

PYTHON FOR DATA SCIENCE CHEAT SHEET

Python Scikit-Learn

Introduction

Scikit-learn: "sklearn" is a machine learning library for the Python programming language. Simple and efficient tool for data mining, Data analysis and Machine Learning.

Importing Convention - import sklearn

Preprocessing

Data Loading

- **Using NumPy:**
>>>import numpy as np
>>>a=np.array([[1,2,3,4),(7,8,9,10)],dtype=int)
>>>data = np.loadtxt('file_name.csv', delimiter=',')
- **Using Pandas:**
>>>import pandas as pd
>>>df=pd.read_csv('file_name.csv',header=0)

Train-Test Data

```
>>>from sklearn.model_selection import train_test_split  
  
>>>X_train,X_test,y_train,y_test = train_test_split(X,y,random_state=0)
```

Data Preparation

- **Standardization**
>>>from sklearn.preprocessing import StandardScaler
>>>get_names = df.columns
>>>scaler = preprocessing.StandardScaler()
>>>scaled_df = scaler.fit_transform(df)
>>>scaled_df = pd.DataFrame(scaled_df, columns=get_names)

- **Normalization**
>>>from sklearn.preprocessing import Normalizer
>>>pd.read_csv("File_name.csv")
>>>x_array = np.array(df['Column1'])
>>>#Normalize Column1
>>>normalized_X = preprocessing.normalize([x_array])

Working On Model

Model Choosing

Supervised Learning Estimator:

- **Linear Regression:**
>>> from sklearn.linear_model import LinearRegression
>>> new_lr = LinearRegression(normalize=True)
- **Support Vector Machine:**
>>> from sklearn.svm import SVC
>>> new_svc = SVC(kernel='linear')

Naive Bayes:

```
>>> from sklearn.naive_bayes import GaussianNB  
>>> new_gnb = GaussianNB()  
  
• KNN:  
>>> from sklearn import neighbors  
>>> knn=neighbors.KNeighborsClassifier(n_neighbors=1)
```

Unsupervised Learning Estimator:

- **Principal Component Analysis (PCA):**
>>> from sklearn.decomposition import PCA
>>> new_pca= PCA(n_components=0.95)
- **K Means:**
>>> from sklearn.cluster import KMeans
>>> k_means = KMeans(n_clusters=5, random_state=0)

Train-Test Data

Supervised:

```
>>>new_lr.fit(X, y)  
>>> knn.fit(X_train, y_train)  
>>>new_svc.fit(X_train, y_train)
```

Unsupervised :

```
>>> k_means.fit(X_train)  
>>> pca_model_fit = new_pca.fit_transform(X_train)
```

Post-Processing

Prediction

Supervised:

```
>>> y_predict = new_svc.predict(np.random.random((3,5)))  
>>> y_predict = new_lr.predict(X_test)  
>>> y_predict = knn.predict_proba(X_test)
```

Unsupervised:

```
>>> y_pred = k_means.predict(X_test)
```

Model Tuning

Grid Search:

```
>>> from sklearn.grid_search import GridSearchCV  
>>> params = {"n_neighbors": np.arange(1,3), "metric": ["euclidean", "cityblock"]}  
>>> grid = GridSearchCV(estimator=knn, param_grid=params)  
>>> grid.fit(X_train, y_train)  
>>> print(grid.best_score_)  
>>> print(grid.best_estimator_.n_neighbors)
```

Randomized Parameter Optimization:

```
>>> from sklearn.grid_search import RandomizedSearchCV  
>>> params = {"n_neighbors": range(1,5), "weights": ["uniform", "distance"]}  
>>> rsearch = RandomizedSearchCV(estimator=knn, param_distributions=params, cv=4, n_iter=8, random_state=5)  
>>> rsearch.fit(X_train, y_train)  
>>> print(rsearch.best_score_)
```

Evaluate Performance

Classification:

1. Confusion Matrix:

```
>>> from sklearn.metrics import confusion_matrix  
>>> print(confusion_matrix(y_test, y_pred))
```

2. Accuracy Score:

```
>>> knn.score(X_test, y_test)  
>>> from sklearn.metrics import accuracy_score  
>>> accuracy_score(y_test, y_pred)
```

Regression:

1. Mean Absolute Error:

```
>>> from sklearn.metrics import mean_absolute_error  
  
>>> y_true = [3, -0.5, 2]  
>>> mean_absolute_error(y_true, y_predict)
```

2. Mean Squared Error:

```
>>> from sklearn.metrics import mean_squared_error  
>>> mean_squared_error(y_test, y_predict)
```

3. R² Score :

```
>>> from sklearn.metrics import r2_score  
>>> r2_score(y_true, y_predict)
```

Clustering:

1. Homogeneity:

```
>>> from sklearn.metrics import homogeneity_score  
>>> homogeneity_score(y_true, y_predict)
```

2. V-measure:

```
>>> from sklearn.metrics import v_measure_score  
>>> metrics.v_measure_score(y_true, y_predict)
```

Cross-validation:

```
>>> from sklearn.cross_validation import cross_val_score  
>>> print(cross_val_score(knn, X_train, y_train, cv=4))  
>>> print(cross_val_score(new_lr, X, y, cv=2))
```