

```
[3]: data = np.random.randn(100, 100)
X = pd.DataFrame(data=data)

# outliers
X.iloc[3, :] = pd.Series(np.random.randn(100) * 10)
X.iloc[25, :] = pd.Series(np.random.randn(100) * 20)
X.iloc[55, :] = pd.Series(np.random.randn(100) * 100)
X.iloc[72, :] = pd.Series(np.random.randn(100) * 100)
```

We then utilize `detect_outliers` to rediscover these outliers.

```
[4]: from evalml.guardrails.utils import detect_outliers

detect_outliers(X)
```

```
[4]: [3, 25, 55, 72]
```

2.6.12 Avoiding Overfitting

The ultimate goal of machine learning is to make accurate predictions on unseen data. EvalML aims to help you build a model that will perform as you expect once it is deployed in to the real world.

One of the benefits of using EvalML to build models is that it provides guardrails to ensure you are building pipelines that will perform reliably in the future. This page describes the various ways EvalML helps you avoid overfitting to your data.

```
[1]: import evalml
```

Detecting Label Leakage

A common problem is having features that include information from your label in your training data. By default, EvalML will provide a warning when it detects this may be the case.

Let's set up a simple example to demonstrate what this looks like

```
[2]: import pandas as pd

X = pd.DataFrame({
    "leaked_feature": [6, 6, 10, 5, 5, 11, 5, 10, 11, 4],
    "leaked_feature_2": [3, 2.5, 5, 2.5, 3, 5.5, 2, 5, 5.5, 2],
    "valid_feature": [3, 1, 3, 2, 4, 6, 1, 3, 3, 11]
})

y = pd.Series([1, 1, 0, 1, 1, 0, 1, 0, 0, 1])

automl = evalml.AutoClassificationSearch(
    max_pipelines=1,
    allowed_model_families=["linear_model"],
)

automl.search(X, y)

*****
* Beginning pipeline search *
*****

Optimizing for Precision. Greater score is better.

Searching up to 1 pipelines.
WARNING: Possible label leakage: leaked_feature, leaked_feature_2

FigureWidget({
  'data': [{'mode': 'lines+markers',
            'name': 'Best Score',
            'type': ...

Logistic Regression Pipeline:      100%| Elapsed:00:01
Optimization finished              100%| Elapsed:00:01
```

In the example above, EvalML warned about the input features `leaked_feature` and `leak_feature_2`, which are both very closely correlated with the label we are trying to predict. If you'd like to turn this check off, set `detect_label_leakage=False`.

The second way to find features that may be leaking label information is to look at the top features of the model. As we can see below, the top features in our model are the 2 leaked features.

```
[3]: best_pipeline = automl.best_pipeline
best_pipeline.feature_importances

[3]:
```

	feature	importance
0	leaked_feature	-1.789393
1	leaked_feature_2	-1.645127
2	valid_feature	-0.398465

Perform cross-validation for pipeline evaluation

By default, EvalML performs 3-fold cross validation when building pipelines. This means that it evaluates each pipeline 3 times using different sets of data for training and testing. In each trial, the data used for testing has no overlap from the data used for training.

While this is a good baseline approach, you can pass your own cross validation object to be used during modeling. The cross validation object can be any of the CV methods defined in [scikit-learn](#) or use a compatible API.

For example, if we wanted to do a time series split:

```
[4]: from sklearn.model_selection import TimeSeriesSplit

X, y = evalml.demos.load_breast_cancer()

automl = evalml.AutoClassificationSearch(
    cv=TimeSeriesSplit(n_splits=6),
    max_pipelines=1
)

automl.search(X, y)

*****
* Beginning pipeline search *
*****

Optimizing for Precision. Greater score is better.

Searching up to 1 pipelines.

FigureWidget({
  'data': [{'mode': 'lines+markers',
            'name': 'Best Score',
            'type': ...

Cat Boost Classification Pipeline:      100%|| Elapsed:00:11
Optimization finished                   100%|| Elapsed:00:11
```

if we describe the 1 pipeline we built, we can see the scores for each of the 6 splits as determined by the cross-validation object we provided. We can also see the number of training examples per fold increased because we were using TimeSeriesSplit

```
[5]: automl.describe_pipeline(0)

*****
* Cat Boost Classification Pipeline *
*****
```

(continues on next page)

(continued from previous page)

```
Supported Problem Types: Binary Classification, Multiclass Classification
Model Family: CatBoost Classifier
Objective to Optimize: Precision (greater is better)
Number of features: 30

Pipeline Steps
=====
1. Simple Imputer
    * impute_strategy : most_frequent
    * fill_value : None
2. CatBoost Classifier
    * n_estimators : 369
    * eta : 0.4236547993389048
    * max_depth : 6

Training
=====
Training for Binary Classification problems.
Total training time (including CV): 11.8 seconds

Cross Validation
-----
```

	Precision	F1	Recall	AUC	Log Loss	MCC	# Training	# Testing
0	0.952	0.851	0.769	0.937	0.870	0.672	83.000	81.000
1	0.976	0.952	0.930	0.998	0.083	0.902	164.000	81.000
2	0.982	0.991	1.000	0.995	0.097	0.972	245.000	81.000
3	1.000	0.973	0.948	1.000	0.049	0.916	326.000	81.000
4	1.000	0.992	0.984	1.000	0.065	0.964	407.000	81.000
5	0.967	0.975	0.983	0.994	0.073	0.903	488.000	81.000
mean	0.980	0.956	0.936	0.987	0.206	0.888	-	-
std	0.019	0.053	0.086	0.025	0.326	0.110	-	-
coef of var	0.019	0.056	0.092	0.025	1.579	0.124	-	-

Detect unstable pipelines

When we perform cross validation we are trying generate an estimate of pipeline performance. EvalML does this by taking the mean of the score across the folds. If the performance across the folds varies greatly, it is indicative the the estimated value may be unreliable.

To protect the user against this, EvalML checks to see if the pipeline’s performance has a variance between the different folds. EvalML triggers a warning if the “coefficient of variance” of the scores (the standard deviation divided by mean) of the pipelines scores exceeds .2.

This warning will appear in the pipeline rankings under `high_variance_cv`.

```
[6]: automl.rankings
[6]:   id          pipeline_name      score  high_variance_cv  \
0   0  Cat Boost Classification Pipeline  0.979504             False

      parameters
0  {'impute_strategy': 'most_frequent', 'n_estima...
```

Create holdout for model validation

EvalML offers a method to quickly create an holdout validation set. A holdout validation set is data that is not used during the process of optimizing or training the model. You should only use this validation set once you've picked the final model you'd like to use.

Below we create a holdout set of 20% of our data

```
[7]: X, y = evalml.demos.load_breast_cancer()
      X_train, X_holdout, y_train, y_holdout = evalml.preprocessing.split_data(X, y, test_
      ↪size=.2)
```

```
[8]: automl = evalml.AutoClassificationSearch(
      objective="recall",
      max_pipelines=3,
      detect_label_leakage=True
    )
    automl.search(X_train, y_train)

*****
* Beginning pipeline search *
*****

Optimizing for Recall. Greater score is better.

Searching up to 3 pipelines.

FigureWidget({
  'data': [{'mode': 'lines+markers',
            'name': 'Best Score',
            'type': ...

Cat Boost Classification Pipeline:      33%|      | Elapsed:00:07
Logistic Regression Pipeline:           67%|      | Elapsed:00:07
Logistic Regression Pipeline:          100%|| Elapsed:00:07
Optimization finished                   100%|| Elapsed:00:07
```

then we can retrain the best pipeline on all of our training data and see how it performs compared to the estimate

```
[9]: pipeline = automl.best_pipeline
      pipeline.fit(X_train, y_train)
      pipeline.score(X_holdout, y_holdout)

[9]: (0.9583333333333334, {})
```

2.6.13 Changelog

Future Releases

- Enhancements
- Fixes
- Changes
- Documentation Changes
- Testing Changes

v0.8.0 Apr. 1, 2020

- **Enhancements**

- Add normalization option and information to confusion matrix [#484](#)
- Add util function to drop rows with NaN values [#487](#)
- Renamed *PipelineBase.name* as *PipelineBase.summary* and redefined *PipelineBase.name* as class property [#491](#)
- Added access to parameters in Pipelines with *PipelineBase.parameters* (used to be return of *PipelineBase.describe*) [#501](#)
- Added *fill_value* parameter for SimpleImputer [#509](#)
- Added functionality to override component hyperparameters and made pipelines take hyperparameters from components [#516](#)
- Allow `numpy.random.RandomState` for *random_state* parameters [#556](#)

- Fixes

- **Changes**

- Undo version cap in XGBoost placed in [#402](#) and allowed all released of XGBoost [#407](#)
- Support pandas 1.0.0 [#486](#)
- Made all references to the logger static [#503](#)
- Refactored *model_type* parameter for components and pipelines to *model_family* [#507](#)
- Refactored *problem_types* for pipelines and components into *supported_problem_types* [#515](#)
- Moved *pipelines/utls.save_pipeline* and *pipelines/utls.load_pipeline* to *PipelineBase.save* and *PipelineBase.load* [#526](#)
- Limit number of categories encoded by OneHotEncoder [#517](#)

- **Documentation Changes**

- Updated API reference to remove PipelinePlot and added moved PipelineBase plotting methods [#483](#)
- Add code style and github issue guides [#463](#) [#512](#)
- Updated API reference for to surface class variables for pipelines and components [#537](#)

- **Testing Changes**

- Added automated dependency check PR [#482](#), [#505](#)
- Updated automated dependency check comment [#497](#)
- Have build_docs job use python executor, so that env vars are set properly [#547](#)
- Run windows unit tests on PRs [#557](#)

Warning: Breaking Changes

- *AutoClassificationSearch* and *AutoRegressionSearch*'s *model_types* parameter has been refactored into *allowed_model_families*
- *ModelTypes* enum has been changed to *ModelFamily*
- Components and Pipelines now have a *model_family* field instead of *model_type*
- *get_pipelines* utility function now accepts *model_families* as an argument instead of *model_types*

- *PipelineBase.name* no longer returns structure of pipeline and has been replaced by *PipelineBase.summary*
- *PipelineBase.problem_types* and *Estimator.problem_types* has been renamed to *supported_problem_types*
- *pipelines/utls.save_pipeline* and *pipelines/utls.load_pipeline* moved to *PipelineBase.save* and *PipelineBase.load*

v0.7.0 Mar. 9, 2020

- **Enhancements**

- Added emacs buffers to .gitignore #350
- Add CatBoost (gradient-boosted trees) classification and regression components and pipelines #247
- Added Tuner abstract base class #351
- Added n_jobs as parameter for AutoClassificationSearch and AutoRegressionSearch #403
- Changed colors of confusion matrix to shades of blue and updated axis order to match scikit-learn's #426
- Added PipelineBase graph and feature_importance_graph methods, moved from previous location #423
- Added support for python 3.8 #462

- **Fixes**

- Fixed ROC and confusion matrix plots not being calculated if user passed own additional_objectives #276
- Fixed ReadtheDocs FileNotFoundError exception for fraud dataset #439

- **Changes**

- Added n_estimators as a tunable parameter for XGBoost #307
- Remove unused parameter ObjectiveBase.fit_needs_proba #320
- Remove extraneous parameter component_type from all components #361
- Remove unused rankings.csv file #397
- Downloaded demo and test datasets so unit tests can run offline #408
- Remove _needs_fitting attribute from Components #398
- Changed plot.feature_importance to show only non-zero feature importances by default, added optional parameter to show all #413
- Refactored PipelineBase to take in parameter dictionary and moved pipeline metadata to class attribute #421
- Dropped support for Python 3.5 #438
- Removed unused apply.py file #449
- Clean up requirements.txt to remove unused deps #451
- Support installation without all required dependencies #459

- **Documentation Changes**

- Update release.md with instructions to release to internal license key #354

- **Testing Changes**

- Added tests for utils (and moved current utils to gen_utils) [#297](#)
- Moved XGBoost install into it's own separate step on Windows using Conda [#313](#)
- Rewind pandas version to before 1.0.0, to diagnose test failures for that version [#325](#)
- Added dependency update checkin test [#324](#)
- Rewind XGBoost version to before 1.0.0 to diagnose test failures for that version [#402](#)
- Update dependency check to use a whitelist [#417](#)
- Update unit test jobs to not install dev deps [#455](#)

Warning: Breaking Changes

- Python 3.5 will not be actively supported.

v0.6.0 Dec. 16, 2019

- **Enhancements**

- Added ability to create a plot of feature importances [#133](#)
- Add early stopping to AutoML using patience and tolerance parameters [#241](#)
- Added ROC and confusion matrix metrics and plot for classification problems and introduce PipelineSearchPlots class [#242](#)
- Enhanced AutoML results with search order [#260](#)

- **Fixes**

- Lower botocore requirement [#235](#)
- Fixed decision_function calculation for FraudCost objective [#254](#)
- Fixed return value of Recall metrics [#264](#)
- Components return *self* on fit [#289](#)

- **Changes**

- Renamed automl classes to AutoRegressionSearch and AutoClassificationSearch [#287](#)
- Updating demo datasets to retain column names [#223](#)
- Moving pipeline visualization to PipelinePlots class [#228](#)
- Standarizing inputs as pd.DataFrame / pd.Series [#130](#)
- Enforcing that pipelines must have an estimator as last component [#277](#)
- Added ipywidgets as a dependency in requirements.txt [#278](#)
- Added Random and Grid Search Tuners [#240](#)

- **Documentation Changes**

- Adding class properties to API reference [#244](#)
- Fix and filter FutureWarnings from scikit-learn [#249](#), [#257](#)
- Adding Linear Regression to API reference and cleaning up some Sphinx warnings [#227](#)

- **Testing Changes**

- Added support for testing on Windows with CircleCI [#226](#)
- Added support for doctests [#233](#)

Warning: Breaking Changes

- The `fit()` method for `AutoClassifier` and `AutoRegressor` has been renamed to `search()`.
- `AutoClassifier` has been renamed to `AutoClassificationSearch`
- `AutoRegressor` has been renamed to `AutoRegressionSearch`
- `AutoClassificationSearch.results` and `AutoRegressionSearch.results` now is a dictionary with `pipeline_results` and `search_order` keys. `pipeline_results` can be used to access a dictionary that is identical to the old `.results` dictionary. Whereas, “`search_order`” returns a list of the search order in terms of pipeline id.
- Pipelines now require an estimator as the last component in `component_list`. Slicing pipelines now throws an `NotImplementedError` to avoid returning Pipelines without an estimator.

v0.5.2 Nov. 18, 2019

- **Enhancements**
 - Adding basic pipeline structure visualization [#211](#)
- **Documentation Changes**
 - Added notebooks to build process [#212](#)

v0.5.1 Nov. 15, 2019

- **Enhancements**
 - Added basic outlier detection guardrail [#151](#)
 - Added basic ID column guardrail [#135](#)
 - Added support for unlimited pipelines with a `max_time` limit [#70](#)
 - Updated `.readthedocs.yaml` to successfully build [#188](#)
- **Fixes**
 - Removed MSLE from default additional objectives [#203](#)
 - Fixed `random_state` passed in pipelines [#204](#)
 - Fixed slow down in `RFRegressor` [#206](#)
- **Changes**
 - Pulled information for `describe_pipeline` from pipeline’s new `describe` method [#190](#)
 - Refactored pipelines [#108](#)
 - Removed guardrails from `Auto(*)` [#202](#), [#208](#)
- **Documentation Changes**
 - Updated documentation to show `max_time` enhancements [#189](#)
 - Updated release instructions for RTD [#193](#)
 - Added notebooks to build process [#212](#)
 - Added contributing instructions [#213](#)

- Added new content #222

v0.5.0 Oct. 29, 2019

- **Enhancements**

- Added basic one hot encoding #73
- Use enums for model_type #110
- Support for splitting regression datasets #112
- Auto-infer multiclass classification #99
- Added support for other units in max_time #125
- Detect highly null columns #121
- Added additional regression objectives #100
- Show an interactive iteration vs. score plot when using fit() #134

- **Fixes**

- Reordered *describe_pipeline* #94
- Added type check for model_type #109
- Fixed *s* units when setting string max_time #132
- Fix objectives not appearing in API documentation #150

- **Changes**

- Reorganized tests #93
- Moved logging to its own module #119
- Show progress bar history #111
- Using cloudpickle instead of pickle to allow unloading of custom objectives #113
- Removed render.py #154

- **Documentation Changes**

- Update release instructions #140
- Include additional_objectives parameter #124
- Added Changelog #136

- **Testing Changes**

- Code coverage #90
- Added CircleCI tests for other Python versions #104
- Added doc notebooks as tests #139
- Test metadata for CircleCI and 2 core parallelism #137

v0.4.1 Sep. 16, 2019

- **Enhancements**

- Added AutoML for classification and regressor using Autobase and Skopt #7 #9
- Implemented standard classification and regression metrics #7
- Added logistic regression, random forest, and XGBoost pipelines #7

- Implemented support for custom objectives #15
- Feature importance for pipelines #18
- Serialization for pipelines #19
- Allow fitting on objectives for optimal threshold #27
- Added detect label leakage #31
- Implemented callbacks #42
- Allow for multiclass classification #21
- Added support for additional objectives #79
- **Fixes**
 - Fixed feature selection in pipelines #13
 - Made random_seed usage consistent #45
- **Documentation Changes**
 - Documentation Changes
 - Added docstrings #6
 - Created notebooks for docs #6
 - Initialized readthedocs EvalML #6
 - Added favicon #38
- **Testing Changes**
 - Added testing for loading data #39

v0.2.0 Aug. 13, 2019

- **Enhancements**
 - Created fraud detection objective #4

v0.1.0 July. 31, 2019

- *First Release*
- **Enhancements**
 - Added lead scoring objective #1
 - Added basic classifier #1
- **Documentation Changes**
 - Initialized Sphinx for docs #1

2.6.14 API Reference

Demo Datasets

<code>load_fraud</code>	Load credit card fraud dataset.
<code>load_wine</code>	Load wine dataset.
<code>load_breast_cancer</code>	Load breast cancer dataset.

Continued on next page

Table 1 – continued from previous page

<code>load_diabetes</code>	Load diabetes dataset.
----------------------------	------------------------

evalml.demos.load_fraud

`evalml.demos.load_fraud(n_rows=None)`

Load credit card fraud dataset. The fraud dataset can be used for binary classification problems.

Parameters `n_rows` (*int*) – number of rows from the dataset to return

Returns `X, y`

Return type `pd.DataFrame, pd.Series`

evalml.demos.load_wine

`evalml.demos.load_wine()`

Load wine dataset. Multiclass problem

Returns `X, y`

Return type `pd.DataFrame, pd.Series`

evalml.demos.load_breast_cancer

`evalml.demos.load_breast_cancer()`

Load breast cancer dataset. Multiclass problem

Returns `X, y`

Return type `pd.DataFrame, pd.Series`

evalml.demos.load_diabetes

`evalml.demos.load_diabetes()`

Load diabetes dataset. Regression problem

Returns `X, y`

Return type `pd.DataFrame, pd.Series`

Preprocessing

<code>load_data</code>	Load features and labels from file(s).
<code>split_data</code>	Splits data into train and test sets.

evalml.preprocessing.load_data

`evalml.preprocessing.load_data(path, index, label, n_rows=None, drop=None, verbose=True, **kwargs)`

Load features and labels from file(s).

Parameters

- **path** (*str*) – path to file or a http/ftp/s3 URL
- **index** (*str*) – column for index
- **label** (*str*) – column for labels
- **n_rows** (*int*) – number of rows to return
- **drop** (*list*) – columns to drop
- **verbose** (*bool*) – whether to print information about features and labels

Returns features and labels

Return type pd.DataFrame, pd.Series

evalml.preprocessing.split_data

evalml.preprocessing.**split_data** (*X*, *y*, *regression=False*, *test_size=0.2*, *random_state=None*)
Splits data into train and test sets.

Parameters

- **X** (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]
- **y** (*pd.Series*) – labels of length [n_samples]
- **regression** (*bool*) – if true, do not use stratified split
- **test_size** (*float*) – percent of train set to holdout for testing
- **random_state** (*int*, *np.random.RandomState*) – seed for the random number generator

Returns features and labels each split into train and test sets

Return type pd.DataFrame, pd.DataFrame, pd.Series, pd.Series

AutoML

<i>AutoClassificationSearch</i>	Automatic pipeline search class for classification problems
<i>AutoRegressionSearch</i>	Automatic pipeline search for regression problems

evalml.AutoClassificationSearch

```
class evalml.AutoClassificationSearch (objective=None, multiclass=False,  
                                       max_pipelines=None, max_time=None,  
                                       patience=None, tolerance=None, allowed_model_families=None, cv=None,  
                                       tuner=None, detect_label_leakage=True,  
                                       start_iteration_callback=None,  
                                       add_result_callback=None, additional_objectives=None, random_state=0, n_jobs=-1,  
                                       verbose=True)
```

Automatic pipeline search class for classification problems

Methods

<code>__init__</code>	Automated classifier pipeline search
-----------------------	--------------------------------------

evalml.AutoClassificationSearch.__init__

```
AutoClassificationSearch.__init__(objective=None, multiclass=False,
                                  max_pipelines=None, max_time=None,
                                  patience=None, tolerance=None, al-
                                  lowed_model_families=None, cv=None,
                                  tuner=None, detect_label_leakage=True,
                                  start_iteration_callback=None,
                                  add_result_callback=None, addi-
                                  tional_objectives=None, random_state=0, n_jobs=-1,
                                  verbose=True)
```

Automated classifier pipeline search

Parameters

- **objective** (*Object*) – the objective to optimize
- **multiclass** (*bool*) – If True, expecting multiclass data. By default: False.
- **max_pipelines** (*int*) – Maximum number of pipelines to search. If max_pipelines and max_time is not set, then max_pipelines will default to max_pipelines of 5.
- **max_time** (*int, str*) – Maximum time to search for pipelines. This will not start a new pipeline search after the duration has elapsed. If it is an integer, then the time will be in seconds. For strings, time can be specified as seconds, minutes, or hours.
- **patience** (*int*) – Number of iterations without improvement to stop search early. Must be positive. If None, early stopping is disabled. Defaults to None.
- **tolerance** (*float*) – Minimum percentage difference to qualify as score improvement for early stopping. Only applicable if patience is not None. Defaults to None.
- **allowed_model_families** (*list*) – The model families to search. By default searches over all model families. Run `evalml.list_model_families("classification")` to see options.
- **cv** – cross validation method to use. By default StratifiedKfold
- **tuner** – the tuner class to use. Defaults to scikit-optimize tuner
- **detect_label_leakage** (*bool*) – If True, check input features for label leakage and warn if found. Defaults to true.
- **start_iteration_callback** (*callable*) – function called before each pipeline training iteration. Passed two parameters: pipeline_class, parameters.
- **add_result_callback** (*callable*) – function called after each pipeline training iteration. Passed two parameters: results, trained_pipeline.
- **additional_objectives** (*list*) – Custom set of objectives to score on. Will override default objectives for problem type if not empty.
- **random_state** (*int, np.random.RandomState*) – The random seed/state. Defaults to 0.

