

## کتابخانه Scikit-Learn

### معرفی کتابخانه Scikit-Learn

«سایکیت‌لرن» (Scikit-learn)، یک کتابخانه «متن‌باز» (Open Source) برای زبان پایتون است که طیفی از الگوریتم‌های «یادگیری ماشین» (Machine Learning)، «پیش‌پردازش داده‌ها» (Data Pre-Processing)، «اعتبارسنجی متقابل» (Cross Validation) و «بصری‌سازی» (Visualization) را با استفاده از یک رابط یکپارچه پیاده‌سازی می‌کند.

### یک مثال پایه‌ای

```
>>> from sklearn import neighbors, datasets, preprocessing
>>> from sklearn.model_selection 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)
```

### بارگذاری داده‌ها (Loading The Data)

داده‌ها باید به صورت عددی و آرایه‌های NumPy یا ماتریس خلوت SciPy ذخیره شده باشند. دیگر انواعی که قابل تبدیل به آرایه‌های عددی هستند، مانند دیتافریم Pandas نیز قابل پذیرش هستند.

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

### داده‌های آموزش و آزمون (Training And 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)
```

## پیش‌پردازش داده‌ها (Preprocessing The Data)

<p><b>رمزنگاری ویژگی‌های طبقه‌ای (Encoding Categorical Features)</b></p> <pre>&gt;&gt;&gt; from sklearn.preprocessing import LabelEncoder &gt;&gt;&gt; enc = LabelEncoder() &gt;&gt;&gt; y = enc.fit_transform(y)</pre>	<p><b>استانداردسازی (Standardization)</b></p> <pre>&gt;&gt;&gt; from sklearn.preprocessing import StandardScaler &gt;&gt;&gt; scaler = StandardScaler().fit(X_train) &gt;&gt;&gt; standardized_X = scaler.transform(X_train) &gt;&gt;&gt; standardized_X_test = scaler.transform(X_test)</pre>
<p><b>جایگذاری مقادیر ناموجود (Missing Values)</b></p> <pre>&gt;&gt;&gt; from sklearn.preprocessing import Imputer &gt;&gt;&gt; imp = Imputer(missing_values=0, strategy='mean', axis=0) &gt;&gt;&gt; imp.fit_transform(X_train)</pre>	<p><b>نرمال‌سازی (Normalization)</b></p> <pre>&gt;&gt;&gt; from sklearn.preprocessing import Normalizer &gt;&gt;&gt; scaler = Normalizer().fit(X_train) &gt;&gt;&gt; normalized_X = scaler.transform(X_train) &gt;&gt;&gt; normalized_X_test = scaler.transform(X_test)</pre>
<p><b>ساخت ویژگی‌های چند جمله‌ای (Polynomial Features)</b></p> <pre>&gt;&gt;&gt; from sklearn.preprocessing import PolynomialFeatures &gt;&gt;&gt; poly = PolynomialFeatures(5) &gt;&gt;&gt; poly.fit_transform(X)</pre>	<p><b>دودویی‌سازی (Binarization)</b></p> <pre>&gt;&gt;&gt; from sklearn.preprocessing import Binarizer &gt;&gt;&gt; binarizer = Binarizer(threshold=0.0).fit(X) &gt;&gt;&gt; binary_X = binarizer.transform(X)</pre>

## ساخت مدل

### برآوردگرهای یادگیری نظارت شده (Supervised Learning Estimators)

<pre>&gt;&gt;&gt; from sklearn.linear_model import LinearRegression &gt;&gt;&gt; lr = LinearRegression(normalize=True)</pre>	<p><b>رگرسیون خطی (Linear Regression)</b></p>
<pre>&gt;&gt;&gt; from sklearn.svm import SVC &gt;&gt;&gt; svc = SVC(kernel='linear')</pre>	<p><b>ماشین بردار پشتیبان (Support Vector Machines   SVM)</b></p>
<pre>&gt;&gt;&gt; from sklearn.naive_bayes import GaussianNB &gt;&gt;&gt; gnb = GaussianNB()</pre>	<p><b>نایو بیز (Naive Bayes)</b></p>
<pre>&gt;&gt;&gt; from sklearn import neighbors &gt;&gt;&gt; knn = neighbors.KNeighborsClassifier(n_neighbors=5)</pre>	<p><b>K نزدیک‌ترین همسایگی (K-Nearest Neighbors   KNN)</b></p>

