

Python's Scikit-learn

Nicolás García-Pedrajas

Computational Intelligence and Bioinformatics Research Group

October 23, 2019



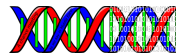
Table of contents

Scikit-learn

Decision trees in Python

Evaluating the models

Optimizing the models



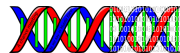
Scikit-learn

What is Scikit-Learn?

Extensions to SciPy (Scientific Python) are called SciKits. SciKit-Learn provides machine learning algorithms:

- ▶ Algorithms for supervised & unsupervised learning
- ▶ Built on SciPy and Numpy
- ▶ Standard Python API interface
- ▶ Sits on top of c libraries, LAPACK, LibSVM, and Cython
- ▶ Open Source: BSD License (part of Linux)

Probably the best general ML framework out there.



Scikit-learn

Where did it come from?

Started as a Google summer of code project in 2007 by David Cournapeau, then used as a thesis project by Matthieu Brucher.

In 2010, INRIA pushed the first public release, and sponsors the project, as do Google, Tinyclues, and the Python Software Foundation.



Scikit-learn

Primary features

- ▶ Generalized Linear Models
- ▶ SVMs, kNN, Bayes, Decision Trees, Ensembles
- ▶ Clustering and Density algorithms
- ▶ Cross Validation
- ▶ Grid Search
- ▶ Pipelining
- ▶ Model Evaluations
- ▶ Dataset Transformations
- ▶ Dataset Loading



Scikit-learn

API

Object-oriented interface centered around the concept of an Estimator:

An estimator is any object that learns from data; it may be a classification, regression or clustering algorithm or a transformer that extracts/filters useful features from raw data.

Scikit-Learn Tutorial



Scikit-learn

Estimator class

Class definition

```
1 class Estimator(object):
2
3     def fit(self, X, y=None):
4         """Fits estimator to data. """
5         # set state of ``self``
6         return self
7
8     def predict(self, X):
9         """Predict response of ``X``. """
10        # compute predictions ``pred``
11        return pred
```



Scikit-learn

Estimator

Estimators

- ▶ `fit(X,y)` sets the state of the estimator.
- ▶ `X` is usually a 2D numpy array of shape `(num_samples, num_features)`.
- ▶ `y` is a 1D array with shape `(n_samples,)`
- ▶ `predict(X)` returns the class or value
- ▶ `predict_proba()` returns a 2D array of shape `(n_samples, n_classes)`

Example:

Estimator (SVM)

```
1 from sklearn import svm
2
3 estimator = svm.SVC(gamma=0.001)
4 estimator.fit(X, y)
5 estimator.predict(x)
```



Scikit-learn

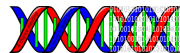
Load and transform data

Load data using appropriate methods.

Transform data:

Transformers

```
1 class Transformer(Estimator):
2     def transform(self, X):
3         """Transforms the input data. """
4         # transform ``X`` to ``X_prime``
5         return X_prime
6
7 from sklearn import preprocessing
8
9 Xt = preprocessing.normalize(X) # Normalizer
10 Xt = preprocessing.scale(X)
11
12 # StandardScaler: Imputation of missing values
13 imputer =Imputer(missing_values='Nan', strategy='mean')
14
15 Xt = imputer.fit_transform(X)
```



Scikit-learn

Classification models

Scikit-learn includes the following models:

- ▶ Generalized linear models.
- ▶ Linear an quadratic discriminant analysis
- ▶ Support vector machines.
- ▶ Nearest neighbors.
- ▶ Decision trees.
- ▶ Naive Bayes.
- ▶ Ensemble methods.
- ▶ Neural networks (deep learning with Keras)
- ▶ Multi-class methods.
- ▶ Multi-label methods.



Decision trees

First step: import required libraries.

Libraries

```
1 # Load libraries
2 import pandas as pd
3 from sklearn.tree import DecisionTreeClassifier # Import Decision Tree
  ↳ Classifier
4 from sklearn.model_selection import train_test_split # Import
  ↳ train_test_split function
5 from sklearn import metrics #Import scikit-learn metrics module for
  ↳ accuracy calculation
```



Decision trees

Second step: load data.