- **n_jobs** (*int or None*) – Non-negative integer describing level of parallelism used for pipelines. None and 1 are equivalent. If set to -1, all CPUs are used. For n_jobs below -1, (n_cpus + 1 + n_jobs) are used.
- **verbose** (*boolean*) – If True, turn verbosity on. Defaults to True

Attributes

<code>best_pipeline</code>	Returns the best model found
<code>rankings</code>	Returns the rankings of the models searched

evalml.AutoRegressionSearch

```
class evalml.AutoRegressionSearch(object=None, max_pipelines=None,
                                     max_time=None, patience=None, tolerance=None, al-
                                     lowed_model_families=None, cv=None, tuner=None, de-
                                     tect_label_leakage=True, start_iteration_callback=None,
                                     add_result_callback=None, additional_objectives=None,
                                     random_state=0, n_jobs=-1, verbose=True)
```

Automatic pipeline search for regression problems

Methods

<code>__init__</code>	Automated regressors pipeline search
-----------------------	--------------------------------------

evalml.AutoRegressionSearch.__init__

```
AutoRegressionSearch.__init__(object=None, max_pipelines=None,
                               max_time=None, patience=None, tolerance=None, al-
                               lowed_model_families=None, cv=None, tuner=None, de-
                               tect_label_leakage=True, start_iteration_callback=None,
                               add_result_callback=None, additional_objectives=None,
                               random_state=0, n_jobs=-1, verbose=True)
```

Automated regressors pipeline search

Parameters

- **objective** (*Object*) – the objective to optimize
- **max_pipelines** (*int*) – Maximum number of pipelines to search. If max_pipelines and max_time is not set, then max_pipelines will default to max_pipelines of 5.
- **max_time** (*int, str*) – Maximum time to search for pipelines. This will not start a new pipeline search after the duration has elapsed. If it is an integer, then the time will be in seconds. For strings, time can be specified as seconds, minutes, or hours.
- **allowed_model_families** (*list*) – The model families to search. By default searches over all model families. Run `evalml.list_model_families("regression")` to see options.
- **patience** (*int*) – Number of iterations without improvement to stop search early. Must be positive. If None, early stopping is disabled. Defaults to None.

- **tolerance** (*float*) – Minimum percentage difference to qualify as score improvement for early stopping. Only applicable if patience is not None. Defaults to None.
- **cv** – cross validation method to use. By default StratifiedKFold
- **tuner** – the tuner class to use. Defaults to scikit-optimize tuner
- **detect_label_leakage** (*bool*) – If True, check input features for label leakage and warn if found. Defaults to true.
- **start_iteration_callback** (*callable*) – function called before each pipeline training iteration. Passed two parameters: pipeline_class, parameters.
- **add_result_callback** (*callable*) – function called after each pipeline training iteration. Passed two parameters: results, trained_pipeline.
- **additional_objectives** (*list*) – Custom set of objectives to score on. Will override default objectives for problem type if not empty.
- **random_state** (*int, np.random.RandomState*) – The random seed/state. Defaults to 0.
- **n_jobs** (*int or None*) – Non-negative integer describing level of parallelism used for pipelines. None and 1 are equivalent. If set to -1, all CPUs are used. For n_jobs below -1, (n_cpus + 1 + n_jobs) are used.
- **verbose** (*boolean*) – If True, turn verbosity on. Defaults to True

Attributes

best_pipeline	Returns the best model found
rankings	Returns the rankings of the models searched

Plotting

<code>AutoClassificationSearch.plot.get_roc_data</code>	Gets data that can be used to create a ROC plot.
<code>AutoClassificationSearch.plot.generate_roc_plot</code>	Generate Receiver Operating Characteristic (ROC) plot for a given pipeline using cross-validation using the data returned from <code>get_roc_data()</code> .
<code>AutoClassificationSearch.plot.get_confusion_matrix_data</code>	Gets data that can be used to create a confusion matrix plot.
<code>AutoClassificationSearch.plot.generate_confusion_matrix</code>	Generate confusion matrix plot for a given pipeline using the data returned from <code>get_confusion_matrix_data()</code> .
<code>AutoClassificationSearch.plot.generate_confusion_matrix</code>	Generate confusion matrix plot for a given pipeline using the data returned from <code>get_confusion_matrix_data()</code> .

evalml.AutoClassificationSearch.plot.get_roc_data

`AutoClassificationSearch.plot.get_roc_data(pipeline_id)`

Gets data that can be used to create a ROC plot.

Returns Dictionary containing metrics used to generate an ROC plot.

evalml.AutoClassificationSearch.plot.generate_roc_plot

AutoClassificationSearch.plot.**generate_roc_plot** (*pipeline_id*)

Generate Receiver Operating Characteristic (ROC) plot for a given pipeline using cross-validation using the data returned from get_roc_data().

Returns plotly.Figure representing the ROC plot generated

evalml.AutoClassificationSearch.plot.get_confusion_matrix_data

AutoClassificationSearch.plot.**get_confusion_matrix_data** (*pipeline_id*, *normalize=None*)

Gets data that can be used to create a confusion matrix plot.

Parameters

- **pipeline_id** (*int*) – ID of pipeline to get confusion matrix data for
- **normalize** ({*'true'*, *'pred'*, *'all'*, *None*}) – Option to normalize over the rows (*'true'*), columns (*'pred'*) or all (*'all'*) values. If option is *None*, returns original confusion matrix. Defaults to *'true'*.

Returns List containing information used to generate a confusion matrix plot. Each element in the list contains the confusion matrix data for that fold.

evalml.AutoClassificationSearch.plot.generate_confusion_matrix

AutoClassificationSearch.plot.**generate_confusion_matrix** (*pipeline_id*, *fold_num=None*, *normalize=None*)

Generate confusion matrix plot for a given pipeline using the data returned from get_confusion_matrix_data().

Parameters

- **pipeline_id** (*int*) – ID of pipeline to get confusion matrix data for
- **fold_num** (*int*) – Fold number of pipeline to get confusion matrix data for
- **option** ({*'true'*, *'pred'*, *'all'*, *None*}) – Option to normalize over the rows (*'true'*), columns (*'pred'*) or all (*'all'*) values. If option is *None*, returns original confusion matrix. Defaults to *'true'*.

Returns plotly.Figure representing the confusion matrix plot generated

Model Family

ModelFamily

Enum for family of machine learning models.

evalml.model_family.ModelFamily

class evalml.model_family.**ModelFamily**

Enum for family of machine learning models.

Components

Transformers

<i>OneHotEncoder</i>	One-hot encoder to encode non-numeric data
<i>RFRegressorSelectFromModel</i>	Selects top features based on importance weights using a Random Forest regressor
<i>RFClassifierSelectFromModel</i>	Selects top features based on importance weights using a Random Forest classifier
<i>SimpleImputer</i>	Imputes missing data according to a specified imputation strategy
<i>StandardScaler</i>	Standardize features: removes mean and scales to unit variance

evalml.pipelines.components.OneHotEncoder

class evalml.pipelines.components.**OneHotEncoder** (*top_n=10, random_state=0*)

One-hot encoder to encode non-numeric data

name = 'One Hot Encoder'

model_family = 'none'

hyperparameter_ranges = {}

Instance attributes

Methods:

<i>__init__</i>	Initializes self.
<i>describe</i>	Describe a component and its parameters
<i>fit</i>	Fits component to data
<i>fit_transform</i>	Fits on X and transforms X
<i>get_feature_names</i>	Returns names of transformed and added columns
<i>transform</i>	One-hot encode the input DataFrame.

evalml.pipelines.components.OneHotEncoder.__init__

OneHotEncoder.**__init__** (*top_n=10, random_state=0*)

Initializes self.

evalml.pipelines.components.OneHotEncoder.describe

OneHotEncoder.**describe** (*print_name=False, return_dict=False*)

Describe a component and its parameters

Parameters

- **print_name** (*bool, optional*) – whether to print name of component
- **return_dict** (*bool, optional*) – whether to return description as dictionary in the format {"name": name, "parameters": parameters}

Returns prints and returns dictionary

Return type None or dict

`evalml.pipelines.components.OneHotEncoder.fit`

`OneHotEncoder.fit(X, y=None)`

Fits component to data

Parameters

- **X** (*pd.DataFrame* or *np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series*, optional) – the target training labels of length [n_samples]

Returns self

`evalml.pipelines.components.OneHotEncoder.fit_transform`

`OneHotEncoder.fit_transform(X, y=None)`

Fits on X and transforms X

Parameters

- **X** (*pd.DataFrame*) – Data to fit and transform
- **y** (*pd.DataFrame*) – Labels to fit and transform

Returns Transformed X

Return type *pd.DataFrame*

`evalml.pipelines.components.OneHotEncoder.get_feature_names`

`OneHotEncoder.get_feature_names()`

Returns names of transformed and added columns

Returns list of feature names not including dropped features

Return type list

`evalml.pipelines.components.OneHotEncoder.transform`

`OneHotEncoder.transform(X, y=None)`

One-hot encode the input DataFrame.

Parameters

- **X** (*pd.DataFrame*) – Dataframe of features.
- **y** (*pd.Series*) – Ignored.

Returns Transformed dataframe, where each categorical feature has been encoded into numerical columns using one-hot encoding.

evalml.pipelines.components.RFRegressorSelectFromModel

```
class evalml.pipelines.components.RFRegressorSelectFromModel (number_features=None,
                                                                n_estimators=10,
                                                                max_depth=None,
                                                                per-
                                                                cent_features=0.5,
                                                                threshold=-inf,
                                                                n_jobs=-1,  ran-
                                                                dom_state=0)
```

Selects top features based on importance weights using a Random Forest regressor

name = 'RF Regressor Select From Model'

model_family = 'none'

hyperparameter_ranges = {'percent_features': Real(low=0.01, high=1, prior='uniform',

Instance attributes**Methods:**

<code>__init__</code>	Initialize self.
<code>describe</code>	Describe a component and its parameters
<code>fit</code>	Fits component to data
<code>fit_transform</code>	Fits feature selector on data X then transforms X by selecting features
<code>get_indices</code>	Get integer index of features selected
<code>get_names</code>	Get names of selected features.
<code>transform</code>	Transforms data X by selecting features

evalml.pipelines.components.RFRegressorSelectFromModel.__init__

```
RFRegressorSelectFromModel.__init__(number_features=None,      n_estimators=10,
                                     max_depth=None,           percent_features=0.5,
                                     threshold=-inf, n_jobs=-1, random_state=0)
```

Initialize self. See help(type(self)) for accurate signature.

evalml.pipelines.components.RFRegressorSelectFromModel.describe

```
RFRegressorSelectFromModel.describe(print_name=False, return_dict=False)
```

Describe a component and its parameters

Parameters

- **print_name** (*bool, optional*) – whether to print name of component
- **return_dict** (*bool, optional*) – whether to return description as dictionary in the format {"name": name, "parameters": parameters}

Returns prints and returns dictionary

Return type None or dict

evalml.pipelines.components.RFRegressorSelectFromModel.fit

`RFRegressorSelectFromModel.fit(X, y=None)`

Fits component to data

Parameters

- **X** (*pd.DataFrame* or *np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series*, *optional*) – the target training labels of length [n_samples]

Returns self

evalml.pipelines.components.RFRegressorSelectFromModel.fit_transform

`RFRegressorSelectFromModel.fit_transform(X, y=None)`

Fits feature selector on data X then transforms X by selecting features

Parameters

- **X** (*pd.DataFrame*) – Data to fit and transform
- **y** (*pd.Series*) – Labels to fit and transform

Returns Transformed X

Return type *pd.DataFrame*

evalml.pipelines.components.RFRegressorSelectFromModel.get_indices

`RFRegressorSelectFromModel.get_indices()`

Get integer index of features selected

Returns list of indices

Return type list

evalml.pipelines.components.RFRegressorSelectFromModel.get_names

`RFRegressorSelectFromModel.get_names()`

Get names of selected features.

Returns list of the names of features selected

evalml.pipelines.components.RFRegressorSelectFromModel.transform

`RFRegressorSelectFromModel.transform(X, y=None)`

Transforms data X by selecting features

Parameters

- **X** (*pd.DataFrame*) – Data to transform
- **y** (*pd.Series*, *optional*) – Input Labels

Returns Transformed X

Return type `pd.DataFrame`

`evalml.pipelines.components.RFClassifierSelectFromModel`

```
class evalml.pipelines.components.RFClassifierSelectFromModel(number_features=None,
                                                             n_estimators=10,
                                                             max_depth=None,
                                                             per-
                                                             cent_features=0.5,
                                                             threshold=-inf,
                                                             n_jobs=-1, ran-
                                                             dom_state=0)
```

Selects top features based on importance weights using a Random Forest classifier

```
name = 'RF Classifier Select From Model'
```

```
model_family = 'none'
```

```
hyperparameter_ranges = {'percent_features': Real(low=0.01, high=1, prior='uniform',
```

Instance attributes

Methods:

<code>__init__</code>	Initialize self.
<code>describe</code>	Describe a component and its parameters
<code>fit</code>	Fits component to data
<code>fit_transform</code>	Fits feature selector on data X then transforms X by selecting features
<code>get_indices</code>	Get integer index of features selected
<code>get_names</code>	Get names of selected features.
<code>transform</code>	Transforms data X by selecting features

`evalml.pipelines.components.RFClassifierSelectFromModel.__init__`

```
RFClassifierSelectFromModel.__init__(number_features=None,      n_estimators=10,
                                     max_depth=None,          percent_features=0.5,
                                     threshold=-inf, n_jobs=-1, random_state=0)
```

Initialize self. See `help(type(self))` for accurate signature.

`evalml.pipelines.components.RFClassifierSelectFromModel.describe`

```
RFClassifierSelectFromModel.describe(print_name=False, return_dict=False)
```

Describe a component and its parameters

Parameters

- **print_name** (*bool, optional*) – whether to print name of component
- **return_dict** (*bool, optional*) – whether to return description as dictionary in the format {"name": name, "parameters": parameters}

Returns prints and returns dictionary

Return type None or dict

evalml.pipelines.components.RFClassifierSelectFromModel.fit

`RFClassifierSelectFromModel.fit(X, y=None)`

Fits component to data

Parameters

- **X** (*pd.DataFrame* or *np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series*, optional) – the target training labels of length [n_samples]

Returns self

evalml.pipelines.components.RFClassifierSelectFromModel.fit_transform

`RFClassifierSelectFromModel.fit_transform(X, y=None)`

Fits feature selector on data X then transforms X by selecting features

Parameters

- **X** (*pd.DataFrame*) – Data to fit and transform
- **y** (*pd.Series*) – Labels to fit and transform

Returns Transformed X

Return type *pd.DataFrame*

evalml.pipelines.components.RFClassifierSelectFromModel.get_indices

`RFClassifierSelectFromModel.get_indices()`

Get integer index of features selected

Returns list of indices

Return type list

evalml.pipelines.components.RFClassifierSelectFromModel.get_names

`RFClassifierSelectFromModel.get_names()`

Get names of selected features.

Returns list of the names of features selected

evalml.pipelines.components.RFClassifierSelectFromModel.transform

`RFClassifierSelectFromModel.transform(X, y=None)`

Transforms data X by selecting features

Parameters

- **x** (*pd.DataFrame*) – Data to transform
- **y** (*pd.Series*, *optional*) – Input Labels

Returns Transformed X

Return type *pd.DataFrame*

evalml.pipelines.components.SimpleImputer

```
class evalml.pipelines.components.SimpleImputer (impute_strategy='most_frequent',  
                                                fill_value=None, random_state=0)  
    Imputes missing data according to a specified imputation strategy  
  
    name = 'Simple Imputer'  
  
    model_family = 'none'  
  
    hyperparameter_ranges = {'impute_strategy': ['mean', 'median', 'most_frequent']}
```

Instance attributes

Methods:

<code>__init__</code>	Initializes an transformer that imputes missing data according to the specified imputation strategy.”
<code>describe</code>	Describe a component and its parameters
<code>fit</code>	Fits component to data
<code>fit_transform</code>	Fits imputer on data X then imputes missing values in X
<code>transform</code>	Transforms data X by imputing missing values

evalml.pipelines.components.SimpleImputer.__init__

`SimpleImputer.__init__` (*impute_strategy='most_frequent', fill_value=None, random_state=0*)
Initializes an transformer that imputes missing data according to the specified imputation strategy.”

Parameters

- **impute_strategy** (*string*) – Impute strategy to use. Valid values include “mean”, “median”, “most_frequent”, “constant” for numerical data, and “most_frequent”, “constant” for object data types.
- **fill_value** (*string*) – When `impute_strategy == “constant”`, `fill_value` is used to replace missing data. Defaults to 0 when imputing numerical data and “missing_value” for strings or object data types.

evalml.pipelines.components.SimpleImputer.describe

`SimpleImputer.describe` (*print_name=False, return_dict=False*)
Describe a component and its parameters

Parameters

- **print_name** (*bool, optional*) – whether to print name of component
- **return_dict** (*bool, optional*) – whether to return description as dictionary in the format {"name": name, "parameters": parameters}

Returns prints and returns dictionary

Return type None or dict

evalml.pipelines.components.SimpleImputer.fit

`SimpleImputer.fit(X, y=None)`

Fits component to data

Parameters

- **X** (*pd.DataFrame or np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series, optional*) – the target training labels of length [n_samples]

Returns self

evalml.pipelines.components.SimpleImputer.fit_transform

`SimpleImputer.fit_transform(X, y=None)`

Fits imputer on data X then imputes missing values in X

Parameters

- **X** (*pd.DataFrame*) – Data to fit and transform
- **y** (*pd.Series*) – Labels to fit and transform

Returns Transformed X

Return type pd.DataFrame

evalml.pipelines.components.SimpleImputer.transform

`SimpleImputer.transform(X, y=None)`

Transforms data X by imputing missing values

Parameters

- **X** (*pd.DataFrame*) – Data to transform
- **y** (*pd.Series, optional*) – Input Labels

Returns Transformed X

Return type pd.DataFrame

evalml.pipelines.components.StandardScaler

class evalml.pipelines.components.**StandardScaler** (*random_state=0*)

Standardize features: removes mean and scales to unit variance

name = 'Standard Scaler'

```
model_family = 'none'
hyperparameter_ranges = {}
```

Instance attributes

Methods:

<code>__init__</code>	Initialize self.
<code>describe</code>	Describe a component and its parameters
<code>fit</code>	Fits component to data
<code>fit_transform</code>	Fits on X and transforms X
<code>transform</code>	Transforms data X

`evalml.pipelines.components.StandardScaler.__init__`

`StandardScaler.__init__(random_state=0)`
Initialize self. See `help(type(self))` for accurate signature.

`evalml.pipelines.components.StandardScaler.describe`

`StandardScaler.describe(print_name=False, return_dict=False)`
Describe a component and its parameters

Parameters

- **print_name** (*bool, optional*) – whether to print name of component
- **return_dict** (*bool, optional*) – whether to return description as dictionary in the format {"name": name, "parameters": parameters}

Returns prints and returns dictionary

Return type None or dict

`evalml.pipelines.components.StandardScaler.fit`

`StandardScaler.fit(X, y=None)`
Fits component to data

Parameters

- **X** (*pd.DataFrame or np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series, optional*) – the target training labels of length [n_samples]

Returns self

evalml.pipelines.components.StandardScaler.fit_transform

`StandardScaler.fit_transform(X, y=None)`

Fits on X and transforms X

Parameters

- **X** (`pd.DataFrame`) – Data to fit and transform
- **y** (`pd.DataFrame`) – Labels to fit and transform

Returns Transformed X

Return type `pd.DataFrame`

evalml.pipelines.components.StandardScaler.transform

`StandardScaler.transform(X, y=None)`

Transforms data X

Parameters

- **X** (`pd.DataFrame`) – Data to transform
- **y** (`pd.Series, optional`) – Input Labels

Returns Transformed X

Return type `pd.DataFrame`

Estimators

<i>LogisticRegressionClassifier</i>	Logistic Regression Classifier
<i>RandomForestClassifier</i>	Random Forest Classifier
<i>XGBoostClassifier</i>	XGBoost Classifier
<i>LinearRegressor</i>	Linear Regressor
<i>RandomForestRegressor</i>	Random Forest Regressor

evalml.pipelines.components.LogisticRegressionClassifier

```
class evalml.pipelines.components.LogisticRegressionClassifier (penalty='l2',  
                                                                C=1.0,  
                                                                n_jobs=-1, ran-  
                                                                dom_state=0)
```

Logistic Regression Classifier

name = 'Logistic Regression Classifier'

model_family = 'linear_model'

supported_problem_types = [`<ProblemTypes.BINARY: 'binary'>`, `<ProblemTypes.MULTICLASS:`

`hyperparameter_ranges` = {'C': `Real(low=0.01, high=10, prior='uniform', transform='iden`

Instance attributes

<code>feature_importances</code>	Returns feature importances.
----------------------------------	------------------------------

Methods:

<code>__init__</code>	Initialize self.
<code>describe</code>	Describe a component and its parameters
<code>fit</code>	Fits component to data
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.

evalml.pipelines.components.LogisticRegressionClassifier.__init__

`LogisticRegressionClassifier.__init__(penalty='l2', C=1.0, n_jobs=-1, random_state=0)`
 Initialize self. See `help(type(self))` for accurate signature.

evalml.pipelines.components.LogisticRegressionClassifier.describe

`LogisticRegressionClassifier.describe(print_name=False, return_dict=False)`
 Describe a component and its parameters

Parameters

- **print_name** (*bool, optional*) – whether to print name of component
- **return_dict** (*bool, optional*) – whether to return description as dictionary in the format {"name": name, "parameters": parameters}

Returns prints and returns dictionary

Return type None or dict

evalml.pipelines.components.LogisticRegressionClassifier.fit

`LogisticRegressionClassifier.fit(X, y=None)`
 Fits component to data

Parameters

- **X** (*pd.DataFrame or np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series, optional*) – the target training labels of length [n_samples]

Returns self

evalml.pipelines.components.LogisticRegressionClassifier.predict

`LogisticRegressionClassifier.predict(X)`
 Make predictions using selected features.