### برآوردگرهای یادگیری نظارت نشده (Unsupervised Learning Estimators)

<p><b>K میانگین (K Means)</b></p> <pre>&gt;&gt;&gt; from sklearn.cluster import KMeans &gt;&gt;&gt; k_means = KMeans(n_clusters=3, random_state=0)</pre>	<p><b>تحلیل مولفه اساسی (Principal Component Analysis)</b></p> <pre>&gt;&gt;&gt; from sklearn.decomposition import PCA &gt;&gt;&gt; pca = PCA(n_components=0.95)</pre>
--	--

## برازش مدل (Model Fitting)

برازش مدل برای داده‌ها	<b>یادگیری نظارت شده (Supervised Learning)</b> <pre>&gt;&gt;&gt; lr.fit(X, y) &gt;&gt;&gt; knn.fit(X_train, y_train)</pre>
برازش مدل برای داده‌ها برازش برای داده‌ها و سپس تبدیل آن	<b>یادگیری نظارت نشده (Unsupervised Learning)</b> <pre>&gt;&gt;&gt; k_means.fit(X_train) &gt;&gt;&gt; pca_model = pca.fit_transform(X_train)</pre>

## پیش‌بینی (Prediction)

پیش‌بینی برجسب‌ها پیش‌بینی برجسب‌ها تخمین احتمال یک برجسب	<b>برآوردگرهای نظارت شده</b> <pre>&gt;&gt;&gt; y_pred = svc.predict(np.random.random((2,5))) &gt;&gt;&gt; y_pred = lr.predict(X_test) &gt;&gt;&gt; y_pred = knn.predict_proba(X_test)</pre>
پیش‌بینی برجسب‌ها در الگوریتم‌های خوشه‌بندی	<b>برآوردگرهای نظارت نشده</b> <pre>&gt;&gt;&gt; y_pred = k_means.predict(X_test)</pre>

## ارزیابی کارایی مدل (Model Performance Evaluation)

### سنجه‌های دسته‌بندی (Classification Metrics)

روش امتیاز برآوردگر توابع امتیازدهی متریک	<b>امتیاز صحت (Accuracy Score)</b> <pre>&gt;&gt;&gt; knn.score(X_test, y_test) &gt;&gt;&gt; from sklearn.metrics import accuracy_score &gt;&gt;&gt; accuracy_score(y_test, y_pred)</pre>
«صحت» (Precision)، «دقت» (Recall)، f <sub>1</sub> -Score و «پشتیبان» (Support)	<b>گزارش دسته‌بندی (Classification Report)</b> <pre>&gt;&gt;&gt; from sklearn.metrics import classification_report &gt;&gt;&gt; print(classification_report(y_test, y_pred))</pre>
	<b>ماتریس درهم‌ریختگی (Confusion Matrix)</b> <pre>&gt;&gt;&gt; from sklearn.metrics import confusion_matrix &gt;&gt;&gt; print(confusion_matrix(y_test, y_pred))</pre>

### سنجه‌های رگرسیون (Regression Metrics)

	<b>امتیاز R<sup>2</sup> (R<sup>2</sup> Score)</b> <pre>&gt;&gt;&gt; from sklearn.metrics import r2_score &gt;&gt;&gt; r2_score(y_true, y_pred)</pre>
--	---

**میانگین قدر مطلق خطا (Mean Absolute Error)**

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

**میانگین خطای مربعات (Mean Squared Error)**

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

**سنجه‌های خوشه‌بندی (Clustering Metrics)****اندیس تصادفی شده راند (Adjusted Rand Index)**

```
>>> from sklearn.metrics import adjusted_rand_score
>>> adjusted_rand_score(y_true, y_pred)
```

**همگن بودن (Homogeneity)**

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

**اندازه V (V-measure)**

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

**اعتبارسنجی متقابل (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(lr, X, y, cv=2))
```

**تنظیم مدل (Model Tuning)****جست‌وجوی شبکه‌ای (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)
```



## بهینه‌سازی پارامتر تصادفی شده (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_)
```

مجموعه آموزش‌های داده‌کاوی فرادرس (+ کلیک کنید)

برای مشاهده دیگر «تقلب‌نامه‌های» مجله فرادرس، به [این لینک](#) مراجعه فرمایید.

جهت آگاهی از آخرین تقلب‌نامه‌های منتشر شده، در [کانال تلگرام](#) مجله فرادرس عضو شوید.

تهیه و تنظیم: مجله فرادرس