Load data

```
1 # load dataset
2 col_names = ['pregnant', 'glucose', 'bp', 'skin', 'insulin', 'bmi',
3             ↪ 'pedigree', 'age', 'label']
4 pima = pd.read_csv("pima-indians-diabetes.csv", header=None,
5             ↪ names=col_names)
6
7 # split dataset in features and target variable
8 feature_cols = ['pregnant', 'glucose', 'bp', 'skin', 'insulin', 'bmi',
9             ↪ 'pedigree', 'age']
10 X = pima[feature_cols] # Features
11 y = pima.label # Target variable
```



Decision trees

Data partitioning

We can split the data randomly (random partition, problem for reproduction):

Random partition

```
1 # Split dataset into training set and test set
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
  ↳ random_state=1) # 70% training and 30% test
```

Or we can have two separate files:

Separate files

```
1 # load dataset
2 col_names = ['pregnant', 'glucose', 'bp', 'skin', 'insulin', 'bmi',
  ↳ 'pedigree', 'age', 'label']
3 pima_train = pd.read_csv("pima-indians-diabetes.train.csv", header=None,
  ↳ names=col_names)
4 pima_test = pd.read_csv("pima-indians-diabetes.test.csv", header=None,
  ↳ names=col_names)
```



Decision tree

Building model

The procedure is common for all classification models.

Building decision tree

```
1 # Create Decision Tree classifier object
2 clf = DecisionTreeClassifier()
3
4 # Train Decision Tree Classifier
5 clf = clf.fit(X_train,y_train)
6
7 #Predict the response for test dataset
8 y_pred = clf.predict(X_test)
```



Visualizing decision trees

Decision trees can be visualize using **graphviz** and **pydotplus**:

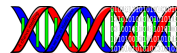
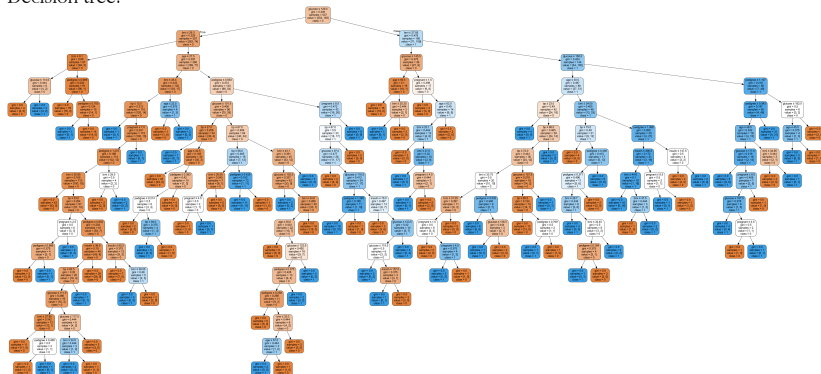
Building decision tree

```
1 from sklearn.tree import export_graphviz
2 from sklearn.externals.six import StringIO
3 from IPython.display import Image
4 import pydotplus
5
6 dot_data = StringIO()
7 export_graphviz(clf, out_file=dot_data,
8                 filled=True, rounded=True,
9                 special_characters=True, feature_names =
10                  ↪ feature_cols, class_names=['0', '1'])
11 graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
12 graph.write_png('diabetes.png')
13 Image(graph.create_png())
```



Visualizing decision trees

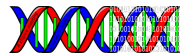
Decision tree:



Optimizing the decision tree

There are a few hyper-parameters:

- ▶ **criterion**: optional (default="gini") or Choose attribute selection measure: This parameter allows us to use the different-different attribute selection measure. Supported criteria are "gini" for the Gini index and "entropy" for the information gain.
- ▶ **splitter**: string, optional (default="best") or Split Strategy: This parameter allows us to choose the split strategy. Supported strategies are "best" to choose the best split and "random" to choose the best random split.
- ▶ **max_depth**: int or None, optional (default=None) or Maximum Depth of a Tree: The maximum depth of the tree. If None, then nodes are expanded until all the leaves contain less than `min_samples_split` samples. The higher value of maximum depth causes overfitting, and a lower value causes underfitting.



Decision trees

Probabilities output

The model can output categories and/or probabilities.

Models output

```
# Get predictions
preds = clf.predict(X_test)
np.savetxt("predictions", preds, fmt='%d')

# Get probabilities
probs = clf.predict_proba(X_test)
np.savetxt("probabilities", probs)
```



Metric module

The evaluation is common for all models.