Parameters **X** (*pd.DataFrame*) – features

Returns estimated labels

Return type `pd.Series`

`evalml.pipelines.components.LogisticRegressionClassifier.predict_proba`

`LogisticRegressionClassifier.predict_proba(X)`

Make probability estimates for labels.

Parameters `X` (`pd.DataFrame`) – features

Returns probability estimates

Return type `pd.DataFrame`

`evalml.pipelines.components.RandomForestClassifier`

```
class evalml.pipelines.components.RandomForestClassifier(n_estimators=10,
                                                         max_depth=None,
                                                         n_jobs=-1,           ran-
                                                         dom_state=0)
```

Random Forest Classifier

`name = 'Random Forest Classifier'`

`model_family = 'random_forest'`

`supported_problem_types = [<ProblemTypes.BINARY: 'binary'>, <ProblemTypes.MULTICLASS:`

`hyperparameter_ranges = {'max_depth': Integer(low=1, high=32, prior='uniform', transf`

Instance attributes

<code>feature_importances</code>	Returns feature importances.
----------------------------------	------------------------------

Methods:

<code>__init__</code>	Initialize self.
<code>describe</code>	Describe a component and its parameters
<code>fit</code>	Fits component to data
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.

`evalml.pipelines.components.RandomForestClassifier.__init__`

`RandomForestClassifier.__init__(n_estimators=10, max_depth=None, n_jobs=-1, ran-`
`dom_state=0`)

Initialize self. See `help(type(self))` for accurate signature.

`evalml.pipelines.components.RandomForestClassifier.describe`

`RandomForestClassifier.describe(print_name=False, return_dict=False)`

Describe a component and its parameters

Parameters

- **print_name** (*bool, optional*) – whether to print name of component
- **return_dict** (*bool, optional*) – whether to return description as dictionary in the format {"name": name, "parameters": parameters}

Returns prints and returns dictionary

Return type None or dict

evalml.pipelines.components.RandomForestClassifier.fit

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

Fits component to data

Parameters

- **X** (*pd.DataFrame or np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series, optional*) – the target training labels of length [n_samples]

Returns self

evalml.pipelines.components.RandomForestClassifier.predict

RandomForestClassifier.**predict** (*X*)

Make predictions using selected features.

Parameters **X** (*pd.DataFrame*) – features

Returns estimated labels

Return type pd.Series

evalml.pipelines.components.RandomForestClassifier.predict_proba

RandomForestClassifier.**predict_proba** (*X*)

Make probability estimates for labels.

Parameters **X** (*pd.DataFrame*) – features

Returns probability estimates

Return type pd.DataFrame

evalml.pipelines.components.XGBoostClassifier

```
class evalml.pipelines.components.XGBoostClassifier(eta=0.1, max_depth=3,
                                                    min_child_weight=1,
                                                    n_estimators=100, ran-
                                                    dom_state=0)
```

XGBoost Classifier

name = 'XGBoost Classifier'

model_family = 'xgboost'

```
supported_problem_types = [<ProblemTypes.BINARY: 'binary'>, <ProblemTypes.MULTICLASS:
hyperparameter_ranges = {'eta': Real(low=0, high=1, prior='uniform', transform='ident
```

Instance attributes

<code>feature_importances</code>	Returns feature importances.
----------------------------------	------------------------------

Methods:

<code>__init__</code>	Initialize self.
<code>describe</code>	Describe a component and its parameters
<code>fit</code>	Fits component to data
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.

`evalml.pipelines.components.XGBoostClassifier.__init__`

`XGBoostClassifier.__init__` (*eta=0.1, max_depth=3, min_child_weight=1, n_estimators=100, random_state=0*)
Initialize self. See `help(type(self))` for accurate signature.

`evalml.pipelines.components.XGBoostClassifier.describe`

`XGBoostClassifier.describe` (*print_name=False, return_dict=False*)
Describe a component and its parameters

Parameters

- **print_name** (*bool, optional*) – whether to print name of component
- **return_dict** (*bool, optional*) – whether to return description as dictionary in the format {"name": name, "parameters": parameters}

Returns prints and returns dictionary

Return type None or dict

`evalml.pipelines.components.XGBoostClassifier.fit`

`XGBoostClassifier.fit` (*X, y=None*)
Fits component to data

Parameters

- **X** (*pd.DataFrame or np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series, optional*) – the target training labels of length [n_samples]

Returns self

evalml.pipelines.components.XGBoostClassifier.predict`XGBoostClassifier.predict(X)`

Make predictions using selected features.

Parameters *X* (`pd.DataFrame`) – features**Returns** estimated labels**Return type** `pd.Series`**evalml.pipelines.components.XGBoostClassifier.predict_proba**`XGBoostClassifier.predict_proba(X)`

Make probability estimates for labels.

Parameters *X* (`pd.DataFrame`) – features**Returns** probability estimates**Return type** `pd.DataFrame`**evalml.pipelines.components.LinearRegressor**

```
class evalml.pipelines.components.LinearRegressor (fit_intercept=True,          normal-
                                                    ize=False,    n_jobs=-1,    ran-
                                                    dom_state=0)
```

Linear Regressor

name = 'Linear Regressor'**model_family** = 'linear_model'**supported_problem_types** = [`<ProblemTypes.REGRESSION: 'regression'>`]**hyperparameter_ranges** = {'fit_intercept': [True, False], 'normalize': [True, False]}**Instance attributes**

<code>feature_importances</code>	Returns feature importances.
----------------------------------	------------------------------

Methods:

<code>__init__</code>	Initialize self.
<code>describe</code>	Describe a component and its parameters
<code>fit</code>	Fits component to data
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.

evalml.pipelines.components.LinearRegressor.__init__`LinearRegressor.__init__(fit_intercept=True, normalize=False, n_jobs=-1, random_state=0)`Initialize self. See `help(type(self))` for accurate signature.

evalml.pipelines.components.LinearRegressor.describe

`LinearRegressor.describe` (*print_name=False, return_dict=False*)

Describe a component and its parameters

Parameters

- **print_name** (*bool, optional*) – whether to print name of component
- **return_dict** (*bool, optional*) – whether to return description as dictionary in the format {"name": name, "parameters": parameters}

Returns prints and returns dictionary

Return type None or dict

evalml.pipelines.components.LinearRegressor.fit

`LinearRegressor.fit` (*X, y=None*)

Fits component to data

Parameters

- **X** (*pd.DataFrame or np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series, optional*) – the target training labels of length [n_samples]

Returns self

evalml.pipelines.components.LinearRegressor.predict

`LinearRegressor.predict` (*X*)

Make predictions using selected features.

Parameters **X** (*pd.DataFrame*) – features

Returns estimated labels

Return type pd.Series

evalml.pipelines.components.LinearRegressor.predict_proba

`LinearRegressor.predict_proba` (*X*)

Make probability estimates for labels.

Parameters **X** (*pd.DataFrame*) – features

Returns probability estimates

Return type pd.DataFrame

evalml.pipelines.components.RandomForestRegressor

```
class evalml.pipelines.components.RandomForestRegressor (n_estimators=10,
                                                         max_depth=None,
                                                         n_jobs=-1,          ran-
                                                         dom_state=0)

Random Forest Regressor

name = 'Random Forest Regressor'
model_family = 'random_forest'
supported_problem_types = [<ProblemTypes.REGRESSION: 'regression'>]
hyperparameter_ranges = {'max_depth': Integer(low=1, high=32, prior='uniform', transf
```

Instance attributes

feature_importances	Returns feature importances.
---------------------	------------------------------

Methods:

<code>__init__</code>	Initialize self.
<code>describe</code>	Describe a component and its parameters
<code>fit</code>	Fits component to data
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.

evalml.pipelines.components.RandomForestRegressor.__init__

RandomForestRegressor.__init__(n_estimators=10, max_depth=None, n_jobs=-1, random_state=0)

Initialize self. See help(type(self)) for accurate signature.

evalml.pipelines.components.RandomForestRegressor.describe

RandomForestRegressor.describe(print_name=False, return_dict=False)

Describe a component and its parameters

Parameters

- **print_name** (*bool, optional*) – whether to print name of component
- **return_dict** (*bool, optional*) – whether to return description as dictionary in the format {"name": name, "parameters": parameters}

Returns prints and returns dictionary

Return type None or dict

evalml.pipelines.components.RandomForestRegressor.fit

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

Fits component to data

Parameters

- **X** (*pd.DataFrame* or *np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series*, optional) – the target training labels of length [n_samples]

Returns self

evalml.pipelines.components.RandomForestRegressor.predict

RandomForestRegressor.**predict** (*X*)

Make predictions using selected features.

Parameters **X** (*pd.DataFrame*) – features

Returns estimated labels

Return type pd.Series

evalml.pipelines.components.RandomForestRegressor.predict_proba

RandomForestRegressor.**predict_proba** (*X*)

Make probability estimates for labels.

Parameters **X** (*pd.DataFrame*) – features

Returns probability estimates

Return type pd.DataFrame

Pipelines

Pipelines

PipelineBase

Base class for all pipelines.

evalml.pipelines.PipelineBase

class evalml.pipelines.**PipelineBase** (*parameters*, *objective*, *random_state=0*)

Base class for all pipelines.

Methods

`__init__`

Machine learning pipeline made out of transformers and a estimator.

Continued on next page

Table 33 – continued from previous page

<code>describe</code>	Outputs pipeline details including component parameters
<code>feature_importance_graph</code>	Generate a bar graph of the pipeline’s feature importances
<code>fit</code>	Build a model
<code>get_component</code>	Returns component by name
<code>graph</code>	Generate an image representing the pipeline graph
<code>load</code>	Loads pipeline at file path
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.
<code>save</code>	Saves pipeline at file path
<code>score</code>	Evaluate model performance on current and additional objectives

evalml.pipelines.PipelineBase.__init__

`PipelineBase.__init__(parameters, objective, random_state=0)`

Machine learning pipeline made out of transformers and a estimator.

Required Class Variables: `component_graph` (list): List of components in order. Accepts strings or `ComponentBase` objects in the list

`supported_problem_types` (list): List of problem types for this pipeline. Accepts strings or `ProblemType` enum in the list.

Parameters

- **objective** (*ObjectiveBase*) – the objective to optimize
- **parameters** (*dict*) – dictionary with component names as keys and dictionary of that component’s parameters as values. An empty dictionary `{}` implies using all default values for component parameters.
- **random_state** (*int*, *np.random.RandomState*) – The random seed/state. Defaults to 0.

evalml.pipelines.PipelineBase.describe

`PipelineBase.describe()`

Outputs pipeline details including component parameters

Parameters `return_dict` (*bool*) – If True, return dictionary of information about pipeline. Defaults to false

Returns dictionary of all component parameters if `return_dict` is True, else None

Return type dict

evalml.pipelines.PipelineBase.feature_importance_graph

`PipelineBase.feature_importance_graph(show_all_features=False)`

Generate a bar graph of the pipeline’s feature importances

Parameters `show_all_features` (*bool*, *optional*) – If true, graph features with an importance value of zero. Defaults to false.

Returns `plotly.Figure`, a bar graph showing features and their importances

`evalml.pipelines.PipelineBase.fit`

`PipelineBase.fit` (*X*, *y*, *objective_fit_size=0.2*)

Build a model

Parameters

- **X** (*pd.DataFrame* or *np.array*) – the input training data of shape [*n_samples*, *n_features*]
- **y** (*pd.Series*) – the target training labels of length [*n_samples*]
- **feature_types** (*list*, *optional*) – list of feature types. either numeric or categorical. categorical features will automatically be encoded

Returns `self`

`evalml.pipelines.PipelineBase.get_component`

`PipelineBase.get_component` (*name*)

Returns component by name

Parameters **name** (*str*) – name of component

Returns component to return

Return type `Component`

`evalml.pipelines.PipelineBase.graph`

`PipelineBase.graph` (*filepath=None*)

Generate an image representing the pipeline graph

Parameters **filepath** (*str*, *optional*) – Path to where the graph should be saved. If set to `None` (as by default), the graph will not be saved.

Returns Graph object that can be directly displayed in Jupyter notebooks.

Return type `graphviz.Digraph`

`evalml.pipelines.PipelineBase.load`

static `PipelineBase.load` (*file_path*)

Loads pipeline at file path

Parameters **file_path** (*str*) – location to load file

Returns `PipelineBase` obj

evalml.pipelines.PipelineBase.predict

PipelineBase.**predict** (*X*)

Make predictions using selected features.

Parameters *X* (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]

Returns estimated labels

Return type *pd.Series*

evalml.pipelines.PipelineBase.predict_proba

PipelineBase.**predict_proba** (*X*)

Make probability estimates for labels.

Parameters *X* (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]

Returns probability estimates

Return type *pd.DataFrame*

evalml.pipelines.PipelineBase.save

PipelineBase.**save** (*file_path*)

Saves pipeline at file path

Parameters *file_path* (*str*) – location to save file

Returns None

evalml.pipelines.PipelineBase.score

PipelineBase.**score** (*X*, *y*, *other_objectives=None*)

Evaluate model performance on current and additional objectives

Parameters

- *X* (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]
- *y* (*pd.Series*) – true labels of length [n_samples]
- **other_objectives** (*list*) – list of other objectives to score

Returns score, ordered dictionary of other objective scores

Return type float, dict

<i>RFClassificationPipeline</i>	Random Forest Pipeline for both binary and multiclass classification
<i>XGBoostPipeline</i>	XGBoost Pipeline for both binary and multiclass classification
<i>CatBoostClassificationPipeline</i>	CatBoost Pipeline for both binary and multiclass classification.
<i>LogisticRegressionPipeline</i>	Logistic Regression Pipeline for both binary and multiclass classification

Continued on next page

Table 34 – continued from previous page

<i>RFRegressionPipeline</i>	Random Forest Pipeline for regression problems
<i>CatBoostRegressionPipeline</i>	CatBoost Pipeline for regression problems.
<i>LinearRegressionPipeline</i>	Linear Regression Pipeline for regression problems

evalml.pipelines.RFClassificationPipeline

```
class evalml.pipelines.RFClassificationPipeline(parameters, objective, random_state=0)
    Random Forest Pipeline for both binary and multiclass classification

    name = 'Random Forest Classification Pipeline'
    summary = 'Random Forest Classifier w/ One Hot Encoder + Simple Imputer + RF Classifier'
    component_graph = ['One Hot Encoder', 'Simple Imputer', 'RF Classifier Select From Model']
    supported_problem_types = ['binary', 'multiclass']
    model_family = 'random_forest'
    hyperparameters = {'impute_strategy': ['mean', 'median', 'most_frequent'], 'max_depth': [5, 10, 15, 20, 25, 30, 35, 40, 45, 50]}
    custom_hyperparameters = None
```

Instance attributes

<i>feature_importances</i>	Return feature importances.
<i>parameters</i>	Returns parameter dictionary for this pipeline

Methods:

<i>__init__</i>	Machine learning pipeline made out of transformers and a estimator.
<i>describe</i>	Outputs pipeline details including component parameters
<i>feature_importance_graph</i>	Generate a bar graph of the pipeline's feature importances
<i>fit</i>	Build a model
<i>get_component</i>	Returns component by name
<i>graph</i>	Generate an image representing the pipeline graph
<i>load</i>	Loads pipeline at file path
<i>predict</i>	Make predictions using selected features.
<i>predict_proba</i>	Make probability estimates for labels.
<i>save</i>	Saves pipeline at file path
<i>score</i>	Evaluate model performance on current and additional objectives

evalml.pipelines.RFClassificationPipeline.__init__

```
RFClassificationPipeline.__init__(parameters, objective, random_state=0)
    Machine learning pipeline made out of transformers and a estimator.
```


Required Class Variables: `component_graph` (list): List of components in order. Accepts strings or `ComponentBase` objects in the list

`supported_problem_types` (list): List of problem types for this pipeline. Accepts strings or `ProblemType` enum in the list.

Parameters

- **objective** (*ObjectiveBase*) – the objective to optimize
- **parameters** (*dict*) – dictionary with component names as keys and dictionary of that component's parameters as values. An empty dictionary {} implies using all default values for component parameters.
- **random_state** (*int*, *np.random.RandomState*) – The random seed/state. Defaults to 0.

evalml.pipelines.RFClassificationPipeline.describe

`RFClassificationPipeline.describe()`

Outputs pipeline details including component parameters

Parameters `return_dict` (*bool*) – If True, return dictionary of information about pipeline. Defaults to false

Returns dictionary of all component parameters if `return_dict` is True, else None

Return type dict

evalml.pipelines.RFClassificationPipeline.feature_importance_graph

`RFClassificationPipeline.feature_importance_graph(show_all_features=False)`

Generate a bar graph of the pipeline's feature importances

Parameters `show_all_features` (*bool*, *optional*) – If true, graph features with an importance value of zero. Defaults to false.

Returns `plotly.Figure`, a bar graph showing features and their importances

evalml.pipelines.RFClassificationPipeline.fit

`RFClassificationPipeline.fit(X, y, objective_fit_size=0.2)`

Build a model

Parameters

- **X** (*pd.DataFrame* or *np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series*) – the target training labels of length [n_samples]
- **feature_types** (*list*, *optional*) – list of feature types. either numeric or categorical. categorical features will automatically be encoded

Returns self

evalml.pipelines.RFClassificationPipeline.get_component

`RFClassificationPipeline.get_component(name)`

Returns component by name

Parameters `name` (*str*) – name of component

Returns component to return

Return type Component

evalml.pipelines.RFClassificationPipeline.graph

`RFClassificationPipeline.graph(filepath=None)`

Generate an image representing the pipeline graph

Parameters `filepath` (*str, optional*) – Path to where the graph should be saved. If set to None (as by default), the graph will not be saved.

Returns Graph object that can be directly displayed in Jupyter notebooks.

Return type graphviz.Digraph

evalml.pipelines.RFClassificationPipeline.load

static `RFClassificationPipeline.load(file_path)`

Loads pipeline at file path

Parameters `file_path` (*str*) – location to load file

Returns PipelineBase obj

evalml.pipelines.RFClassificationPipeline.predict

`RFClassificationPipeline.predict(X)`

Make predictions using selected features.

Parameters `X` (*pd.DataFrame or np.array*) – data of shape [n_samples, n_features]

Returns estimated labels

Return type pd.Series

evalml.pipelines.RFClassificationPipeline.predict_proba

`RFClassificationPipeline.predict_proba(X)`

Make probability estimates for labels.

Parameters `X` (*pd.DataFrame or np.array*) – data of shape [n_samples, n_features]

Returns probability estimates

Return type pd.DataFrame

evalml.pipelines.RFClassificationPipeline.save

`RFClassificationPipeline.save(file_path)`

Saves pipeline at file path

Parameters `file_path` (*str*) – location to save file

Returns None

evalml.pipelines.RFClassificationPipeline.score

`RFClassificationPipeline.score(X, y, other_objectives=None)`

Evaluate model performance on current and additional objectives

Parameters

- `X` (*pd.DataFrame* or *np.array*) – data of shape `[n_samples, n_features]`
- `y` (*pd.Series*) – true labels of length `[n_samples]`
- `other_objectives` (*list*) – list of other objectives to score

Returns score, ordered dictionary of other objective scores

Return type float, dict

evalml.pipelines.XGBoostPipeline

class `evalml.pipelines.XGBoostPipeline(parameters, objective, random_state=0)`

XGBoost Pipeline for both binary and multiclass classification

`name = 'XGBoost Classification Pipeline'`

`summary = 'XGBoost Classifier w/ One Hot Encoder + Simple Imputer + RF Classifier Select From Model'`

`component_graph = ['One Hot Encoder', 'Simple Imputer', 'RF Classifier Select From Model']`

`supported_problem_types = [<ProblemTypes.BINARY: 'binary'>, <ProblemTypes.MULTICLASS: 'multiclass'>]`

`model_family = 'xgboost'`

`hyperparameters = {'eta': Real(low=0, high=1, prior='uniform', transform='identity')}`

`custom_hyperparameters = None`

Instance attributes

<code>feature_importances</code>	Return feature importances.
<code>parameters</code>	Returns parameter dictionary for this pipeline

Methods:

<code>__init__</code>	Machine learning pipeline made out of transformers and a estimator.
<code>describe</code>	Outputs pipeline details including component parameters

Continued on next page

Table 38 – continued from previous page

<code>feature_importance_graph</code>	Generate a bar graph of the pipeline’s feature importances
<code>fit</code>	Build a model
<code>get_component</code>	Returns component by name
<code>graph</code>	Generate an image representing the pipeline graph
<code>load</code>	Loads pipeline at file path
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.
<code>save</code>	Saves pipeline at file path
<code>score</code>	Evaluate model performance on current and additional objectives

`evalml.pipelines.XGBoostPipeline.__init__`

`XGBoostPipeline.__init__(parameters, objective, random_state=0)`

Machine learning pipeline made out of transformers and a estimator.

Required Class Variables: `component_graph` (list): List of components in order. Accepts strings or `ComponentBase` objects in the list

`supported_problem_types` (list): List of problem types for this pipeline. Accepts strings or `ProblemType` enum in the list.

Parameters

- **objective** (*ObjectiveBase*) – the objective to optimize
- **parameters** (*dict*) – dictionary with component names as keys and dictionary of that component’s parameters as values. An empty dictionary `{}` implies using all default values for component parameters.
- **random_state** (*int*, *np.random.RandomState*) – The random seed/state. Defaults to 0.

`evalml.pipelines.XGBoostPipeline.describe`

`XGBoostPipeline.describe()`

Outputs pipeline details including component parameters

Parameters `return_dict` (*bool*) – If True, return dictionary of information about pipeline. Defaults to false

Returns dictionary of all component parameters if `return_dict` is True, else None

Return type dict

`evalml.pipelines.XGBoostPipeline.feature_importance_graph`

`XGBoostPipeline.feature_importance_graph(show_all_features=False)`

Generate a bar graph of the pipeline’s feature importances

Parameters `show_all_features` (*bool*, *optional*) – If true, graph features with an importance value of zero. Defaults to false.

