

# Evaluate Your Model's Performance

## Classification Metrics

<b>Accuracy Score</b> >>> knn.score(X_test, y_test) >>> from sklearn.metrics import accuracy_score >>> accuracy_score(y_test, y_pred)	#Estimator score method
<b>Classification Report</b> >>> from sklearn.metrics import classification_report >>> print(classification_report(y_test, y_pred))	#Metric scoring functions
<b>Confusion Matrix</b> >>> from sklearn.metrics import confusion_matrix >>> print(confusion_matrix(y_test, y_pred))	#Precision, recall, f1-score and support

## Regression Metrics

**Mean Absolute Error**  
>>> from sklearn.metrics import mean\_absolute\_error  
>>> y\_true = [3, -0.5, 2]  
>>> mean\_absolute\_error(y\_true, y\_pred)

**Mean Squared Error**  
>>> from sklearn.metrics import mean\_squared\_error  
>>> mean\_squared\_error(y\_test, y\_pred)

**R<sup>2</sup> Score**  
>>> from sklearn.metrics import r2\_score  
>>> r2\_score(y\_true, y\_pred)

## Clustering Metrics

**Adjusted Rand Index**  
>>> from sklearn.metrics import adjusted\_rand\_score  
>>> adjusted\_rand\_score(y\_true, y\_pred)

**Homogeneity**  
>>> from sklearn.metrics import homogeneity\_score  
>>> homogeneity\_score(y\_true, y\_pred)

**V-measure**  
>>> from sklearn.metrics import v\_measure\_score  
>>> metrics.v\_measure\_score(y\_true, y\_pred)

## Cross-Validation

**Adjusted Rand Index**  
>>> from sklearn.cross\_validation import cross\_val\_score  
>>> print(cross\_val\_score(knn, X\_train, y\_train, cv=4))  
>>> print(cross\_val\_score(lr, X, y, cv=2))

## Tune Your Model

### 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 Your Model's Performance

## Classification Metrics

<b>Accuracy Score</b> >>> knn.score(X_test, y_test) >>> from sklearn.metrics import accuracy_score >>> accuracy_score(y_test, y_pred)	#Estimator score method
<b>Classification Report</b> >>> from sklearn.metrics import classification_report >>> print(classification_report(y_test, y_pred))	#Metric scoring functions
<b>Confusion Matrix</b> >>> from sklearn.metrics import confusion_matrix >>> print(confusion_matrix(y_test, y_pred))	#Precision, recall, f1-score and support

## Regression Metrics

**Mean Absolute Error**  
>>> from sklearn.metrics import mean\_absolute\_error  
>>> y\_true = [3, -0.5, 2]  
>>> mean\_absolute\_error(y\_true, y\_pred)

**Mean Squared Error**  
>>> from sklearn.metrics import mean\_squared\_error  
>>> mean\_squared\_error(y\_test, y\_pred)

**R<sup>2</sup> Score**  
>>> from sklearn.metrics import r2\_score  
>>> r2\_score(y\_true, y\_pred)

## Clustering Metrics

**Adjusted Rand Index**  
>>> from sklearn.metrics import adjusted\_rand\_score  
>>> adjusted\_rand\_score(y\_true, y\_pred)

**Homogeneity**  
>>> from sklearn.metrics import homogeneity\_score  
>>> homogeneity\_score(y\_true, y\_pred)

**V-measure**  
>>> from sklearn.metrics import v\_measure\_score  
>>> metrics.v\_measure\_score(y\_true, y\_pred)

## Cross-Validation

**Adjusted Rand Index**  
>>> from sklearn.cross\_validation import cross\_val\_score  
>>> print(cross\_val\_score(knn, X\_train, y\_train, cv=4))  
>>> print(cross\_val\_score(lr, X, y, cv=2))

## Tune Your Model

### 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_)
```

# Scikit-learn

## Standardization

```
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler().fit(X_train)
>>> standardized_X = scaler.transform(X_train)
>>> standardized_X_test = scaler.transform(X_test)
```

## Normalization

```
>>> from sklearn.preprocessing import Normalizer
>>> scaler = Normalizer().fit(X_train)
>>> normalized_X = scaler.transform(X_train)
>>> normalized_X_test = scaler.transform(X_test)
```

## Binarization

```
>>> from sklearn.preprocessing import Binarizer
>>> binarizer = Binarizer(threshold=0.0).fit(X)
>>> binary_X = binarizer.transform(X)
```

## Encoding Categorical Features

```
>>> from sklearn.preprocessing import LabelEncoder  
>>> enc = LabelEncoder()  
>>> y = enc.fit_transform(y)
```

## Imputing Missing Values

```
>>> from sklearn.preprocessing import Imputer  
>>> imp = Imputer(missing_values=0, strategy='mean', axis=0)  
>>> imp.fit_transform(X_train)
```

## Generating Polynomial Features

```
>>> from sklearn.preprocessing import PolynomialFeatures  
>>> poly = PolynomialFeatures(5)  
>>> poly.fit_transform(X)
```



# Scikit-learn

**Scikit-learn** is an open source Python library that implements a range of machine learning, preprocessing, cross-validation and visualization algorithms using a unified interface.



## A Basic Example

```
>>> from sklearn import neighbors, datasets, preprocessing
>>> from sklearn.cross_validation import train_test_split
>>> from sklearn.metrics import accuracy_score
>>> iris = datasets.load_iris()
>>> X, y = iris.data[:, :2], iris.target
>>> X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=33)
>>> scaler = preprocessing.StandardScaler().fit(X_train)
>>> X_train = scaler.transform(X_train)
>>> X_test = scaler.transform(X_test)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn.fit(X_train, y_train)
>>> y_pred = knn.predict(X_test)
>>> accuracy_score(y_test, y_pred)
```

## Loading The Data

Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrices. Other types that are convertible to numeric arrays, such as Pandas DataFrame, are also acceptable.

```
>>> import numpy as np
>>> X = np.random.random((10,5))
>>> y = np.array(['M', 'M', 'F', 'F', 'M', 'F', 'M', 'M', 'F', 'F'])
>>> X[X < 0.7] = 0
```

## Training And Test Data

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