Metric module implements many classification performance metrics:

<code>metrics.accuracy_score(y_true, y_pred[, ...])</code>	<code>metrics.hinge_loss(y_true, pred_decision[, ...])</code>
<code>metrics.auc(x, y[, reorder])</code>	<code>metrics.jaccard_score(y_true, y_pred[, ...])</code>
<code>metrics.average_precision_score(y_true, y_score)</code>	<code>metrics.log_loss(y_true, y_pred[, eps, ...])</code>
<code>metrics.balanced_accuracy_score(y_true, y_pred)</code>	<code>metrics.matthews_corrcoef(y_true, y_pred[, ...])</code>
<code>metrics.brier_score_loss(y_true, y_prob[, ...])</code>	<code>metrics.multilabel_confusion_matrix(y_true, y_pred[, ...])</code>
<code>metrics.classification_report(y_true, y_pred)</code>	<code>metrics.precision_recall_curve(y_true, y_score[, ...])</code>
<code>metrics.cohen_kappa_score(y1, y2[, labels, ...])</code>	<code>metrics.precision_recall_fscore_support(y_true, y_pred[, ...])</code>
<code>metrics.confusion_matrix(y_true, y_pred[, ...])</code>	<code>metrics.precision_score(y_true, y_pred[, ...])</code>
<code>metrics.f1_score(y_true, y_pred[, labels, ...])</code>	<code>metrics.recall_score(y_true, y_pred[, ...])</code>
<code>metrics.fbeta_score(y_true, y_pred, beta[, ...])</code>	<code>metrics.roc_auc_score(y_true, y_score[, ...])</code>
<code>metrics.hamming_loss(y_true, y_pred[, ...])</code>	<code>metrics.roc_curve(y_true, y_score[, ...])</code>
	<code>metrics.zero_one_loss(y_true, y_pred[, ...])</code>



Example of evaluation

Once trained and tested the model can be evaluated:

Evaluating the decision tree

```
1 # Model Accuracy, how often is the classifier correct?  
2 print("Accuracy:", metrics.accuracy_score(y_test, y_pred))  
3  
4 Accuracy: 0.6753246753246753
```

The confusion matrix can also be obtained from the predictions:

Confusion matrix for the decision tree

```
1 cm = confusion_matrix(y_test, y_pred)
```



Evaluating the models

Cross-validation

Python implements cross-validation as an option for evaluating the models

The following example demonstrates how to estimate the accuracy of a linear kernel support vector machine on the iris dataset by splitting the data, fitting a model and computing the score 5 consecutive times (with different splits each time):

k-Fold cross-validation

```
>>> from sklearn.model_selection import cross_val_score
>>> clf = svm.SVC(kernel='linear', C=1)
>>> scores = cross_val_score(clf, iris.data, iris.target, cv=5, scoring='accuracy')
>>> scores
array([0.96..., 1. ..., 0.96..., 0.96..., 1. ...])
```



Optimizing hyper-parameters

Hyper-parameters are critic for many classification models (e.g: Support vector machines)

A step of obtaining the best set of hyper-parameters is usually needed

Scikit provides a grid search using cross-validation



Optimizing hyper-parameters

Grid search

Grid search

```
1 #Grid Search
2 from sklearn.model_selection import GridSearchCV
3 clf = LogisticRegression()
4 grid_values = {'penalty': ['l1',
5   ↪  'l2'], 'C': [0.001, .009, 0.01, .09, 1, 5, 10, 25]}
6 grid_clf_acc = GridSearchCV(clf, param_grid = grid_values, scoring =
7   ↪  'recall')
8 grid_clf_acc.fit(X_train, y_train)
9
10 #Predict values based on new parameters
11 y_pred_acc = grid_clf_acc.predict(X_test)
12
13 # New Model Evaluation metrics
14 print('Accuracy Score : ' + str(accuracy_score(y_test, y_pred_acc)))
15 print('Precision Score : ' + str(precision_score(y_test, y_pred_acc)))
16 print('Recall Score : ' + str(recall_score(y_test, y_pred_acc)))
17 print('F1 Score : ' + str(f1_score(y_test, y_pred_acc)))
18
19 #Logistic Regression (Grid Search) Confusion matrix
20 confusion_matrix(y_test, y_pred_acc)
```