Returns `plotly.Figure`, a bar graph showing features and their importances

evalml.pipelines.XGBoostPipeline.fit

XGBoostPipeline.**fit** (*X*, *y*, *objective_fit_size=0.2*)

Build a model

Parameters

- **X** (*pd.DataFrame* or *np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series*) – the target training labels of length [n_samples]
- **feature_types** (*list*, *optional*) – list of feature types. either numeric or categorical. categorical features will automatically be encoded

Returns self

evalml.pipelines.XGBoostPipeline.get_component

XGBoostPipeline.**get_component** (*name*)

Returns component by name

Parameters **name** (*str*) – name of component

Returns component to return

Return type Component

evalml.pipelines.XGBoostPipeline.graph

XGBoostPipeline.**graph** (*filepath=None*)

Generate an image representing the pipeline graph

Parameters **filepath** (*str*, *optional*) – Path to where the graph should be saved. If set to None (as by default), the graph will not be saved.

Returns Graph object that can be directly displayed in Jupyter notebooks.

Return type graphviz.Digraph

evalml.pipelines.XGBoostPipeline.load

static XGBoostPipeline.**load** (*file_path*)

Loads pipeline at file path

Parameters **file_path** (*str*) – location to load file

Returns PipelineBase obj

evalml.pipelines.XGBoostPipeline.predict

XGBoostPipeline.**predict** (*X*)

Make predictions using selected features.

Parameters **X** (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]

Returns estimated labels

Return type `pd.Series`

`evalml.pipelines.XGBoostPipeline.predict_proba`

`XGBoostPipeline.predict_proba(X)`

Make probability estimates for labels.

Parameters `X` (`pd.DataFrame` or `np.array`) – data of shape `[n_samples, n_features]`

Returns probability estimates

Return type `pd.DataFrame`

`evalml.pipelines.XGBoostPipeline.save`

`XGBoostPipeline.save(file_path)`

Saves pipeline at file path

Parameters `file_path` (`str`) – location to save file

Returns `None`

`evalml.pipelines.XGBoostPipeline.score`

`XGBoostPipeline.score(X, y, other_objectives=None)`

Evaluate model performance on current and additional objectives

Parameters

- `X` (`pd.DataFrame` or `np.array`) – data of shape `[n_samples, n_features]`
- `y` (`pd.Series`) – true labels of length `[n_samples]`
- `other_objectives` (`list`) – list of other objectives to score

Returns score, ordered dictionary of other objective scores

Return type `float, dict`

`evalml.pipelines.CatBoostClassificationPipeline`

`class evalml.pipelines.CatBoostClassificationPipeline(parameters, objective, random_state=0)`

CatBoost Pipeline for both binary and multiclass classification. CatBoost is an open-source library and natively supports categorical features.

For more information, check out <https://catboost.ai/> Note: `impute_strategy` must support both string and numeric data

`name = 'Cat Boost Classification Pipeline'`

`summary = 'CatBoost Classifier w/ Simple Imputer'`

`component_graph = ['Simple Imputer', <evalml.pipelines.components.estimators.classifiers.CatBoostClassifier>']`

`supported_problem_types = ['binary', 'multiclass']`

`model_family = 'catboost'`

```
hyperparameters = {'eta': Real(low=0, high=1, prior='uniform', transform='identity'),
custom_hyperparameters = {'impute_strategy': ['most_frequent']}}
```

Instance attributes

<code>feature_importances</code>	Return feature importances.
<code>parameters</code>	Returns parameter dictionary for this pipeline

Methods:

<code>__init__</code>	Machine learning pipeline made out of transformers and a estimator.
<code>describe</code>	Outputs pipeline details including component parameters
<code>feature_importance_graph</code>	Generate a bar graph of the pipeline's feature importances
<code>fit</code>	Build a model
<code>get_component</code>	Returns component by name
<code>graph</code>	Generate an image representing the pipeline graph
<code>load</code>	Loads pipeline at file path
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.
<code>save</code>	Saves pipeline at file path
<code>score</code>	Evaluate model performance on current and additional objectives

evalml.pipelines.CatBoostClassificationPipeline.__init__

CatBoostClassificationPipeline.__init__(parameters, objective, random_state=0)

Machine learning pipeline made out of transformers and a estimator.

Required Class Variables: `component_graph` (list): List of components in order. Accepts strings or ComponentBase objects in the list

`supported_problem_types` (list): List of problem types for this pipeline. Accepts strings or ProblemType enum in the list.

Parameters

- **objective** (*ObjectiveBase*) – the objective to optimize
- **parameters** (*dict*) – dictionary with component names as keys and dictionary of that component's parameters as values. An empty dictionary {} implies using all default values for component parameters.
- **random_state** (*int*, *np.random.RandomState*) – The random seed/state. Defaults to 0.

`evalml.pipelines.CatBoostClassificationPipeline.describe`

`CatBoostClassificationPipeline.describe()`

Outputs pipeline details including component parameters

Parameters `return_dict` (*bool*) – If True, return dictionary of information about pipeline.
Defaults to false

Returns dictionary of all component parameters if `return_dict` is True, else None

Return type dict

`evalml.pipelines.CatBoostClassificationPipeline.feature_importance_graph`

`CatBoostClassificationPipeline.feature_importance_graph` (*show_all_features=False*)

Generate a bar graph of the pipeline's feature importances

Parameters `show_all_features` (*bool, optional*) – If true, graph features with an importance value of zero. Defaults to false.

Returns plotly.Figure, a bar graph showing features and their importances

`evalml.pipelines.CatBoostClassificationPipeline.fit`

`CatBoostClassificationPipeline.fit` (*X, y, objective_fit_size=0.2*)

Build a model

Parameters

- **X** (*pd.DataFrame or np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series*) – the target training labels of length [n_samples]
- **feature_types** (*list, optional*) – list of feature types. either numeric or categorical. categorical features will automatically be encoded

Returns self

`evalml.pipelines.CatBoostClassificationPipeline.get_component`

`CatBoostClassificationPipeline.get_component` (*name*)

Returns component by name

Parameters `name` (*str*) – name of component

Returns component to return

Return type Component

`evalml.pipelines.CatBoostClassificationPipeline.graph`

`CatBoostClassificationPipeline.graph` (*filepath=None*)

Generate an image representing the pipeline graph

Parameters `filepath` (*str, optional*) – Path to where the graph should be saved. If set to None (as by default), the graph will not be saved.

Returns Graph object that can be directly displayed in Jupyter notebooks.

Return type graphviz.Digraph

`evalml.pipelines.CatBoostClassificationPipeline.load`

static `CatBoostClassificationPipeline.load(file_path)`

Loads pipeline at file path

Parameters `file_path` (*str*) – location to load file

Returns PipelineBase obj

`evalml.pipelines.CatBoostClassificationPipeline.predict`

`CatBoostClassificationPipeline.predict(X)`

Make predictions using selected features.

Parameters `X` (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]

Returns estimated labels

Return type *pd.Series*

`evalml.pipelines.CatBoostClassificationPipeline.predict_proba`

`CatBoostClassificationPipeline.predict_proba(X)`

Make probability estimates for labels.

Parameters `X` (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]

Returns probability estimates

Return type *pd.DataFrame*

`evalml.pipelines.CatBoostClassificationPipeline.save`

`CatBoostClassificationPipeline.save(file_path)`

Saves pipeline at file path

Parameters `file_path` (*str*) – location to save file

Returns None

`evalml.pipelines.CatBoostClassificationPipeline.score`

`CatBoostClassificationPipeline.score(X, y, other_objectives=None)`

Evaluate model performance on current and additional objectives

Parameters

- `X` (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]
- `y` (*pd.Series*) – true labels of length [n_samples]
- `other_objectives` (*list*) – list of other objectives to score

Returns score, ordered dictionary of other objective scores

Return type float, dict

evalml.pipelines.LogisticRegressionPipeline

```
class evalml.pipelines.LogisticRegressionPipeline(parameters, objective, random_state=0)
    Logistic Regression Pipeline for both binary and multiclass classification

    name = 'Logistic Regression Pipeline'
    summary = 'Logistic Regression Classifier w/ One Hot Encoder + Simple Imputer + Standard Scaler'
    component_graph = ['One Hot Encoder', 'Simple Imputer', 'Standard Scaler', <evalml.pipelines.LogisticRegressionPipeline>]
    supported_problem_types = ['binary', 'multiclass']
    model_family = 'linear_model'
    hyperparameters = {'C': Real(low=0.01, high=10, prior='uniform', transform='identity')}
    custom_hyperparameters = None
```

Instance attributes

<code>feature_importances</code>	Return feature importances.
<code>parameters</code>	Returns parameter dictionary for this pipeline

Methods:

<code>__init__</code>	Machine learning pipeline made out of transformers and a estimator.
<code>describe</code>	Outputs pipeline details including component parameters
<code>feature_importance_graph</code>	Generate a bar graph of the pipeline's feature importances
<code>fit</code>	Build a model
<code>get_component</code>	Returns component by name
<code>graph</code>	Generate an image representing the pipeline graph
<code>load</code>	Loads pipeline at file path
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.
<code>save</code>	Saves pipeline at file path
<code>score</code>	Evaluate model performance on current and additional objectives

evalml.pipelines.LogisticRegressionPipeline.__init__

LogisticRegressionPipeline.__init__(parameters, objective, random_state=0)
Machine learning pipeline made out of transformers and a estimator.

Required Class Variables: component_graph (list): List of components in order. Accepts strings or ComponentBase objects in the list

`supported_problem_types` (list): List of problem types for this pipeline. Accepts strings or Problem-Type enum in the list.

Parameters

- **objective** (*ObjectiveBase*) – the objective to optimize
- **parameters** (*dict*) – dictionary with component names as keys and dictionary of that component's parameters as values. An empty dictionary {} implies using all default values for component parameters.
- **random_state** (*int*, *np.random.RandomState*) – The random seed/state. Defaults to 0.

evalml.pipelines.LogisticRegressionPipeline.describe

`LogisticRegressionPipeline.describe()`

Outputs pipeline details including component parameters

Parameters `return_dict` (*bool*) – If True, return dictionary of information about pipeline. Defaults to false

Returns dictionary of all component parameters if `return_dict` is True, else None

Return type dict

evalml.pipelines.LogisticRegressionPipeline.feature_importance_graph

`LogisticRegressionPipeline.feature_importance_graph(show_all_features=False)`

Generate a bar graph of the pipeline's feature importances

Parameters `show_all_features` (*bool*, *optional*) – If true, graph features with an importance value of zero. Defaults to false.

Returns `plotly.Figure`, a bar graph showing features and their importances

evalml.pipelines.LogisticRegressionPipeline.fit

`LogisticRegressionPipeline.fit(X, y, objective_fit_size=0.2)`

Build a model

Parameters

- **X** (*pd.DataFrame* or *np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series*) – the target training labels of length [n_samples]
- **feature_types** (*list*, *optional*) – list of feature types. either numeric or categorical. categorical features will automatically be encoded

Returns self

evalml.pipelines.LogisticRegressionPipeline.get_component

`LogisticRegressionPipeline.get_component(name)`

Returns component by name

Parameters `name` (*str*) – name of component

Returns component to return

Return type Component

evalml.pipelines.LogisticRegressionPipeline.graph

`LogisticRegressionPipeline.graph(filepath=None)`

Generate an image representing the pipeline graph

Parameters `filepath` (*str, optional*) – Path to where the graph should be saved. If set to None (as by default), the graph will not be saved.

Returns Graph object that can be directly displayed in Jupyter notebooks.

Return type graphviz.Digraph

evalml.pipelines.LogisticRegressionPipeline.load

static `LogisticRegressionPipeline.load(file_path)`

Loads pipeline at file path

Parameters `file_path` (*str*) – location to load file

Returns PipelineBase obj

evalml.pipelines.LogisticRegressionPipeline.predict

`LogisticRegressionPipeline.predict(X)`

Make predictions using selected features.

Parameters `X` (*pd.DataFrame or np.array*) – data of shape [n_samples, n_features]

Returns estimated labels

Return type pd.Series

evalml.pipelines.LogisticRegressionPipeline.predict_proba

`LogisticRegressionPipeline.predict_proba(X)`

Make probability estimates for labels.

Parameters `X` (*pd.DataFrame or np.array*) – data of shape [n_samples, n_features]

Returns probability estimates

Return type pd.DataFrame

evalml.pipelines.LogisticRegressionPipeline.save

LogisticRegressionPipeline.**save** (*file_path*)

Saves pipeline at file path

Parameters **file_path** (*str*) – location to save file

Returns None

evalml.pipelines.LogisticRegressionPipeline.score

LogisticRegressionPipeline.**score** (*X*, *y*, *other_objectives=None*)

Evaluate model performance on current and additional objectives

Parameters

- **X** (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]
- **y** (*pd.Series*) – true labels of length [n_samples]
- **other_objectives** (*list*) – list of other objectives to score

Returns score, ordered dictionary of other objective scores

Return type float, dict

evalml.pipelines.RFRegressionPipeline

class evalml.pipelines.**RFRegressionPipeline** (*parameters*, *objective*, *random_state=0*)

Random Forest Pipeline for regression problems

name = 'Random Forest Regression Pipeline'

summary = 'Random Forest Regressor w/ One Hot Encoder + Simple Imputer + RF Regressor'

component_graph = ['One Hot Encoder', 'Simple Imputer', 'RF Regressor Select From Mode']

supported_problem_types = ['regression']

model_family = 'random_forest'

hyperparameters = {'impute_strategy': ['mean', 'median', 'most_frequent'], 'max_depth': [5, 10, 15, 20, 25, 30, 35, 40, 45, 50]}

custom_hyperparameters = None

Instance attributes

<code>feature_importances</code>	Return feature importances.
<code>parameters</code>	Returns parameter dictionary for this pipeline

Methods:

<code>__init__</code>	Machine learning pipeline made out of transformers and a estimator.
<code>describe</code>	Outputs pipeline details including component parameters

Continued on next page

Table 44 – continued from previous page

<code>feature_importance_graph</code>	Generate a bar graph of the pipeline’s feature importances
<code>fit</code>	Build a model
<code>get_component</code>	Returns component by name
<code>graph</code>	Generate an image representing the pipeline graph
<code>load</code>	Loads pipeline at file path
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.
<code>save</code>	Saves pipeline at file path
<code>score</code>	Evaluate model performance on current and additional objectives

`evalml.pipelines.RFRegressionPipeline.__init__`

`RFRegressionPipeline.__init__(parameters, objective, random_state=0)`

Machine learning pipeline made out of transformers and a estimator.

Required Class Variables: `component_graph` (list): List of components in order. Accepts strings or `ComponentBase` objects in the list

`supported_problem_types` (list): List of problem types for this pipeline. Accepts strings or `ProblemType` enum in the list.

Parameters

- **objective** (*ObjectiveBase*) – the objective to optimize
- **parameters** (*dict*) – dictionary with component names as keys and dictionary of that component’s parameters as values. An empty dictionary `{}` implies using all default values for component parameters.
- **random_state** (*int*, *np.random.RandomState*) – The random seed/state. Defaults to 0.

`evalml.pipelines.RFRegressionPipeline.describe`

`RFRegressionPipeline.describe()`

Outputs pipeline details including component parameters

Parameters `return_dict` (*bool*) – If True, return dictionary of information about pipeline. Defaults to false

Returns dictionary of all component parameters if `return_dict` is True, else None

Return type dict

`evalml.pipelines.RFRegressionPipeline.feature_importance_graph`

`RFRegressionPipeline.feature_importance_graph(show_all_features=False)`

Generate a bar graph of the pipeline’s feature importances

Parameters `show_all_features` (*bool*, *optional*) – If true, graph features with an importance value of zero. Defaults to false.

Returns `plotly.Figure`, a bar graph showing features and their importances

evalml.pipelines.RFRegressionPipeline.fit

RFRegressionPipeline.**fit** (*X*, *y*, *objective_fit_size=0.2*)

Build a model

Parameters

- **X** (*pd.DataFrame* or *np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series*) – the target training labels of length [n_samples]
- **feature_types** (*list*, *optional*) – list of feature types. either numeric or categorical. categorical features will automatically be encoded

Returns self

evalml.pipelines.RFRegressionPipeline.get_component

RFRegressionPipeline.**get_component** (*name*)

Returns component by name

Parameters **name** (*str*) – name of component

Returns component to return

Return type Component

evalml.pipelines.RFRegressionPipeline.graph

RFRegressionPipeline.**graph** (*filepath=None*)

Generate an image representing the pipeline graph

Parameters **filepath** (*str*, *optional*) – Path to where the graph should be saved. If set to None (as by default), the graph will not be saved.

Returns Graph object that can be directly displayed in Jupyter notebooks.

Return type graphviz.Digraph

evalml.pipelines.RFRegressionPipeline.load

static RFRegressionPipeline.**load** (*file_path*)

Loads pipeline at file path

Parameters **file_path** (*str*) – location to load file

Returns PipelineBase obj

evalml.pipelines.RFRegressionPipeline.predict

RFRegressionPipeline.**predict** (*X*)

Make predictions using selected features.

Parameters **X** (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]

Returns estimated labels

Return type `pd.Series`

`evalml.pipelines.RFRegressionPipeline.predict_proba`

`RFRegressionPipeline.predict_proba(X)`

Make probability estimates for labels.

Parameters `X` (`pd.DataFrame` or `np.array`) – data of shape `[n_samples, n_features]`

Returns probability estimates

Return type `pd.DataFrame`

`evalml.pipelines.RFRegressionPipeline.save`

`RFRegressionPipeline.save(file_path)`

Saves pipeline at file path

Parameters `file_path` (`str`) – location to save file

Returns `None`

`evalml.pipelines.RFRegressionPipeline.score`

`RFRegressionPipeline.score(X, y, other_objectives=None)`

Evaluate model performance on current and additional objectives

Parameters

- `X` (`pd.DataFrame` or `np.array`) – data of shape `[n_samples, n_features]`
- `y` (`pd.Series`) – true labels of length `[n_samples]`
- `other_objectives` (`list`) – list of other objectives to score

Returns score, ordered dictionary of other objective scores

Return type `float, dict`

`evalml.pipelines.CatBoostRegressionPipeline`

`class evalml.pipelines.CatBoostRegressionPipeline(parameters, objective, random_state=0)`

CatBoost Pipeline for regression problems. CatBoost is an open-source library and natively supports categorical features.

For more information, check out <https://catboost.ai/>

Note: `impute_strategy` must support both string and numeric data

`name = 'Cat Boost Regression Pipeline'`

`summary = 'CatBoost Regressor w/ Simple Imputer'`

`component_graph = ['Simple Imputer', <evalml.pipelines.components.estimators.regressor`

`supported_problem_types = ['regression']`

`model_family = 'catboost'`


```
hyperparameters = {'eta': Real(low=0, high=1, prior='uniform', transform='identity'),
custom_hyperparameters = {'impute_strategy': ['most_frequent']}}
```

Instance attributes

<code>feature_importances</code>	Return feature importances.
<code>parameters</code>	Returns parameter dictionary for this pipeline

Methods:

<code>__init__</code>	Machine learning pipeline made out of transformers and a estimator.
<code>describe</code>	Outputs pipeline details including component parameters
<code>feature_importance_graph</code>	Generate a bar graph of the pipeline's feature importances
<code>fit</code>	Build a model
<code>get_component</code>	Returns component by name
<code>graph</code>	Generate an image representing the pipeline graph
<code>load</code>	Loads pipeline at file path
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.
<code>save</code>	Saves pipeline at file path
<code>score</code>	Evaluate model performance on current and additional objectives

evalml.pipelines.CatBoostRegressionPipeline.__init__

CatBoostRegressionPipeline.__init__(parameters, objective, random_state=0)

Machine learning pipeline made out of transformers and a estimator.

Required Class Variables: `component_graph` (list): List of components in order. Accepts strings or ComponentBase objects in the list

`supported_problem_types` (list): List of problem types for this pipeline. Accepts strings or ProblemType enum in the list.

Parameters

- **objective** (*ObjectiveBase*) – the objective to optimize
- **parameters** (*dict*) – dictionary with component names as keys and dictionary of that component's parameters as values. An empty dictionary {} implies using all default values for component parameters.
- **random_state** (*int*, *np.random.RandomState*) – The random seed/state. Defaults to 0.

`evalml.pipelines.CatBoostRegressionPipeline.describe`

`CatBoostRegressionPipeline.describe()`

Outputs pipeline details including component parameters

Parameters `return_dict` (*bool*) – If True, return dictionary of information about pipeline.
Defaults to false

Returns dictionary of all component parameters if `return_dict` is True, else None

Return type dict

`evalml.pipelines.CatBoostRegressionPipeline.feature_importance_graph`

`CatBoostRegressionPipeline.feature_importance_graph(show_all_features=False)`

Generate a bar graph of the pipeline's feature importances

Parameters `show_all_features` (*bool, optional*) – If true, graph features with an importance value of zero. Defaults to false.

Returns plotly.Figure, a bar graph showing features and their importances

`evalml.pipelines.CatBoostRegressionPipeline.fit`

`CatBoostRegressionPipeline.fit(X, y, objective_fit_size=0.2)`

Build a model

Parameters

- `X` (*pd.DataFrame* or *np.array*) – the input training data of shape [n_samples, n_features]
- `y` (*pd.Series*) – the target training labels of length [n_samples]
- `feature_types` (*list, optional*) – list of feature types. either numeric or categorical. categorical features will automatically be encoded

Returns self

`evalml.pipelines.CatBoostRegressionPipeline.get_component`

`CatBoostRegressionPipeline.get_component(name)`

Returns component by name

Parameters `name` (*str*) – name of component

Returns component to return

Return type Component

`evalml.pipelines.CatBoostRegressionPipeline.graph`

`CatBoostRegressionPipeline.graph(filepath=None)`

Generate an image representing the pipeline graph

Parameters `filepath` (*str, optional*) – Path to where the graph should be saved. If set to None (as by default), the graph will not be saved.

Returns Graph object that can be directly displayed in Jupyter notebooks.

Return type graphviz.Digraph

evalml.pipelines.CatBoostRegressionPipeline.load

static CatBoostRegressionPipeline.**load** (*file_path*)

Loads pipeline at file path

Parameters **file_path** (*str*) – location to load file

Returns PipelineBase obj

evalml.pipelines.CatBoostRegressionPipeline.predict

CatBoostRegressionPipeline.**predict** (*X*)

Make predictions using selected features.

Parameters **X** (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]

Returns estimated labels

Return type pd.Series

evalml.pipelines.CatBoostRegressionPipeline.predict_proba

CatBoostRegressionPipeline.**predict_proba** (*X*)

Make probability estimates for labels.

Parameters **X** (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]

Returns probability estimates

Return type pd.DataFrame

evalml.pipelines.CatBoostRegressionPipeline.save

CatBoostRegressionPipeline.**save** (*file_path*)

Saves pipeline at file path

Parameters **file_path** (*str*) – location to save file

Returns None

evalml.pipelines.CatBoostRegressionPipeline.score

CatBoostRegressionPipeline.**score** (*X*, *y*, *other_objectives=None*)

Evaluate model performance on current and additional objectives

Parameters

- **X** (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]
- **y** (*pd.Series*) – true labels of length [n_samples]
- **other_objectives** (*list*) – list of other objectives to score

Returns score, ordered dictionary of other objective scores

Return type float, dict

evalml.pipelines.LinearRegressionPipeline

```
class evalml.pipelines.LinearRegressionPipeline(parameters, objective, random_state=0)
    Linear Regression Pipeline for regression problems
    name = 'Linear Regression Pipeline'
    summary = 'Linear Regressor w/ One Hot Encoder + Simple Imputer + Standard Scaler'
    component_graph = ['One Hot Encoder', 'Simple Imputer', 'Standard Scaler', <evalml.pipelines.LinearRegressionPipeline>]
    supported_problem_types = ['regression']
    model_family = 'linear_model'
    hyperparameters = {'fit_intercept': [True, False], 'impute_strategy': ['mean', 'median']}
    custom_hyperparameters = None
```

Instance attributes

<code>feature_importances</code>	Return feature importances.
<code>parameters</code>	Returns parameter dictionary for this pipeline

Methods:

<code>__init__</code>	Machine learning pipeline made out of transformers and a estimator.
<code>describe</code>	Outputs pipeline details including component parameters
<code>feature_importance_graph</code>	Generate a bar graph of the pipeline's feature importances
<code>fit</code>	Build a model
<code>get_component</code>	Returns component by name
<code>graph</code>	Generate an image representing the pipeline graph
<code>load</code>	Loads pipeline at file path
<code>predict</code>	Make predictions using selected features.
<code>predict_proba</code>	Make probability estimates for labels.
<code>save</code>	Saves pipeline at file path
<code>score</code>	Evaluate model performance on current and additional objectives

evalml.pipelines.LinearRegressionPipeline.__init__

LinearRegressionPipeline.__init__(parameters, objective, random_state=0)
Machine learning pipeline made out of transformers and a estimator.

Required Class Variables: component_graph (list): List of components in order. Accepts strings or ComponentBase objects in the list

`supported_problem_types` (list): List of problem types for this pipeline. Accepts strings or Problem-Type enum in the list.

Parameters

- **objective** (*ObjectiveBase*) – the objective to optimize
- **parameters** (*dict*) – dictionary with component names as keys and dictionary of that component's parameters as values. An empty dictionary {} implies using all default values for component parameters.
- **random_state** (*int*, *np.random.RandomState*) – The random seed/state. Defaults to 0.

evalml.pipelines.LinearRegressionPipeline.describe

`LinearRegressionPipeline.describe()`

Outputs pipeline details including component parameters

Parameters **return_dict** (*bool*) – If True, return dictionary of information about pipeline. Defaults to false

Returns dictionary of all component parameters if return_dict is True, else None

Return type dict

evalml.pipelines.LinearRegressionPipeline.feature_importance_graph

`LinearRegressionPipeline.feature_importance_graph(show_all_features=False)`

Generate a bar graph of the pipeline's feature importances

Parameters **show_all_features** (*bool*, *optional*) – If true, graph features with an importance value of zero. Defaults to false.

Returns plotly.Figure, a bar graph showing features and their importances

evalml.pipelines.LinearRegressionPipeline.fit

`LinearRegressionPipeline.fit(X, y, objective_fit_size=0.2)`

Build a model

Parameters

- **X** (*pd.DataFrame* or *np.array*) – the input training data of shape [n_samples, n_features]
- **y** (*pd.Series*) – the target training labels of length [n_samples]
- **feature_types** (*list*, *optional*) – list of feature types. either numeric or categorical. categorical features will automatically be encoded

Returns self

evalml.pipelines.LinearRegressionPipeline.get_component

`LinearRegressionPipeline.get_component(name)`

Returns component by name

Parameters `name` (*str*) – name of component

Returns component to return

Return type Component

evalml.pipelines.LinearRegressionPipeline.graph

`LinearRegressionPipeline.graph(filepath=None)`

Generate an image representing the pipeline graph

Parameters `filepath` (*str, optional*) – Path to where the graph should be saved. If set to None (as by default), the graph will not be saved.

Returns Graph object that can be directly displayed in Jupyter notebooks.

Return type graphviz.Digraph

evalml.pipelines.LinearRegressionPipeline.load

static `LinearRegressionPipeline.load(file_path)`

Loads pipeline at file path

Parameters `file_path` (*str*) – location to load file

Returns PipelineBase obj

evalml.pipelines.LinearRegressionPipeline.predict

`LinearRegressionPipeline.predict(X)`

Make predictions using selected features.

Parameters `X` (*pd.DataFrame or np.array*) – data of shape [n_samples, n_features]

Returns estimated labels

Return type pd.Series

evalml.pipelines.LinearRegressionPipeline.predict_proba

`LinearRegressionPipeline.predict_proba(X)`

Make probability estimates for labels.

Parameters `X` (*pd.DataFrame or np.array*) – data of shape [n_samples, n_features]

Returns probability estimates

Return type pd.DataFrame

evalml.pipelines.LinearRegressionPipeline.save`LinearRegressionPipeline.save(file_path)`

Saves pipeline at file path

Parameters `file_path` (*str*) – location to save file**Returns** None**evalml.pipelines.LinearRegressionPipeline.score**`LinearRegressionPipeline.score(X, y, other_objectives=None)`

Evaluate model performance on current and additional objectives

Parameters

- **X** (*pd.DataFrame* or *np.array*) – data of shape [n_samples, n_features]
- **y** (*pd.Series*) – true labels of length [n_samples]
- **other_objectives** (*list*) – list of other objectives to score

Returns score, ordered dictionary of other objective scores**Return type** float, dict**Pipeline Utils**

<code>get_pipelines</code>	Returns the pipelines allowed for a particular problem type.
<code>list_model_families</code>	List model type for a particular problem type

evalml.pipelines.get_pipelines`evalml.pipelines.get_pipelines(problem_type, model_families=None)`

Returns the pipelines allowed for a particular problem type.

Can also optionally filter by a list of model types.

Arguments:

Returns a list of pipeline classes**Return type** list[*PipelineBase*]**evalml.pipelines.list_model_families**`evalml.pipelines.list_model_families(problem_type)`

List model type for a particular problem type

Parameters `problem_types` (*ProblemTypes* or *str*) – binary, multiclass, or regression**Returns** a list of model families**Return type** list[*ModelFamily*]

Plotting

<code>PipelineBase.graph(filepath)</code>	Generate an image representing the pipeline graph
<code>PipelineBase.feature_importance_graph(filepath)</code>	Generate a bar graph of the pipeline's feature importances

evalml.pipelines.PipelineBase.graph

`pipelines.PipelineBase.graph(**kwargs)`
Call self as a function.

evalml.pipelines.PipelineBase.feature_importance_graph

`pipelines.PipelineBase.feature_importance_graph(**kwargs)`
Call self as a function.

Objective Functions

Domain Specific

<code>FraudCost</code>	Score the percentage of money lost of the total transaction amount process due to fraud
<code>LeadScoring</code>	Lead scoring

evalml.objectives.FraudCost

class `evalml.objectives.FraudCost` (`retry_percentage=0.5`, `interchange_fee=0.02`,
`fraud_payout_percentage=1.0`, `amount_col='amount'`,
`verbose=False`)
Score the percentage of money lost of the total transaction amount process due to fraud

Methods

<code>__init__</code>	Create instance of FraudCost
<code>decision_function</code>	Determine if transaction is fraud given predicted probabilities, dataframe with transaction amount, and threshold
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>objective_function</code>	Calculate amount lost to fraud per transaction given predictions, true values, and dataframe with transaction amount
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values

Continued on next page

Table 52 – continued from previous page

<code>supports_problem_type</code>	Checks if objective supports given ProblemType
------------------------------------	--

evalml.objectives.FraudCost.__init__

`FraudCost.__init__(retry_percentage=0.5, interchange_fee=0.02, fraud_payout_percentage=1.0, amount_col='amount', verbose=False)`
 Create instance of FraudCost

Parameters

- **retry_percentage** (*float*) – what percentage of customers will retry a transaction if it is declined? Between 0 and 1. Defaults to .5
- **interchange_fee** (*float*) – how much of each successful transaction do you collect? Between 0 and 1. Defaults to .02
- **fraud_payout_percentage** (*float*) – how percentage of fraud will you be unable to collect. Between 0 and 1. Defaults to 1.0
- **amount_col** (*str*) – name of column in data that contains the amount. defaults to “amount”

evalml.objectives.FraudCost.decision_function

`FraudCost.decision_function(y_predicted, extra_cols, threshold)`

Determine if transaction is fraud given predicted probabilities, dataframe with transaction amount, and threshold

Parameters

- **y_predicted** (*pd.Series*) – predicted labels
- **extra_cols** (*pd.DataFrame*) – extra data needed
- **threshold** (*float*) – dollar threshold to determine if transaction is fraud

Returns series of predicted fraud label using extra cols and threshold

Return type *pd.Series*

evalml.objectives.FraudCost.fit

`FraudCost.fit(y_predicted, y_true, extra_cols=None)`

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won’t be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns *self*

`evalml.objectives.FraudCost.objective_function`

`FraudCost.objective_function` (*y_predicted*, *y_true*, *extra_cols*)

Calculate amount lost to fraud per transaction given predictions, true values, and dataframe with transaction amount

Parameters

- **y_predicted** (*pd.Series*) – predicted fraud labels
- **y_true** (*pd.Series*) – true fraud labels
- **extra_cols** (*pd.DataFrame*) – extra data needed

Returns amount lost to fraud per transaction

Return type float

`evalml.objectives.FraudCost.predict`

`FraudCost.predict` (*y_predicted*, *extra_cols=None*)

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

`evalml.objectives.FraudCost.score`

`FraudCost.score` (*y_predicted*, *y_true*, *extra_cols=None*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

`evalml.objectives.FraudCost.supports_problem_type`

classmethod `FraudCost.supports_problem_type` (*problem_type*)

Checks if objective supports given `ProblemType`

Parameters **problem_type** (*str* or *ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

evalml.objectives.LeadScoring

class evalml.objectives.LeadScoring(*true_positives=1, false_positives=-1, verbose=False*)
Lead scoring

Methods

<code>__init__</code>	Create instance.
<code>decision_function</code>	
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>objective_function</code>	
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.LeadScoring.__init__

LeadScoring.__init__(*true_positives=1, false_positives=-1, verbose=False*)
Create instance.

Parameters

- **label** (*int*) – label to optimize threshold for
- **true_positives** (*int*) – reward for a true positive
- **false_positives** (*int*) – cost for a false positive. Should be negative.

evalml.objectives.LeadScoring.decision_function

LeadScoring.decision_function(*y_predicted, threshold*)

evalml.objectives.LeadScoring.fit

LeadScoring.fit(*y_predicted, y_true, extra_cols=None*)
Learn the objective function based on the predictions from a model.
If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.LeadScoring.objective_function

`LeadScoring.objective_function(y_predicted, y_true)`

evalml.objectives.LeadScoring.predict

`LeadScoring.predict(y_predicted, extra_cols=None)`

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called

Parameters `y_predicted` – the prediction to transform to final prediction

Returns predictions

evalml.objectives.LeadScoring.score

`LeadScoring.score(y_predicted, y_true, extra_cols=None)`

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

evalml.objectives.LeadScoring.supports_problem_type

classmethod `LeadScoring.supports_problem_type(problem_type)`

Checks if objective supports given `ProblemType`

Parameters `problem_type` (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

Classification

<i>F1</i>	F1 score for binary classification
<i>F1Micro</i>	F1 score for multiclass classification using micro averaging
<i>F1Macro</i>	F1 score for multiclass classification using macro averaging
<i>F1Weighted</i>	F1 score for multiclass classification using weighted averaging

Continued on next page

Table 54 – continued from previous page

<i>Precision</i>	Precision score for binary classification
<i>PrecisionMicro</i>	Precision score for multiclass classification using micro averaging
<i>PrecisionMacro</i>	Precision score for multiclass classification using macro averaging
<i>PrecisionWeighted</i>	Precision score for multiclass classification using weighted averaging
<i>Recall</i>	Recall score for binary classification
<i>RecallMicro</i>	Recall score for multiclass classification using micro averaging
<i>RecallMacro</i>	Recall score for multiclass classification using macro averaging
<i>RecallWeighted</i>	Recall score for multiclass classification using weighted averaging
<i>AUC</i>	AUC score for binary classification
<i>AUCMicro</i>	AUC score for multiclass classification using micro averaging
<i>AUCMacro</i>	AUC score for multiclass classification using macro averaging
<i>AUCWeighted</i>	AUC Score for multiclass classification using weighted averaging
<i>LogLoss</i>	Log Loss for both binary and multiclass classification
<i>MCC</i>	Matthews correlation coefficient for both binary and multiclass classification
<i>ROC</i>	Receiver Operating Characteristic score for binary classification.
<i>ConfusionMatrix</i>	Confusion matrix for classification problems

evalml.objectives.F1

class evalml.objectives.F1 (*verbose=False*)
F1 score for binary classification

Methods

<i>__init__</i>	Initialize self.
<i>fit</i>	Learn the objective function based on the predictions from a model.
<i>predict</i>	Apply the learned objective function to the output of a model.
<i>score</i>	Calculate score from applying fitted objective to predicted values
<i>supports_problem_type</i>	Checks if objective supports given ProblemType

evalml.objectives.F1.__init__

F1.__init__ (*verbose=False*)
Initialize self. See help(type(self)) for accurate signature.

`evalml.objectives.F1.fit`

`F1.fit(y_predicted, y_true, extra_cols=None)`

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns `self`

`evalml.objectives.F1.predict`

`F1.predict(y_predicted, extra_cols=None)`

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

`evalml.objectives.F1.score`

`F1.score(y_predicted, y_true)`

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

`evalml.objectives.F1.supports_problem_type`

classmethod `F1.supports_problem_type(problem_type)`

Checks if objective supports given `ProblemType`

Parameters **problem_type** (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

evalml.objectives.F1Micro

class evalml.objectives.**F1Micro** (*verbose=False*)
 F1 score for multiclass classification using micro averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.F1Micro.__init__

F1Micro.__init__ (*verbose=False*)
 Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.F1Micro.fit

F1Micro.fit (*y_predicted, y_true, extra_cols=None*)
 Learn the objective function based on the predictions from a model.
 If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.F1Micro.predict

F1Micro.predict (*y_predicted, extra_cols=None*)
 Apply the learned objective function to the output of a model.
 If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.F1Micro.score**F1Micro.score** (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to `True`**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is `True`, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is `True`.

Returns score**evalml.objectives.F1Micro.supports_problem_type****classmethod** **F1Micro.supports_problem_type** (*problem_type*)Checks if objective supports given `ProblemType`**Parameters** **problem_type** (*str or ProblemType*) – problem type to check**Returns** whether objective supports `ProblemType`**Return type** bool**evalml.objectives.F1Macro****class** **evalml.objectives.F1Macro** (*verbose=False*)

F1 score for multiclass classification using macro averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given <code>ProblemType</code>

evalml.objectives.F1Macro.__init__**F1Macro.__init__** (*verbose=False*)Initialize self. See `help(type(self))` for accurate signature.

evalml.objectives.F1Macro.fit**F1Macro.fit** (*y_predicted*, *y_true*, *extra_cols=None*)

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won't be called**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns `self`**evalml.objectives.F1Macro.predict****F1Macro.predict** (*y_predicted*, *extra_cols=None*)

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called**Parameters** **y_predicted** – the prediction to transform to final prediction**Returns** predictions**evalml.objectives.F1Macro.score****F1Macro.score** (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score**evalml.objectives.F1Macro.supports_problem_type****classmethod** **F1Macro.supports_problem_type** (*problem_type*)Checks if objective supports given `ProblemType`**Parameters** **problem_type** (*str* or *ProblemType*) – problem type to check**Returns** whether objective supports `ProblemType`**Return type** bool

evalml.objectives.F1Weighted

class evalml.objectives.**F1Weighted** (*verbose=False*)
F1 score for multiclass classification using weighted averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.F1Weighted.__init__

F1Weighted.__init__ (*verbose=False*)
Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.F1Weighted.fit

F1Weighted.fit (*y_predicted, y_true, extra_cols=None*)
Learn the objective function based on the predictions from a model.
If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.F1Weighted.predict

F1Weighted.predict (*y_predicted, extra_cols=None*)
Apply the learned objective function to the output of a model.
If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.F1Weighted.score**F1Weighted.score** (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to `True`**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is `True`, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is `True`.

Returns score**evalml.objectives.F1Weighted.supports_problem_type****classmethod** **F1Weighted.supports_problem_type** (*problem_type*)Checks if objective supports given `ProblemType`**Parameters** **problem_type** (*str or ProblemType*) – problem type to check**Returns** whether objective supports `ProblemType`**Return type** bool**evalml.objectives.Precision****class** **evalml.objectives.Precision** (*verbose=False*)

Precision score for binary classification

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given <code>ProblemType</code>

evalml.objectives.Precision.__init__**Precision.__init__** (*verbose=False*)Initialize self. See `help(type(self))` for accurate signature.

`evalml.objectives.Precision.fit`

`Precision.fit(y_predicted, y_true, extra_cols=None)`

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns self

`evalml.objectives.Precision.predict`

`Precision.predict(y_predicted, extra_cols=None)`

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

`evalml.objectives.Precision.score`

`Precision.score(y_predicted, y_true)`

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

`evalml.objectives.Precision.supports_problem_type`

classmethod `Precision.supports_problem_type(problem_type)`

Checks if objective supports given `ProblemType`

Parameters **problem_type** (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

evalml.objectives.PrecisionMicro

class evalml.objectives.**PrecisionMicro** (*verbose=False*)
 Precision score for multiclass classification using micro averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.PrecisionMicro.__init__

PrecisionMicro.__init__ (*verbose=False*)
 Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.PrecisionMicro.fit

PrecisionMicro.fit (*y_predicted, y_true, extra_cols=None*)
 Learn the objective function based on the predictions from a model.
 If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.PrecisionMicro.predict

PrecisionMicro.predict (*y_predicted, extra_cols=None*)
 Apply the learned objective function to the output of a model.
 If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

`evalml.objectives.PrecisionMicro.score`

`PrecisionMicro.score` (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to `True`

Parameters

- **`y_predicted`** (*list*) – the predictions from the model. If `needs_proba` is `True`, it is the probability estimates
- **`y_true`** (*list*) – the ground truth for the predictions.
- **`extra_cols`** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is `True`.

Returns score

`evalml.objectives.PrecisionMicro.supports_problem_type`

classmethod `PrecisionMicro.supports_problem_type` (*problem_type*)

Checks if objective supports given `ProblemType`

Parameters **`problem_type`** (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

`evalml.objectives.PrecisionMacro`

class `evalml.objectives.PrecisionMacro` (*verbose=False*)

Precision score for multiclass classification using macro averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given <code>ProblemType</code>

`evalml.objectives.PrecisionMacro.__init__`

`PrecisionMacro.__init__` (*verbose=False*)

Initialize self. See `help(type(self))` for accurate signature.

evalml.objectives.PrecisionMacro.fit`PrecisionMacro.fit(y_predicted, y_true, extra_cols=None)`

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns self

evalml.objectives.PrecisionMacro.predict`PrecisionMacro.predict(y_predicted, extra_cols=None)`

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.PrecisionMacro.score`PrecisionMacro.score(y_predicted, y_true)`

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

evalml.objectives.PrecisionMacro.supports_problem_type`classmethod PrecisionMacro.supports_problem_type(problem_type)`

Checks if objective supports given `ProblemType`

Parameters **problem_type** (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

evalml.objectives.PrecisionWeighted

class evalml.objectives.PrecisionWeighted (*verbose=False*)
Precision score for multiclass classification using weighted averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.PrecisionWeighted.__init__

PrecisionWeighted.__init__ (*verbose=False*)
Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.PrecisionWeighted.fit

PrecisionWeighted.**fit** (*y_predicted, y_true, extra_cols=None*)
Learn the objective function based on the predictions from a model.
If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.PrecisionWeighted.predict

PrecisionWeighted.**predict** (*y_predicted, extra_cols=None*)
Apply the learned objective function to the output of a model.
If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.PrecisionWeighted.score

PrecisionWeighted.**score**(*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

evalml.objectives.PrecisionWeighted.supports_problem_type

classmethod PrecisionWeighted.**supports_problem_type**(*problem_type*)

Checks if objective supports given ProblemType

Parameters **problem_type** (*str or ProblemType*) – problem type to check

Returns whether objective supports ProblemType

Return type bool

evalml.objectives.Recall

class evalml.objectives.**Recall**(*verbose=False*)

Recall score for binary classification

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.Recall.__init__

Recall.**__init__**(*verbose=False*)

Initialize self. See `help(type(self))` for accurate signature.

`evalml.objectives.Recall.fit`

`Recall.fit(y_predicted, y_true, extra_cols=None)`

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns `self`

`evalml.objectives.Recall.predict`

`Recall.predict(y_predicted, extra_cols=None)`

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

`evalml.objectives.Recall.score`

`Recall.score(y_predicted, y_true)`

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

`evalml.objectives.Recall.supports_problem_type`

classmethod `Recall.supports_problem_type(problem_type)`

Checks if objective supports given `ProblemType`

Parameters **problem_type** (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

evalml.objectives.RecallMicro

class evalml.objectives.RecallMicro (*verbose=False*)
 Recall score for multiclass classification using micro averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.RecallMicro.__init__

RecallMicro.**__init__** (*verbose=False*)
 Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.RecallMicro.fit

RecallMicro.**fit** (*y_predicted, y_true, extra_cols=None*)
 Learn the objective function based on the predictions from a model.
 If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.RecallMicro.predict

RecallMicro.**predict** (*y_predicted, extra_cols=None*)
 Apply the learned objective function to the output of a model.
 If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.RecallMicro.score

`RecallMicro.score(y_predicted, y_true)`

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to `True`

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is `True`, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is `True`.

Returns score

evalml.objectives.RecallMicro.supports_problem_type

classmethod `RecallMicro.supports_problem_type(problem_type)`

Checks if objective supports given `ProblemType`

Parameters **problem_type** (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

evalml.objectives.RecallMacro

class `evalml.objectives.RecallMacro(verbose=False)`

Recall score for multiclass classification using macro averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given <code>ProblemType</code>

evalml.objectives.RecallMacro.__init__

`RecallMacro.__init__(verbose=False)`

Initialize self. See `help(type(self))` for accurate signature.

evalml.objectives.RecallMacro.fit`RecallMacro.fit(y_predicted, y_true, extra_cols=None)`

Learn the objective function based on the predictions from a model.

If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.RecallMacro.predict`RecallMacro.predict(y_predicted, extra_cols=None)`

Apply the learned objective function to the output of a model.

If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.RecallMacro.score`RecallMacro.score(y_predicted, y_true)`

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set greater_is_better attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns score

evalml.objectives.RecallMacro.supports_problem_type`classmethod RecallMacro.supports_problem_type(problem_type)`

Checks if objective supports given ProblemType

Parameters **problem_type** (*str or ProblemType*) – problem type to check

Returns whether objective supports ProblemType

Return type bool

evalml.objectives.RecallWeighted

class evalml.objectives.**RecallWeighted** (*verbose=False*)
Recall score for multiclass classification using weighted averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.RecallWeighted.__init__

RecallWeighted.**__init__** (*verbose=False*)
Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.RecallWeighted.fit

RecallWeighted.**fit** (*y_predicted, y_true, extra_cols=None*)
Learn the objective function based on the predictions from a model.
If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.RecallWeighted.predict

RecallWeighted.**predict** (*y_predicted, extra_cols=None*)
Apply the learned objective function to the output of a model.
If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.RecallWeighted.score

RecallWeighted.**score** (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to `True`

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is `True`, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is `True`.

Returns score

evalml.objectives.RecallWeighted.supports_problem_type

classmethod RecallWeighted.**supports_problem_type** (*problem_type*)

Checks if objective supports given `ProblemType`

Parameters **problem_type** (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

evalml.objectives.AUC

class evalml.objectives.**AUC** (*verbose=False*)

AUC score for binary classification

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given <code>ProblemType</code>

evalml.objectives.AUC.__init__

AUC.**__init__** (*verbose=False*)

Initialize self. See `help(type(self))` for accurate signature.

`evalml.objectives.AUC.fit`

`AUC.fit` (*y_predicted*, *y_true*, *extra_cols=None*)

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won't be called

Parameters

- **`y_predicted`** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **`y_true`** (*list*) – the ground truth for the predictions.
- **`extra_cols`** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns `self`

`evalml.objectives.AUC.predict`

`AUC.predict` (*y_predicted*, *extra_cols=None*)

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called

Parameters **`y_predicted`** – the prediction to transform to final prediction

Returns predictions

`evalml.objectives.AUC.score`

`AUC.score` (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- **`y_predicted`** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **`y_true`** (*list*) – the ground truth for the predictions.
- **`extra_cols`** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

`evalml.objectives.AUC.supports_problem_type`

classmethod `AUC.supports_problem_type` (*problem_type*)

Checks if objective supports given `ProblemType`

Parameters **`problem_type`** (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

evalml.objectives.AUCMicro

class evalml.objectives.**AUCMicro** (*verbose=False*)
 AUC score for multiclass classification using micro averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.AUCMicro.__init__

AUCMicro.__init__ (*verbose=False*)
 Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.AUCMicro.fit

AUCMicro.fit (*y_predicted, y_true, extra_cols=None*)
 Learn the objective function based on the predictions from a model.
 If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.AUCMicro.predict

AUCMicro.predict (*y_predicted, extra_cols=None*)
 Apply the learned objective function to the output of a model.
 If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.AUCMicro.score**AUCMicro.score** (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to `True`**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is `True`, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is `True`.

Returns score**evalml.objectives.AUCMicro.supports_problem_type****classmethod** **AUCMicro.supports_problem_type** (*problem_type*)Checks if objective supports given `ProblemType`**Parameters** **problem_type** (*str or ProblemType*) – problem type to check**Returns** whether objective supports `ProblemType`**Return type** bool**evalml.objectives.AUCMacro****class** **evalml.objectives.AUCMacro** (*verbose=False*)

AUC score for multiclass classification using macro averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given <code>ProblemType</code>

evalml.objectives.AUCMacro.__init__**AUCMacro.__init__** (*verbose=False*)Initialize self. See `help(type(self))` for accurate signature.

evalml.objectives.AUCMacro.fit**AUCMacro.fit** (*y_predicted*, *y_true*, *extra_cols=None*)

Learn the objective function based on the predictions from a model.

If *needs_fitting* is false, this method won't be called**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If *needs_proba* is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if *uses_extra_columns* is True.

Returns self**evalml.objectives.AUCMacro.predict****AUCMacro.predict** (*y_predicted*, *extra_cols=None*)

Apply the learned objective function to the output of a model.

If *needs_fitting* is false, this method won't be called**Parameters** **y_predicted** – the prediction to transform to final prediction**Returns** predictions**evalml.objectives.AUCMacro.score****AUCMacro.score** (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set *greater_is_better* attribute to True**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If *needs_proba* is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if *uses_extra_columns* is True.

Returns score**evalml.objectives.AUCMacro.supports_problem_type****classmethod** **AUCMacro.supports_problem_type** (*problem_type*)Checks if objective supports given *ProblemType***Parameters** **problem_type** (*str* or *ProblemType*) – problem type to check**Returns** whether objective supports *ProblemType***Return type** bool

evalml.objectives.AUCWeighted

class evalml.objectives.**AUCWeighted** (*verbose=False*)
AUC Score for multiclass classification using weighted averaging

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.AUCWeighted.__init__

AUCWeighted.**__init__** (*verbose=False*)
Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.AUCWeighted.fit

AUCWeighted.**fit** (*y_predicted, y_true, extra_cols=None*)
Learn the objective function based on the predictions from a model.
If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.AUCWeighted.predict

AUCWeighted.**predict** (*y_predicted, extra_cols=None*)
Apply the learned objective function to the output of a model.
If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.AUCWeighted.score**AUCWeighted.score** (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to `True`**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is `True`, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is `True`.

Returns score**evalml.objectives.AUCWeighted.supports_problem_type****classmethod** **AUCWeighted.supports_problem_type** (*problem_type*)Checks if objective supports given `ProblemType`**Parameters** **problem_type** (*str or ProblemType*) – problem type to check**Returns** whether objective supports `ProblemType`**Return type** bool**evalml.objectives.LogLoss****class** **evalml.objectives.LogLoss** (*verbose=False*)

Log Loss for both binary and multiclass classification

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given <code>ProblemType</code>

evalml.objectives.LogLoss.__init__**LogLoss.__init__** (*verbose=False*)Initialize self. See `help(type(self))` for accurate signature.

`evalml.objectives.LogLoss.fit`

`LogLoss.fit` (*y_predicted*, *y_true*, *extra_cols=None*)

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won't be called

Parameters

- **`y_predicted`** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **`y_true`** (*list*) – the ground truth for the predictions.
- **`extra_cols`** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns `self`

`evalml.objectives.LogLoss.predict`

`LogLoss.predict` (*y_predicted*, *extra_cols=None*)

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called

Parameters **`y_predicted`** – the prediction to transform to final prediction

Returns predictions

`evalml.objectives.LogLoss.score`

`LogLoss.score` (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- **`y_predicted`** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **`y_true`** (*list*) – the ground truth for the predictions.
- **`extra_cols`** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

`evalml.objectives.LogLoss.supports_problem_type`

classmethod `LogLoss.supports_problem_type` (*problem_type*)

Checks if objective supports given `ProblemType`

Parameters **`problem_type`** (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

evalml.objectives.MCC

class evalml.objectives.MCC (*verbose=False*)
 Matthews correlation coefficient for both binary and multiclass classification

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.MCC.__init__

MCC.**__init__** (*verbose=False*)
 Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.MCC.fit

MCC.**fit** (*y_predicted, y_true, extra_cols=None*)
 Learn the objective function based on the predictions from a model.
 If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.MCC.predict

MCC.**predict** (*y_predicted, extra_cols=None*)
 Apply the learned objective function to the output of a model.
 If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.MCC.score**MCC.score** (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to `True`**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is `True`, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is `True`.

Returns score**evalml.objectives.MCC.supports_problem_type****classmethod** **MCC.supports_problem_type** (*problem_type*)Checks if objective supports given `ProblemType`**Parameters** **problem_type** (*str or ProblemType*) – problem type to check**Returns** whether objective supports `ProblemType`**Return type** bool**evalml.objectives.ROC****class** **evalml.objectives.ROC** (*verbose=False*)

Receiver Operating Characteristic score for binary classification.

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given <code>ProblemType</code>

evalml.objectives.ROC.__init__**ROC.__init__** (*verbose=False*)Initialize self. See `help(type(self))` for accurate signature.

evalml.objectives.ROC.fit**ROC.fit** (*y_predicted*, *y_true*, *extra_cols=None*)

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won't be called**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns `self`**evalml.objectives.ROC.predict****ROC.predict** (*y_predicted*, *extra_cols=None*)

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called**Parameters** **y_predicted** – the prediction to transform to final prediction**Returns** predictions**evalml.objectives.ROC.score****ROC.score** (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score**evalml.objectives.ROC.supports_problem_type****classmethod** **ROC.supports_problem_type** (*problem_type*)Checks if objective supports given `ProblemType`**Parameters** **problem_type** (*str or ProblemType*) – problem type to check**Returns** whether objective supports `ProblemType`**Return type** bool

evalml.objectives.ConfusionMatrix

class evalml.objectives.**ConfusionMatrix** (*verbose=False*)
Confusion matrix for classification problems

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.ConfusionMatrix.__init__

ConfusionMatrix.**__init__** (*verbose=False*)
Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.ConfusionMatrix.fit

ConfusionMatrix.**fit** (*y_predicted, y_true, extra_cols=None*)
Learn the objective function based on the predictions from a model.
If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.ConfusionMatrix.predict

ConfusionMatrix.**predict** (*y_predicted, extra_cols=None*)
Apply the learned objective function to the output of a model.
If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.ConfusionMatrix.score

`ConfusionMatrix.score(y_predicted, y_true)`

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to `True`

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is `True`, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is `True`.

Returns score

evalml.objectives.ConfusionMatrix.supports_problem_type

classmethod `ConfusionMatrix.supports_problem_type(problem_type)`

Checks if objective supports given `ProblemType`

Parameters **problem_type** (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

Regression

<i>R2</i>	Coefficient of determination for regression
<i>MAE</i>	Mean absolute error for regression
<i>MSE</i>	Mean squared error for regression
<i>MSLE</i>	Mean squared log error for regression
<i>MedianAE</i>	Median absolute error for regression
<i>MaxError</i>	Maximum residual error for regression
<i>ExpVariance</i>	Explained variance score for regression

evalml.objectives.R2

class `evalml.objectives.R2(verbose=False)`

Coefficient of determination for regression

Methods

<i>__init__</i>	Initialize self.
<i>fit</i>	Learn the objective function based on the predictions from a model.
<i>predict</i>	Apply the learned objective function to the output of a model.

Continued on next page

Table 76 – continued from previous page

<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

`evalml.objectives.R2.__init__`

`R2.__init__` (*verbose=False*)

Initialize self. See `help(type(self))` for accurate signature.

`evalml.objectives.R2.fit`

`R2.fit` (*y_predicted, y_true, extra_cols=None*)

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won't be called

Parameters

- **`y_predicted`** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **`y_true`** (*list*) – the ground truth for the predictions.
- **`extra_cols`** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns self

`evalml.objectives.R2.predict`

`R2.predict` (*y_predicted, extra_cols=None*)

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called

Parameters **`y_predicted`** – the prediction to transform to final prediction

Returns predictions

`evalml.objectives.R2.score`

`R2.score` (*y_predicted, y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- **`y_predicted`** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **`y_true`** (*list*) – the ground truth for the predictions.
- **`extra_cols`** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

evalml.objectives.R2.supports_problem_type**classmethod** `R2.supports_problem_type(problem_type)`

Checks if objective supports given ProblemType

Parameters `problem_type` (*str* or *ProblemType*) – problem type to check**Returns** whether objective supports ProblemType**Return type** bool**evalml.objectives.MAE****class** `evalml.objectives.MAE(verbose=False)`

Mean absolute error for regression

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.MAE.__init__`MAE.__init__(verbose=False)`Initialize self. See `help(type(self))` for accurate signature.**evalml.objectives.MAE.fit**`MAE.fit(y_predicted, y_true, extra_cols=None)`

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won't be called**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns self

evalml.objectives.MAE.predict

MAE.predict (*y_predicted*, *extra_cols=None*)

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.MAE.score

MAE.score (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

evalml.objectives.MAE.supports_problem_type

classmethod **MAE.supports_problem_type** (*problem_type*)

Checks if objective supports given ProblemType

Parameters **problem_type** (*str or ProblemType*) – problem type to check

Returns whether objective supports ProblemType

Return type bool

evalml.objectives.MSE

class **evalml.objectives.MSE** (*verbose=False*)

Mean squared error for regression

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.

Continued on next page

Table 78 – continued from previous page

<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.MSE.__init__**MSE.__init__** (*verbose=False*)

Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.MSE.fit**MSE.fit** (*y_predicted, y_true, extra_cols=None*)

Learn the objective function based on the predictions from a model.

If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self**evalml.objectives.MSE.predict****MSE.predict** (*y_predicted, extra_cols=None*)

Apply the learned objective function to the output of a model.

If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction**Returns** predictions**evalml.objectives.MSE.score****MSE.score** (*y_predicted, y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set greater_is_better attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns score

evalml.objectives.MSE.supports_problem_type**classmethod** `MSE.supports_problem_type(problem_type)`

Checks if objective supports given ProblemType

Parameters `problem_type` (*str or ProblemType*) – problem type to check**Returns** whether objective supports ProblemType**Return type** bool**evalml.objectives.MSLE****class** `evalml.objectives.MSLE(verbose=False)`

Mean squared log error for regression

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.MSLE.__init__`MSLE.__init__(verbose=False)`Initialize self. See `help(type(self))` for accurate signature.**evalml.objectives.MSLE.fit**`MSLE.fit(y_predicted, y_true, extra_cols=None)`

Learn the objective function based on the predictions from a model.

If `needs_fitting` is false, this method won't be called**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns self

evalml.objectives.MSLE.predict**MSLE.predict** (*y_predicted*, *extra_cols=None*)

Apply the learned objective function to the output of a model.

If `needs_fitting` is `false`, this method won't be called**Parameters** `y_predicted` – the prediction to transform to final prediction**Returns** predictions**evalml.objectives.MSLE.score****MSLE.score** (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to `True`**Parameters**

- **y_predicted** (*list*) – the predictions from the model. If `needs_proba` is `True`, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is `True`.

Returns score**evalml.objectives.MSLE.supports_problem_type****classmethod** **MSLE.supports_problem_type** (*problem_type*)Checks if objective supports given `ProblemType`**Parameters** `problem_type` (*str or ProblemType*) – problem type to check**Returns** whether objective supports `ProblemType`**Return type** `bool`**evalml.objectives.MedianAE****class** **evalml.objectives.MedianAE** (*verbose=False*)

Median absolute error for regression

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.

Continued on next page

Table 80 – continued from previous page

<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.MedianAE.__init__

MedianAE.**__init__** (*verbose=False*)

Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.MedianAE.fit

MedianAE.**fit** (*y_predicted, y_true, extra_cols=None*)

Learn the objective function based on the predictions from a model.

If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.MedianAE.predict

MedianAE.**predict** (*y_predicted, extra_cols=None*)

Apply the learned objective function to the output of a model.

If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.MedianAE.score

MedianAE.**score** (*y_predicted, y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set greater_is_better attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns score

evalml.objectives.MedianAE.supports_problem_type**classmethod** MedianAE.supports_problem_type(*problem_type*)

Checks if objective supports given ProblemType

Parameters *problem_type* (*str* or *ProblemType*) – problem type to check**Returns** whether objective supports ProblemType**Return type** bool**evalml.objectives.MaxError****class** evalml.objectives.MaxError(*verbose=False*)

Maximum residual error for regression

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.
<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.MaxError.__init__MaxError.__init__(*verbose=False*)

Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.MaxError.fitMaxError.fit(*y_predicted*, *y_true*, *extra_cols=None*)

Learn the objective function based on the predictions from a model.

If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.MaxError.predict

`MaxError.predict` (*y_predicted*, *extra_cols=None*)

Apply the learned objective function to the output of a model.

If `needs_fitting` is false, this method won't be called

Parameters `y_predicted` – the prediction to transform to final prediction

Returns predictions

evalml.objectives.MaxError.score

`MaxError.score` (*y_predicted*, *y_true*)

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set `greater_is_better` attribute to True

Parameters

- `y_predicted` (*list*) – the predictions from the model. If `needs_proba` is True, it is the probability estimates
- `y_true` (*list*) – the ground truth for the predictions.
- `extra_cols` (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if `uses_extra_columns` is True.

Returns score

evalml.objectives.MaxError.supports_problem_type

classmethod `MaxError.supports_problem_type` (*problem_type*)

Checks if objective supports given `ProblemType`

Parameters `problem_type` (*str or ProblemType*) – problem type to check

Returns whether objective supports `ProblemType`

Return type bool

evalml.objectives.ExpVariance

class `evalml.objectives.ExpVariance` (*verbose=False*)

Explained variance score for regression

Methods

<code>__init__</code>	Initialize self.
<code>fit</code>	Learn the objective function based on the predictions from a model.
<code>predict</code>	Apply the learned objective function to the output of a model.

Continued on next page

Table 82 – continued from previous page

<code>score</code>	Calculate score from applying fitted objective to predicted values
<code>supports_problem_type</code>	Checks if objective supports given ProblemType

evalml.objectives.ExpVariance.__init__

`ExpVariance.__init__(verbose=False)`

Initialize self. See help(type(self)) for accurate signature.

evalml.objectives.ExpVariance.fit

`ExpVariance.fit(y_predicted, y_true, extra_cols=None)`

Learn the objective function based on the predictions from a model.

If needs_fitting is false, this method won't be called

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns self

evalml.objectives.ExpVariance.predict

`ExpVariance.predict(y_predicted, extra_cols=None)`

Apply the learned objective function to the output of a model.

If needs_fitting is false, this method won't be called

Parameters **y_predicted** – the prediction to transform to final prediction

Returns predictions

evalml.objectives.ExpVariance.score

`ExpVariance.score(y_predicted, y_true)`

Calculate score from applying fitted objective to predicted values

If a higher score is better than a lower score, set greater_is_better attribute to True

Parameters

- **y_predicted** (*list*) – the predictions from the model. If needs_proba is True, it is the probability estimates
- **y_true** (*list*) – the ground truth for the predictions.
- **extra_cols** (*pd.DataFrame*) – any extra columns that are needed from training data to fit. Only provided if uses_extra_columns is True.

Returns score

evalml.objectives.ExpVariance.supports_problem_type**classmethod** `ExpVariance.supports_problem_type(problem_type)`

Checks if objective supports given ProblemType

Parameters `problem_type` (*str* or *ProblemType*) – problem type to check**Returns** whether objective supports ProblemType**Return type** bool**Problem Types**

<i>ProblemTypes</i>	Enum for type of machine learning problem: BINARY, MULTICLASS, or REGRESSION
---------------------	--

evalml.problem_types.ProblemTypes**class** `evalml.problem_types.ProblemTypes`

Enum for type of machine learning problem: BINARY, MULTICLASS, or REGRESSION

<i>handle_problem_types</i>	Handles problem_type by either returning the ProblemTypes or converting from a str
-----------------------------	--

evalml.problem_types.handle_problem_types`evalml.problem_types.handle_problem_types(problem_type)`

Handles problem_type by either returning the ProblemTypes or converting from a str

Parameters `problem_types` (*str* or *ProblemTypes*) – problem type that needs to be handled**Returns** ProblemTypes**Tuners**

<i>Tuner</i>	Defines API for Tuners
<i>SKOptTuner</i>	Bayesian Optimizer
<i>GridSearchTuner</i>	Grid Search Optimizer
<i>RandomSearchTuner</i>	Random Search Optimizer

evalml.tuners.Tuner**class** `evalml.tuners.Tuner(space, random_state=0)`

Defines API for Tuners

Tuners implement different strategies for sampling from a search space. They're used in EvalML to search the space of pipeline hyperparameters.

Methods

<code>__init__</code>	Init Tuner
<code>add</code>	Register a set of hyperparameters with the score obtained from training a pipeline with those hyperparameters.
<code>is_search_space_exhausted</code>	Optional.
<code>propose</code>	Returns a set of hyperparameters to train a pipeline with, based off the search space dimensions and prior samples

`evalml.tuners.Tuner.__init__`

`Tuner.__init__(space, random_state=0)`

Init Tuner

Parameters

- **space** (*dict*) – search space for hyperparameters
- **random_state** (*int*, *np.random.RandomState*) – The random state

Returns `self`

Return type *Tuner*

`evalml.tuners.Tuner.add`

`Tuner.add(parameters, score)`

Register a set of hyperparameters with the score obtained from training a pipeline with those hyperparameters.

Parameters

- **parameters** (*dict*) – hyperparameters
- **score** (*float*) – associated score

Returns `None`

`evalml.tuners.Tuner.is_search_space_exhausted`

`Tuner.is_search_space_exhausted()`

Optional. If possible search space for tuner is finite, this method indicates whether or not all possible parameters have been scored.

Returns Returns true if all possible parameters in a search space has been scored.

Return type `bool`

`evalml.tuners.Tuner.propose`

`Tuner.propose()`

Returns a set of hyperparameters to train a pipeline with, based off the search space dimensions and prior samples

Returns proposed hyperparameters

Return type dict

evalml.tuners.SKOptTuner

class evalml.tuners.SKOptTuner(*space*, *random_state*=0)
Bayesian Optimizer

Methods

<code>__init__</code>	Init SKOptTuner
<code>add</code>	Add score to sample
<code>is_search_space_exhausted</code>	Optional.
<code>propose</code>	Returns hyperparameters based off search space and samples

evalml.tuners.SKOptTuner.__init__

SKOptTuner.__init__(*space*, *random_state*=0)
Init SKOptTuner

Parameters

- **space** (*dict*) – search space for hyperparameters
- **random_state** (*int*, *np.random.RandomState*) – The random state

Returns self

Return type SKOptTuner

evalml.tuners.SKOptTuner.add

SKOptTuner.add(*parameters*, *score*)
Add score to sample

Parameters

- **parameters** (*dict*) – hyperparameters
- **score** (*float*) – associated score

Returns None

evalml.tuners.SKOptTuner.is_search_space_exhausted

SKOptTuner.is_search_space_exhausted()
Optional. If possible search space for tuner is finite, this method indicates whether or not all possible parameters have been scored.

Returns Returns true if all possible parameters in a search space has been scored.

Return type bool

evalml.tuners.SKOptTuner.propose`SKOptTuner.propose()`

Returns hyperparameters based off search space and samples

Returns proposed hyperparameters**Return type** dict**evalml.tuners.GridSearchTuner**

class evalml.tuners.**GridSearchTuner** (*space, n_points=10, random_state=0*)
 Grid Search Optimizer

Example

```
>>> tuner = GridSearchTuner([(1,10), ['A', 'B']], n_points=5)
>>> print(tuner.propose())
(1.0, 'A')
>>> print(tuner.propose())
(1.0, 'B')
>>> print(tuner.propose())
(3.25, 'A')
```

Methods

<code>__init__</code>	Generate all of the possible points to search for in the grid
<code>add</code>	Not applicable to grid search tuner as generated parameters are not dependent on scores of previous parameters.
<code>is_search_space_exhausted</code>	Checks if it is possible to generate a set of valid parameters.
<code>propose</code>	Returns hyperparameters from <code>_grid_points</code> iterations

evalml.tuners.GridSearchTuner.__init__

`GridSearchTuner.__init__` (*space, n_points=10, random_state=0*)
 Generate all of the possible points to search for in the grid

Parameters

- **space** – A list of all dimensions available to tune
- **n_points** – The number of points to sample from along each dimension defined in the `space` argument
- **random_state** – Unused in this class

`evalml.tuners.GridSearchTuner.add`

`GridSearchTuner.add(parameters, score)`

Not applicable to grid search tuner as generated parameters are not dependent on scores of previous parameters.

Parameters

- **parameters** – Hyperparameters used
- **score** – Associated score

`evalml.tuners.GridSearchTuner.is_search_space_exhausted`

`GridSearchTuner.is_search_space_exhausted()`

Checks if it is possible to generate a set of valid parameters. Stores generated parameters in `self.curr_params` to be returned by `propose()`.

Raises `NoParamsException` – If a search space is exhausted, then this exception is thrown.

Returns If no more valid parameters exists in the search space, return false.

Return type bool

`evalml.tuners.GridSearchTuner.propose`

`GridSearchTuner.propose()`

Returns hyperparameters from `_grid_points` iterations

If all possible combinations of parameters have been scored, then `NoParamsException` is raised.

Returns proposed hyperparameters

Return type dict

`evalml.tuners.RandomSearchTuner`

class `evalml.tuners.RandomSearchTuner` (*space, random_state=0, with_replacement=False, replacement_max_attempts=10*)

Random Search Optimizer

Example

```
>>> tuner = RandomSearchTuner([(1,10), ['A', 'B']], random_state=0)
>>> print(tuner.propose())
(6, 'B')
>>> print(tuner.propose())
(4, 'B')
>>> print(tuner.propose())
(5, 'A')
```

Methods

<code>__init__</code>	Sets up check for duplication if needed.
<code>add</code>	Not applicable to random search tuner as generated parameters are not dependent on scores of previous parameters.
<code>is_search_space_exhausted</code>	Checks if it is possible to generate a set of valid parameters.
<code>propose</code>	Generate a unique set of parameters.

`evalml.tuners.RandomSearchTuner.__init__`

`RandomSearchTuner.__init__(space, random_state=0, with_replacement=False, replacement_max_attempts=10)`

Sets up check for duplication if needed.

Parameters

- **space** – A list of all dimensions available to tune
- **random_state** – Unused in this class
- **with_replacement** – If false, only unique hyperparameters will be shown
- **replacement_max_attempts** – The maximum number of tries to get a unique set of random parameters. Only used if tuner is initialized with `with_replacement=True`

`evalml.tuners.RandomSearchTuner.add`

`RandomSearchTuner.add(parameters, score)`

Not applicable to random search tuner as generated parameters are not dependent on scores of previous parameters.

Parameters

- **parameters** – Hyperparameters used
- **score** – Associated score

`evalml.tuners.RandomSearchTuner.is_search_space_exhausted`

`RandomSearchTuner.is_search_space_exhausted()`

Checks if it is possible to generate a set of valid parameters. Stores generated parameters in `self.curr_params` to be returned by `propose()`.

Raises `NoParamsException` – If a search space is exhausted, then this exception is thrown.

Returns If no more valid parameters exists in the search space, return false.

Return type bool

`evalml.tuners.RandomSearchTuner.propose`

`RandomSearchTuner.propose()`

Generate a unique set of parameters.

If tuner was initialized with `with_replacement=True` and the tuner is unable to generate a unique set of parameters after `replacement_max_attempts` tries, then `NoParamsException` is raised.

Returns A list of unique parameters

Guardrails

<code>detect_highly_null</code>	Checks if there are any highly-null columns in a dataframe.
<code>detect_label_leakage</code>	Check if any of the features are highly correlated with the target.
<code>detect_outliers</code>	Checks if there are any outliers in a dataframe by using first Isolation Forest to obtain the anomaly score of each index and then using IQR to determine score anomalies.
<code>detect_id_columns</code>	Check if any of the features are ID columns.

evalml.guardrails.detect_highly_null

`evalml.guardrails.detect_highly_null(X, percent_threshold=0.95)`

Checks if there are any highly-null columns in a dataframe.

Parameters

- **X** (`pd.DataFrame`) – features
- **percent_threshold** (`float`) – Require that percentage of null values to be considered “highly-null”, defaults to .95

Returns A dictionary of features with column name or index and their percentage of null values

Example

```
>>> df = pd.DataFrame({
...     'lots_of_null': [None, None, None, None, 5],
...     'no_null': [1, 2, 3, 4, 5]
... })
>>> detect_highly_null(df, percent_threshold=0.8)
{'lots_of_null': 0.8}
```

evalml.guardrails.detect_label_leakage

`evalml.guardrails.detect_label_leakage(X, y, threshold=0.95)`

Check if any of the features are highly correlated with the target.

Currently only supports binary and numeric targets and features

Parameters

- **X** (`pd.DataFrame`) – The input features to check
- **y** (`pd.Series`) – the labels
- **threshold** (`float`) – the correlation threshold to be considered leakage. Defaults to .95

Returns leakage, dictionary of features with leakage and corresponding threshold

Example

```
>>> X = pd.DataFrame({
...     'leak': [10, 42, 31, 51, 61],
...     'x': [42, 54, 12, 64, 12],
...     'y': [12, 5, 13, 74, 24],
... })
>>> y = pd.Series([10, 42, 31, 51, 40])
>>> detect_label_leakage(X, y, threshold=0.8)
{'leak': 0.8827072320669518}
```

evalml.guardrails.detect_outliers

`evalml.guardrails.detect_outliers(X, random_state=0)`

Checks if there are any outliers in a dataframe by using first Isolation Forest to obtain the anomaly score of each index and then using IQR to determine score anomalies. Indices with score anomalies are considered outliers.

Parameters *X* (`pd.DataFrame`) – features

Returns A set of indices that may have outlier data.

Example

```
>>> df = pd.DataFrame({
...     'x': [1, 2, 3, 40, 5],
...     'y': [6, 7, 8, 990, 10],
...     'z': [-1, -2, -3, -1201, -4]
... })
>>> detect_outliers(df)
[3]
```

evalml.guardrails.detect_id_columns

`evalml.guardrails.detect_id_columns(X, threshold=1.0)`

Check if any of the features are ID columns. Currently performs these simple checks:

- column name is “id”
- column name ends in “_id”
- column contains all unique values (and is not float / boolean)

Parameters

- *X* (`pd.DataFrame`) – The input features to check
- **threshold** (`float`) – the probability threshold to be considered an ID column. Defaults to 1.0

Returns A dictionary of features with column name or index and their probability of being ID columns

Example

```
>>> df = pd.DataFrame({
...     'df_id': [0, 1, 2, 3, 4],
...     'x': [10, 42, 31, 51, 61],
...     'y': [42, 54, 12, 64, 12]
... })
>>> detect_id_columns(df)
{'df_id': 1.0}
```

2.6.15 FAQ

What is the difference between EvalML and other AutoML libraries?

EvalML optimizes machine learning pipelines on *custom practical objectives* instead of vague machine learning loss functions so that it will find the best pipelines for your specific needs. Furthermore, EvalML *pipelines* are able to take in all kinds of data (missing values, categorical, etc.) as long as the data are in a single table. EvalML also allows you to build your own pipelines with existing or custom components so you can have more control over the AutoML process. Moreover, EvalML also provides you with support in the form of *guardrails* to ensure that you are aware of potential issues your data may cause with machine learning algorithms”.

How does EvalML handle missing values?

EvalML contains imputation components in its pipelines so that missing values are taken care of. EvalML optimizes over different types of imputation to search for the best possible pipeline. You can find more information about components [here](#) and in the API reference [here](#).

How does EvalML handle categorical encoding?

EvalML provides a *one-hot-encoding component* in its pipelines for categorical variables. EvalML plans to support other encoders in the future.

How does EvalML handle feature selection?

EvalML currently utilizes scikit-learn’s [SelectFromModel](#) with a Random Forest classifier/regressor to handle feature selection. EvalML plans on supporting more feature selectors in the future. You can find more information in the API reference [here](#).

How are feature importances calculated?

Feature importance depends on the estimator used. Variable coefficients are used for regression-based estimators (Logistic Regression and Linear Regression) and Gini importance is used for tree-based estimators (Random Forest and XGBoost).

How does hyperparameter tuning work?

EvalML tunes hyperparameters for its pipelines through Bayesian optimization. In the future we plan to support more optimization techniques such as random search.

Can I create my own objective metric?

Yes you can! You can *create your own custom objective* so that EvalML optimizes the best model for your needs.

How does EvalML avoid overfitting?

EvalML provides *guardrails* to combat overfitting. Such guardrails include detecting label leakage, unstable pipelines, hold-out datasets and cross validation. EvalML defaults to using Stratified K-Fold cross-validation for classification problems and K-Fold cross-validation for regression problems but allows you to utilize your own cross-validation methods as well.

Can I create my own pipeline for EvalML?

Yes! EvalML allows you to create *custom pipelines* using modular components. This allows you to customize EvalML pipelines for your own needs or for AutoML.

Does EvalML work with X algorithm?

EvalML is constantly improving and adding new components and will allow your own algorithms to be used as components in our pipelines.

Symbols

`__init__()` (*evalml.AutoClassificationSearch method*), 59
`__init__()` (*evalml.AutoRegressionSearch method*), 60
`__init__()` (*evalml.objectives.AUC method*), 131
`__init__()` (*evalml.objectives.AUCMacro method*), 134
`__init__()` (*evalml.objectives.AUCMicro method*), 133
`__init__()` (*evalml.objectives.AUCWeighted method*), 136
`__init__()` (*evalml.objectives.ConfusionMatrix method*), 142
`__init__()` (*evalml.objectives.ExpVariance method*), 153
`__init__()` (*evalml.objectives.F1 method*), 113
`__init__()` (*evalml.objectives.F1Macro method*), 116
`__init__()` (*evalml.objectives.F1Micro method*), 115
`__init__()` (*evalml.objectives.F1Weighted method*), 118
`__init__()` (*evalml.objectives.FraudCost method*), 109
`__init__()` (*evalml.objectives.LeadScoring method*), 111
`__init__()` (*evalml.objectives.LogLoss method*), 137
`__init__()` (*evalml.objectives.MAE method*), 145
`__init__()` (*evalml.objectives.MCC method*), 139
`__init__()` (*evalml.objectives.MSE method*), 147
`__init__()` (*evalml.objectives.MSLE method*), 148
`__init__()` (*evalml.objectives.MaxError method*), 151
`__init__()` (*evalml.objectives.MedianAE method*), 150
`__init__()` (*evalml.objectives.Precision method*), 119
`__init__()` (*evalml.objectives.PrecisionMacro method*), 122
`__init__()` (*evalml.objectives.PrecisionMicro method*), 121
`__init__()` (*evalml.objectives.PrecisionWeighted method*), 124
`__init__()` (*evalml.objectives.R2 method*), 144
`__init__()` (*evalml.objectives.ROC method*), 140
`__init__()` (*evalml.objectives.Recall method*), 125
`__init__()` (*evalml.objectives.RecallMacro method*), 128
`__init__()` (*evalml.objectives.RecallMicro method*), 127
`__init__()` (*evalml.objectives.RecallWeighted method*), 130
`__init__()` (*evalml.pipelines.CatBoostClassificationPipeline method*), 91
`__init__()` (*evalml.pipelines.CatBoostRegressionPipeline method*), 101
`__init__()` (*evalml.pipelines.LinearRegressionPipeline method*), 104
`__init__()` (*evalml.pipelines.LogisticRegressionPipeline method*), 94
`__init__()` (*evalml.pipelines.PipelineBase method*), 81
`__init__()` (*evalml.pipelines.RFClassificationPipeline method*), 84
`__init__()` (*evalml.pipelines.RFRegressionPipeline method*), 98
`__init__()` (*evalml.pipelines.XGBoostPipeline method*), 88
`__init__()` (*evalml.pipelines.components.LinearRegressor method*), 77
`__init__()` (*evalml.pipelines.components.LogisticRegressionClassifier method*), 73
`__init__()` (*evalml.pipelines.components.OneHotEncoder method*), 63
`__init__()` (*evalml.pipelines.components.RFClassifierSelectFromModel method*), 67
`__init__()` (*evalml.pipelines.components.RFRegressorSelectFromModel method*), 65
`__init__()` (*evalml.pipelines.components.RandomForestClassifier method*), 74
`__init__()` (*evalml.pipelines.components.RandomForestRegressor method*), 79
`__init__()` (*evalml.pipelines.components.SimpleImputer method*), 69
`__init__()` (*evalml.pipelines.components.StandardScaler method*), 71

`__init__()` (*evalml.pipelines.components.XGBoostClassifier* custom_hyperparameters method), 76

`__init__()` (*evalml.tuners.GridSearchTuner* method), 157

`__init__()` (*evalml.tuners.RandomSearchTuner* method), 159

`__init__()` (*evalml.tuners.SKOptTuner* method), 156

`__init__()` (*evalml.tuners.Tuner* method), 155

`custom_hyperparameters` (*evalml.pipelines.RFClassificationPipeline* attribute), 84

`custom_hyperparameters` (*evalml.pipelines.RFRegressionPipeline* attribute), 97

`custom_hyperparameters` (*evalml.pipelines.XGBoostPipeline* attribute), 87

A

`add()` (*evalml.tuners.GridSearchTuner* method), 158

`add()` (*evalml.tuners.RandomSearchTuner* method), 159

`add()` (*evalml.tuners.SKOptTuner* method), 156

`add()` (*evalml.tuners.Tuner* method), 155

`AUC` (class in *evalml.objectives*), 131

`AUCMacro` (class in *evalml.objectives*), 134

`AUCMicro` (class in *evalml.objectives*), 133

`AUCWeighted` (class in *evalml.objectives*), 136

`AutoClassificationSearch` (class in *evalml*), 58

`AutoRegressionSearch` (class in *evalml*), 60

C

`CatBoostClassificationPipeline` (class in *evalml.pipelines*), 90

`CatBoostRegressionPipeline` (class in *evalml.pipelines*), 100

`component_graph` (*evalml.pipelines.CatBoostClassificationPipeline* attribute), 90

`component_graph` (*evalml.pipelines.CatBoostRegressionPipeline* attribute), 100

`component_graph` (*evalml.pipelines.LinearRegressionPipeline* attribute), 104

`component_graph` (*evalml.pipelines.LogisticRegressionPipeline* attribute), 94

`component_graph` (*evalml.pipelines.RFClassificationPipeline* attribute), 84

`component_graph` (*evalml.pipelines.RFRegressionPipeline* attribute), 97

`component_graph` (*evalml.pipelines.XGBoostPipeline* attribute), 87

`ConfusionMatrix` (class in *evalml.objectives*), 142

`custom_hyperparameters` (*evalml.pipelines.CatBoostClassificationPipeline* attribute), 91

`custom_hyperparameters` (*evalml.pipelines.CatBoostRegressionPipeline* attribute), 101

`custom_hyperparameters` (*evalml.pipelines.LinearRegressionPipeline* attribute), 104

`custom_hyperparameters` (*evalml.pipelines.LogisticRegressionPipeline* attribute), 94

D

`decision_function()` (*evalml.objectives.FraudCost* method), 109

`decision_function()` (*evalml.objectives.LeadScoring* method), 111

`describe()` (*evalml.pipelines.CatBoostClassificationPipeline* method), 92

`describe()` (*evalml.pipelines.CatBoostRegressionPipeline* method), 102

`describe()` (*evalml.pipelines.components.LinearRegressor* method), 78

`describe()` (*evalml.pipelines.components.LogisticRegressionClassifier* method), 73

`describe()` (*evalml.pipelines.components.OneHotEncoder* method), 63

`describe()` (*evalml.pipelines.components.RandomForestClassifier* method), 74

`describe()` (*evalml.pipelines.components.RandomForestRegressor* method), 79

`describe()` (*evalml.pipelines.components.RFClassifierSelectFromModel* method), 67

`describe()` (*evalml.pipelines.components.RFRegressorSelectFromModel* method), 65

`describe()` (*evalml.pipelines.components.SimpleImputer* method), 69

`describe()` (*evalml.pipelines.components.StandardScaler* method), 71

`describe()` (*evalml.pipelines.components.XGBoostClassifier* method), 76

`describe()` (*evalml.pipelines.LinearRegressionPipeline* method), 105

`describe()` (*evalml.pipelines.LogisticRegressionPipeline* method), 95

`describe()` (*evalml.pipelines.PipelineBase* method), 81

`describe()` (*evalml.pipelines.RFClassificationPipeline* method), 85

`describe()` (*evalml.pipelines.RFRegressionPipeline* method), 98

`describe()` (*evalml.pipelines.XGBoostPipeline* method), 88

`detect_highly_null()` (in module *evalml.guardrails*), 160

`detect_id_columns()` (in module `evalml.guardrails`), 161
`detect_label_leakage()` (in module `evalml.guardrails`), 160
`detect_outliers()` (in module `evalml.guardrails`), 161

E

`ExpVariance` (class in `evalml.objectives`), 152

F

`F1` (class in `evalml.objectives`), 113
`F1Macro` (class in `evalml.objectives`), 116
`F1Micro` (class in `evalml.objectives`), 115
`F1Weighted` (class in `evalml.objectives`), 118
`feature_importance_graph()` (`evalml.pipelines.CatBoostClassificationPipeline` method), 92
`feature_importance_graph()` (`evalml.pipelines.CatBoostRegressionPipeline` method), 102
`feature_importance_graph()` (`evalml.pipelines.LinearRegressionPipeline` method), 105
`feature_importance_graph()` (`evalml.pipelines.LogisticRegressionPipeline` method), 95
`feature_importance_graph()` (`evalml.pipelines.PipelineBase` method), 81, 108
`feature_importance_graph()` (`evalml.pipelines.RFClassificationPipeline` method), 85
`feature_importance_graph()` (`evalml.pipelines.RFRegressionPipeline` method), 98
`feature_importance_graph()` (`evalml.pipelines.XGBoostPipeline` method), 88
`fit()` (`evalml.objectives.AUC` method), 132
`fit()` (`evalml.objectives.AUCMacro` method), 135
`fit()` (`evalml.objectives.AUCMicro` method), 133
`fit()` (`evalml.objectives.AUCWeighted` method), 136
`fit()` (`evalml.objectives.ConfusionMatrix` method), 142
`fit()` (`evalml.objectives.ExpVariance` method), 153
`fit()` (`evalml.objectives.F1` method), 114
`fit()` (`evalml.objectives.F1Macro` method), 117
`fit()` (`evalml.objectives.F1Micro` method), 115
`fit()` (`evalml.objectives.F1Weighted` method), 118
`fit()` (`evalml.objectives.FraudCost` method), 109
`fit()` (`evalml.objectives.LeadScoring` method), 111
`fit()` (`evalml.objectives.LogLoss` method), 138
`fit()` (`evalml.objectives.MAE` method), 145
`fit()` (`evalml.objectives.MaxError` method), 151

`fit()` (`evalml.objectives.MCC` method), 139
`fit()` (`evalml.objectives.MedianAE` method), 150
`fit()` (`evalml.objectives.MSE` method), 147
`fit()` (`evalml.objectives.MSLE` method), 148
`fit()` (`evalml.objectives.Precision` method), 120
`fit()` (`evalml.objectives.PrecisionMacro` method), 123
`fit()` (`evalml.objectives.PrecisionMicro` method), 121
`fit()` (`evalml.objectives.PrecisionWeighted` method), 124
`fit()` (`evalml.objectives.R2` method), 144
`fit()` (`evalml.objectives.Recall` method), 126
`fit()` (`evalml.objectives.RecallMacro` method), 129
`fit()` (`evalml.objectives.RecallMicro` method), 127
`fit()` (`evalml.objectives.RecallWeighted` method), 130
`fit()` (`evalml.objectives.ROC` method), 141
`fit()` (`evalml.pipelines.CatBoostClassificationPipeline` method), 92
`fit()` (`evalml.pipelines.CatBoostRegressionPipeline` method), 102
`fit()` (`evalml.pipelines.components.LinearRegressor` method), 78
`fit()` (`evalml.pipelines.components.LogisticRegressionClassifier` method), 73
`fit()` (`evalml.pipelines.components.OneHotEncoder` method), 64
`fit()` (`evalml.pipelines.components.RandomForestClassifier` method), 75
`fit()` (`evalml.pipelines.components.RandomForestRegressor` method), 80
`fit()` (`evalml.pipelines.components.RFClassifierSelectFromModel` method), 68
`fit()` (`evalml.pipelines.components.RFRegressorSelectFromModel` method), 66
`fit()` (`evalml.pipelines.components.SimpleImputer` method), 70
`fit()` (`evalml.pipelines.components.StandardScaler` method), 71
`fit()` (`evalml.pipelines.components.XGBoostClassifier` method), 76
`fit()` (`evalml.pipelines.LinearRegressionPipeline` method), 105
`fit()` (`evalml.pipelines.LogisticRegressionPipeline` method), 95
`fit()` (`evalml.pipelines.PipelineBase` method), 82
`fit()` (`evalml.pipelines.RFClassificationPipeline` method), 85
`fit()` (`evalml.pipelines.RFRegressionPipeline` method), 99
`fit()` (`evalml.pipelines.XGBoostPipeline` method), 89
`fit_transform()` (`evalml.pipelines.components.OneHotEncoder` method), 64
`fit_transform()` (`evalml.pipelines.components.RFClassifierSelectFromModel` method), 68
`fit_transform()` (`evalml.pipelines.components.RFRegressorSelectFromModel` method), 66

method), 66
 fit_transform() (*evalml.pipelines.components.SimpleImputer*
method), 70
 fit_transform() (*evalml.pipelines.components.StandardScaler*
method), 72
 FraudCost (*class in evalml.objectives*), 108

G

generate_confusion_matrix()
 (*evalml.AutoClassificationSearch.plot method*),
 62
 generate_roc_plot()
 (*evalml.AutoClassificationSearch.plot method*),
 62
 get_component() (*evalml.pipelines.CatBoostClassificationPipeline*
method), 92
 get_component() (*evalml.pipelines.CatBoostRegressionPipeline*
method), 102
 get_component() (*evalml.pipelines.LinearRegressionPipeline*
method), 106
 get_component() (*evalml.pipelines.LogisticRegressionPipeline*
method), 96
 get_component() (*evalml.pipelines.PipelineBase*
method), 82
 get_component() (*evalml.pipelines.RFClassificationPipeline*
method), 86
 get_component() (*evalml.pipelines.RFRegressionPipeline*
method), 99
 get_component() (*evalml.pipelines.XGBoostPipeline*
method), 89
 get_confusion_matrix_data()
 (*evalml.AutoClassificationSearch.plot method*),
 62
 get_feature_names()
 (*evalml.pipelines.components.OneHotEncoder*
method), 64
 get_indices() (*evalml.pipelines.components.RFClassifierSelectFromModel*
method), 68
 get_indices() (*evalml.pipelines.components.RFRegressorSelectFromModel*
method), 66
 get_names() (*evalml.pipelines.components.RFClassifierSelectFromModel*
method), 68
 get_names() (*evalml.pipelines.components.RFRegressorSelectFromModel*
method), 66
 get_pipelines() (*in module evalml.pipelines*), 107
 get_roc_data() (*evalml.AutoClassificationSearch.plot*
method), 61
 graph() (*evalml.pipelines.CatBoostClassificationPipeline*
method), 92
 graph() (*evalml.pipelines.CatBoostRegressionPipeline*
method), 102
 graph() (*evalml.pipelines.LinearRegressionPipeline*
method), 106

graph() (*evalml.pipelines.LogisticRegressionPipeline*
method), 96
 graph() (*evalml.pipelines.PipelineBase method*), 82,
 108
 graph() (*evalml.pipelines.RFClassificationPipeline*
method), 86
 graph() (*evalml.pipelines.RFRegressionPipeline*
method), 99
 graph() (*evalml.pipelines.XGBoostPipeline method*),
 89
 GridSearchTuner (*class in evalml.tuners*), 157

H

handle_problem_types() (*in module*
evalml.problem_types), 154
 hyperparameter_ranges
 (*evalml.pipelines.components.LinearRegressor*
attribute), 77
 hyperparameter_ranges
 (*evalml.pipelines.components.LogisticRegressionClassifier*
attribute), 72
 hyperparameter_ranges
 (*evalml.pipelines.components.OneHotEncoder*
attribute), 63
 hyperparameter_ranges
 (*evalml.pipelines.components.RandomForestClassifier*
attribute), 74
 hyperparameter_ranges
 (*evalml.pipelines.components.RandomForestRegressor*
attribute), 79
 hyperparameter_ranges
 (*evalml.pipelines.components.RFClassifierSelectFromModel*
attribute), 67
 hyperparameter_ranges
 (*evalml.pipelines.components.RFRegressorSelectFromModel*
attribute), 65
 hyperparameter_ranges
 (*evalml.pipelines.components.SimpleImputer*
attribute), 69
 hyperparameter_ranges
 (*evalml.pipelines.components.StandardScaler*
attribute), 71
 hyperparameter_ranges
 (*evalml.pipelines.components.XGBoostClassifier*
attribute), 76
 hyperparameters (*evalml.pipelines.CatBoostClassificationPipeline*
attribute), 90
 hyperparameters (*evalml.pipelines.CatBoostRegressionPipeline*
attribute), 100
 hyperparameters (*evalml.pipelines.LinearRegressionPipeline*
attribute), 104
 hyperparameters (*evalml.pipelines.LogisticRegressionPipeline*
attribute), 94

hyperparameters (*evalml.pipelines.RFClassificationPipeline* attribute), 84

hyperparameters (*evalml.pipelines.RFRegressionPipeline* attribute), 97

hyperparameters (*evalml.pipelines.XGBoostPipeline* attribute), 87

I

is_search_space_exhausted() (*evalml.tuners.GridSearchTuner* method), 158

is_search_space_exhausted() (*evalml.tuners.RandomSearchTuner* method), 159

is_search_space_exhausted() (*evalml.tuners.SKOptTuner* method), 156

is_search_space_exhausted() (*evalml.tuners.Tuner* method), 155

L

LeadScoring (class in *evalml.objectives*), 111

LinearRegressionPipeline (class in *evalml.pipelines*), 104

LinearRegressor (class in *evalml.pipelines.components*), 77

list_model_families() (in module *evalml.pipelines*), 107

load() (*evalml.pipelines.CatBoostClassificationPipeline* static method), 93

load() (*evalml.pipelines.CatBoostRegressionPipeline* static method), 103

load() (*evalml.pipelines.LinearRegressionPipeline* static method), 106

load() (*evalml.pipelines.LogisticRegressionPipeline* static method), 96

load() (*evalml.pipelines.PipelineBase* static method), 82

load() (*evalml.pipelines.RFClassificationPipeline* static method), 86

load() (*evalml.pipelines.RFRegressionPipeline* static method), 99

load() (*evalml.pipelines.XGBoostPipeline* static method), 89

load_breast_cancer() (in module *evalml.demos*), 57

load_data() (in module *evalml.preprocessing*), 57

load_diabetes() (in module *evalml.demos*), 57

load_fraud() (in module *evalml.demos*), 57

load_wine() (in module *evalml.demos*), 57

LogisticRegressionClassifier (class in *evalml.pipelines.components*), 72

LogisticRegressionPipeline (class in *evalml.pipelines*), 94

LogLoss (class in *evalml.objectives*), 137

MAE (class in *evalml.objectives*), 145

MaxError (class in *evalml.objectives*), 151

MCC (class in *evalml.objectives*), 139

MedianAE (class in *evalml.objectives*), 149

model_family (*evalml.pipelines.CatBoostClassificationPipeline* attribute), 90

model_family (*evalml.pipelines.CatBoostRegressionPipeline* attribute), 100

model_family (*evalml.pipelines.components.LinearRegressor* attribute), 77

model_family (*evalml.pipelines.components.LogisticRegressionClassifier* attribute), 72

model_family (*evalml.pipelines.components.OneHotEncoder* attribute), 63

model_family (*evalml.pipelines.components.RandomForestClassifier* attribute), 74

model_family (*evalml.pipelines.components.RandomForestRegressor* attribute), 79

model_family (*evalml.pipelines.components.RFClassifierSelectFromModel* attribute), 67

model_family (*evalml.pipelines.components.RFRegressorSelectFromModel* attribute), 65

model_family (*evalml.pipelines.components.SimpleImputer* attribute), 69

model_family (*evalml.pipelines.components.StandardScaler* attribute), 70

model_family (*evalml.pipelines.components.XGBoostClassifier* attribute), 75

model_family (*evalml.pipelines.LinearRegressionPipeline* attribute), 104

model_family (*evalml.pipelines.LogisticRegressionPipeline* attribute), 94

model_family (*evalml.pipelines.RFClassificationPipeline* attribute), 84

model_family (*evalml.pipelines.RFRegressionPipeline* attribute), 97

model_family (*evalml.pipelines.XGBoostPipeline* attribute), 87

ModelFamily (class in *evalml.model_family*), 62

MSE (class in *evalml.objectives*), 146

MSLE (class in *evalml.objectives*), 148

N

name (*evalml.pipelines.CatBoostClassificationPipeline* attribute), 90

name (*evalml.pipelines.CatBoostRegressionPipeline* attribute), 100

name (*evalml.pipelines.components.LinearRegressor* attribute), 77

name (*evalml.pipelines.components.LogisticRegressionClassifier* attribute), 72

name (*evalml.pipelines.components.OneHotEncoder* attribute), 63

name (*evalml.pipelines.components.RandomForestClassifier* attribute), 74
 name (*evalml.pipelines.components.RandomForestRegressor* attribute), 79
 name (*evalml.pipelines.components.RFClassifierSelectFromModel* attribute), 67
 name (*evalml.pipelines.components.RFRegressorSelectFromModel* attribute), 65
 name (*evalml.pipelines.components.SimpleImputer* attribute), 69
 name (*evalml.pipelines.components.StandardScaler* attribute), 70
 name (*evalml.pipelines.components.XGBoostClassifier* attribute), 75
 name (*evalml.pipelines.LinearRegressionPipeline* attribute), 104
 name (*evalml.pipelines.LogisticRegressionPipeline* attribute), 94
 name (*evalml.pipelines.RFClassificationPipeline* attribute), 84
 name (*evalml.pipelines.RFRegressionPipeline* attribute), 97
 name (*evalml.pipelines.XGBoostPipeline* attribute), 87

O

objective_function() (*evalml.objectives.FraudCost* method), 110
 objective_function() (*evalml.objectives.LeadScoring* method), 112
 OneHotEncoder (class in *evalml.pipelines.components*), 63

P

PipelineBase (class in *evalml.pipelines*), 80
 Precision (class in *evalml.objectives*), 119
 PrecisionMacro (class in *evalml.objectives*), 122
 PrecisionMicro (class in *evalml.objectives*), 121
 PrecisionWeighted (class in *evalml.objectives*), 124
 predict() (*evalml.objectives.AUC* method), 132
 predict() (*evalml.objectives.AUCMacro* method), 135
 predict() (*evalml.objectives.AUCMicro* method), 133
 predict() (*evalml.objectives.AUCWeighted* method), 136
 predict() (*evalml.objectives.ConfusionMatrix* method), 142
 predict() (*evalml.objectives.ExpVariance* method), 153
 predict() (*evalml.objectives.F1* method), 114
 predict() (*evalml.objectives.F1Macro* method), 117
 predict() (*evalml.objectives.F1Micro* method), 115
 predict() (*evalml.objectives.F1Weighted* method), 118
 predict() (*evalml.objectives.FraudCost* method), 110
 predict() (*evalml.objectives.LeadScoring* method), 112
 predict() (*evalml.objectives.LogLoss* method), 138
 predict() (*evalml.objectives.MAE* method), 146
 predict() (*evalml.objectives.MaxError* method), 152
 predict() (*evalml.objectives.MCC* method), 139
 predict() (*evalml.objectives.MedianAE* method), 150
 predict() (*evalml.objectives.MSE* method), 147
 predict() (*evalml.objectives.MSLE* method), 149
 predict() (*evalml.objectives.Precision* method), 120
 predict() (*evalml.objectives.PrecisionMacro* method), 123
 predict() (*evalml.objectives.PrecisionMicro* method), 121
 predict() (*evalml.objectives.PrecisionWeighted* method), 124
 predict() (*evalml.objectives.R2* method), 144
 predict() (*evalml.objectives.Recall* method), 126
 predict() (*evalml.objectives.RecallMacro* method), 129
 predict() (*evalml.objectives.RecallMicro* method), 127
 predict() (*evalml.objectives.RecallWeighted* method), 130
 predict() (*evalml.objectives.ROC* method), 141
 predict() (*evalml.pipelines.CatBoostClassificationPipeline* method), 93
 predict() (*evalml.pipelines.CatBoostRegressionPipeline* method), 103
 predict() (*evalml.pipelines.components.LinearRegressor* method), 78
 predict() (*evalml.pipelines.components.LogisticRegressionClassifier* method), 73
 predict() (*evalml.pipelines.components.RandomForestClassifier* method), 75
 predict() (*evalml.pipelines.components.RandomForestRegressor* method), 80
 predict() (*evalml.pipelines.components.XGBoostClassifier* method), 77
 predict() (*evalml.pipelines.LinearRegressionPipeline* method), 106
 predict() (*evalml.pipelines.LogisticRegressionPipeline* method), 96
 predict() (*evalml.pipelines.PipelineBase* method), 83
 predict() (*evalml.pipelines.RFClassificationPipeline* method), 86
 predict() (*evalml.pipelines.RFRegressionPipeline* method), 99
 predict() (*evalml.pipelines.XGBoostPipeline* method), 89
 predict_proba() (*evalml.pipelines.CatBoostClassificationPipeline*

method), 93

predict_proba() (evalml.pipelines.CatBoostRegressionPipeline method), 103

predict_proba() (evalml.pipelines.components.LinearRegressor method), 78

predict_proba() (evalml.pipelines.components.LogisticRegressionClassifier method), 74

predict_proba() (evalml.pipelines.components.RandomForestClassifier method), 75

predict_proba() (evalml.pipelines.components.RandomForestRegressor method), 80

predict_proba() (evalml.pipelines.components.XGBoostClassifier method), 77

predict_proba() (evalml.pipelines.LinearRegressionPipeline method), 106

predict_proba() (evalml.pipelines.LogisticRegressionPipeline method), 96

predict_proba() (evalml.pipelines.PipelineBase method), 83

predict_proba() (evalml.pipelines.RFClassificationPipeline method), 86

predict_proba() (evalml.pipelines.RFRegressionPipeline method), 100

predict_proba() (evalml.pipelines.XGBoostPipeline method), 90

ProblemTypes (class in evalml.problem_types), 154

propose() (evalml.tuners.GridSearchTuner method), 158

propose() (evalml.tuners.RandomSearchTuner method), 159

propose() (evalml.tuners.SKOptTuner method), 157

propose() (evalml.tuners.Tuner method), 155

R

R2 (class in evalml.objectives), 143

RandomForestClassifier (class in evalml.pipelines.components), 74

RandomForestRegressor (class in evalml.pipelines.components), 79

RandomSearchTuner (class in evalml.tuners), 158

Recall (class in evalml.objectives), 125

RecallMacro (class in evalml.objectives), 128

RecallMicro (class in evalml.objectives), 127

RecallWeighted (class in evalml.objectives), 130

RFClassificationPipeline (class in evalml.pipelines), 84

RFClassifierSelectFromModel (class in evalml.pipelines.components), 67

RFRegressionPipeline (class in evalml.pipelines), 97

RFRegressorSelectFromModel (class in evalml.pipelines.components), 65

ROC (class in evalml.objectives), 140

S

save() (evalml.pipelines.CatBoostClassificationPipeline method), 93

save() (evalml.pipelines.CatBoostRegressionPipeline method), 103

save() (evalml.pipelines.LinearRegressionPipeline method), 107

save() (evalml.pipelines.LogisticRegressionPipeline method), 97

save() (evalml.pipelines.PipelineBase method), 83

save() (evalml.pipelines.RFClassificationPipeline method), 87

save() (evalml.pipelines.RFRegressionPipeline method), 100

save() (evalml.pipelines.XGBoostPipeline method), 90

score() (evalml.objectives.AUC method), 132

score() (evalml.objectives.AUCMacro method), 135

score() (evalml.objectives.AUCMicro method), 134

score() (evalml.objectives.AUCWeighted method), 137

score() (evalml.objectives.ConfusionMatrix method), 143

score() (evalml.objectives.ExpVariance method), 153

score() (evalml.objectives.F1 method), 114

score() (evalml.objectives.F1Macro method), 117

score() (evalml.objectives.F1Micro method), 116

score() (evalml.objectives.F1Weighted method), 119

score() (evalml.objectives.FraudCost method), 110

score() (evalml.objectives.LeadScoring method), 112

score() (evalml.objectives.LogLoss method), 138

score() (evalml.objectives.MAE method), 146

score() (evalml.objectives.MaxError method), 152

score() (evalml.objectives.MCC method), 140

score() (evalml.objectives.MedianAE method), 150

score() (evalml.objectives.MSE method), 147

score() (evalml.objectives.MSLE method), 149

score() (evalml.objectives.Precision method), 120

score() (evalml.objectives.PrecisionMacro method), 123

score() (evalml.objectives.PrecisionMicro method), 122

score() (evalml.objectives.PrecisionWeighted method), 125

score() (evalml.objectives.R2 method), 144

score() (evalml.objectives.Recall method), 126

score() (evalml.objectives.RecallMacro method), 129

score() (evalml.objectives.RecallMicro method), 128

score() (evalml.objectives.RecallWeighted method), 131

score() (evalml.objectives.ROC method), 141

score() (evalml.pipelines.CatBoostClassificationPipeline method), 93

score() (evalml.pipelines.CatBoostRegressionPipeline method), 103

`score()` (*evalml.pipelines.LinearRegressionPipeline* method), 107
`score()` (*evalml.pipelines.LogisticRegressionPipeline* method), 97
`score()` (*evalml.pipelines.PipelineBase* method), 83
`score()` (*evalml.pipelines.RFClassificationPipeline* method), 87
`score()` (*evalml.pipelines.RFRegressionPipeline* method), 100
`score()` (*evalml.pipelines.XGBoostPipeline* method), 90
`SimpleImputer` (class in *evalml.pipelines.components*), 69
`SKOptTuner` (class in *evalml.tuners*), 156
`split_data()` (in module *evalml.preprocessing*), 58
`StandardScaler` (class in *evalml.pipelines.components*), 70
`summary` (*evalml.pipelines.CatBoostClassificationPipeline* attribute), 90
`summary` (*evalml.pipelines.CatBoostRegressionPipeline* attribute), 100
`summary` (*evalml.pipelines.LinearRegressionPipeline* attribute), 104
`summary` (*evalml.pipelines.LogisticRegressionPipeline* attribute), 94
`summary` (*evalml.pipelines.RFClassificationPipeline* attribute), 84
`summary` (*evalml.pipelines.RFRegressionPipeline* attribute), 97
`summary` (*evalml.pipelines.XGBoostPipeline* attribute), 87
`supported_problem_types` (*evalml.pipelines.CatBoostClassificationPipeline* attribute), 90
`supported_problem_types` (*evalml.pipelines.CatBoostRegressionPipeline* attribute), 100
`supported_problem_types` (*evalml.pipelines.components.LinearRegressor* attribute), 77
`supported_problem_types` (*evalml.pipelines.components.LogisticRegressionClassifier* attribute), 72
`supported_problem_types` (*evalml.pipelines.components.RandomForestClassifier* attribute), 74
`supported_problem_types` (*evalml.pipelines.components.RandomForestRegressor* attribute), 79
`supported_problem_types` (*evalml.pipelines.components.XGBoostClassifier* attribute), 75
`supported_problem_types` (*evalml.pipelines.LinearRegressionPipeline* attribute), 104
`supported_problem_types` (*evalml.pipelines.LogisticRegressionPipeline* attribute), 94
`supported_problem_types` (*evalml.pipelines.RFClassificationPipeline* attribute), 84
`supported_problem_types` (*evalml.pipelines.RFRegressionPipeline* attribute), 97
`supported_problem_types` (*evalml.pipelines.XGBoostPipeline* attribute), 87
`supports_problem_type()` (*evalml.objectives.AUC* class method), 132
`supports_problem_type()` (*evalml.objectives.AUCMacro* class method), 135
`supports_problem_type()` (*evalml.objectives.AUCMicro* class method), 134
`supports_problem_type()` (*evalml.objectives.AUCWeighted* class method), 137
`supports_problem_type()` (*evalml.objectives.ConfusionMatrix* class method), 143
`supports_problem_type()` (*evalml.objectives.ExpVariance* class method), 154
`supports_problem_type()` (*evalml.objectives.F1* class method), 114
`supports_problem_type()` (*evalml.objectives.F1Macro* class method), 117
`supports_problem_type()` (*evalml.objectives.F1Micro* class method), 116
`supports_problem_type()` (*evalml.objectives.F1Weighted* class method), 119
`supports_problem_type()` (*evalml.objectives.FraudCost* class method), 110
`supports_problem_type()` (*evalml.objectives.LeadScoring* class method), 112
`supports_problem_type()` (*evalml.objectives.LogLoss* class method), 138
`supports_problem_type()` (*evalml.objectives.MAE* class method), 146
`supports_problem_type()` (*evalml.objectives.MaxError* class method), 146

152
`supports_problem_type()`
 (`evalml.objectives.MCC` class method), 140
`supports_problem_type()`
 (`evalml.objectives.MedianAE` class method),
 151
`supports_problem_type()`
 (`evalml.objectives.MSE` class method), 148
`supports_problem_type()`
 (`evalml.objectives.MSLE` class method),
 149
`supports_problem_type()`
 (`evalml.objectives.Precision` class method),
 120
`supports_problem_type()`
 (`evalml.objectives.PrecisionMacro` class
 method), 123
`supports_problem_type()`
 (`evalml.objectives.PrecisionMicro` class
 method), 122
`supports_problem_type()`
 (`evalml.objectives.PrecisionWeighted` class
 method), 125
`supports_problem_type()` (`evalml.objectives.R2`
 class method), 145
`supports_problem_type()`
 (`evalml.objectives.Recall` class method),
 126
`supports_problem_type()`
 (`evalml.objectives.RecallMacro` class method),
 129
`supports_problem_type()`
 (`evalml.objectives.RecallMicro` class method),
 128
`supports_problem_type()`
 (`evalml.objectives.RecallWeighted` class
 method), 131
`supports_problem_type()`
 (`evalml.objectives.ROC` class method), 141

T

`transform()` (`evalml.pipelines.components.OneHotEncoder`
 method), 64
`transform()` (`evalml.pipelines.components.RFClassifierSelectFromModel`
 method), 68
`transform()` (`evalml.pipelines.components.RFRegressorSelectFromModel`
 method), 66
`transform()` (`evalml.pipelines.components.SimpleImputer`
 method), 70
`transform()` (`evalml.pipelines.components.StandardScaler`
 method), 72
 Tuner (class in `evalml.tuners`), 154

X

XGBoostClassifier (class in
`evalml.pipelines.components`), 75
 XGBoostPipeline (class in `evalml.pipelines`), 87