

Workflow of a basic example

```
1 from sklearn import datasets
2 from sklearn.model_selection import train_test_split
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.neighbors import KNeighborsClassifier
5 from sklearn.metrics import accuracy_score
6
7 iris = datasets.load_iris() #1.Preprocessing
8 X, y = iris.data[:, :3], iris.target
9 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state = 0)
10 X_train_std = StandardScaler().fit_transform(X_train)
11 X_test_std = StandardScaler().fit_transform(X_test)
12 knc = KNeighborsClassifier(n_neighbors=3) #2.Model creation
13 knc.fit(X_train_std, y_train) #3.Model fitting
14 y_pred = knc.predict(X_test_std) #4.Prediction
15 accuracy_score(y_test, y_pred) #5.Performance evaluation
```

1. Data Preprocessing

Loading data

```
1 # Import data
2 from sklearn import datasets
3 iris = datasets.load_iris()
4 X, y = iris.data[:, :4], iris.target
```

Training and test data

```
1 # Split training and test data
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state = 0)
```

Data preparation

```
1 # Standardization (z score = (x - μ) / σ )
2 scaler = StandardScaler().fit(X_train)
3 X_train_standardized = scaler.transform(X_train)
4 X_test_standardized = scaler.transform(X_test)
1 # Normalization ((X - X_min)/(X_max - X_min))
2 from sklearn.preprocessing import Normalizer
3 scaler = Normalizer().fit(X_train)
4 X_train_normalized = scaler.transform(X_train)
5 X_test_normalized = scaler.transform(X_test)
1 # Binarization (numerical features to boolean values)
2 from sklearn.preprocessing import Binarizer
3 binarizer = Binarizer(threshold = 3).fit(X_train)
4 X_binarized = binarizer.transform(X_test)
1 # Imputation of Missing Values
2 from sklearn.impute import SimpleImputer
3 imp = SimpleImputer(missing_values=np.nan, strategy='mean')
4 imp.fit_transform(X_train)
1 # Generating polynomial features
2 from sklearn.preprocessing import PolynomialFeatures
3 poly = PolynomialFeatures(2)
4 X_poly = poly.fit_transform(X_train)
1 # Custom transformers
2 from sklearn.preprocessing import FunctionTransformer
3 transformer = FunctionTransformer(np.log1p, validate=True)
4 X_transformed = transformer.transform(X_train)
1 # Binning
2 import numpy as np
3 from sklearn.preprocessing import KBinsDiscretizer
4 k_bins = KBinsDiscretizer(n_bins=3, encode='ordinal', strategy='kmeans').fit(X)
5 est = k_bins.transform(X_train)
1 # One-Hot Encoder
2 from sklearn.preprocessing import OneHotEncoder
3 enc = OneHotEncoder()
4 enc.fit(y.reshape(-1,1))
5 enc.transform(y.reshape(-1,1)).toarray()
```

```
1 # Label Encoder
2 from sklearn.preprocessing import LabelEncoder
3 enc = LabelEncoder()
4 y_encoded = enc.fit_transform(y)
```

2. Model Creation

Supervised learning - regression

```
1 # Supervised Learning Estimators - regression
2 # 1. Linear Regression
3 from sklearn.linear_model import LinearRegression
4 lr = LinearRegression()
5 lr.fit([[0, 0], [1, 1], [2, 2]], [0, 1, 2])
6 # 2. K Nearest Neighbor
7 from sklearn.neighbors import KNeighborsRegressor
8 knr = KNeighborsRegressor(n_neighbors=2)
```

Supervised learning - classification

```
1 # Supervised Learning Estimators - classification
2 # 1. Support Vector Machines (SVM)
3 from sklearn.svm import SVC
4 svc = SVC()
5 # 2. Naive Bayes
6 from sklearn.naive_bayes import GaussianNB
7 gnb = GaussianNB()
8 # 3. K Nearest Neighbor
9 from sklearn.neighbors import KNeighborsClassifier
10 knc = KNeighborsClassifier(n_neighbors=3)
# Unsupervised Learning Estimators
# K means
from sklearn.cluster import KMeans
k_means = KMeans(n_clusters = 3, random_state= 0)
# PCA
# Reduce number of attributes, while preserving as much info as possible
from sklearn.decomposition import PCA
pca = PCA(n_components=2)
pca.fit(X_train)
```

3. Model Fitting

```
1 # Model fitting
2 # Supervised learning
3 clf = svc.fit(X_train, y_train) # Fit the model to data
4 clf = knc.fit(X_train, y_train) # Fit the model to data
5 clf = gnb.fit(X_train, y_train) # Fit the model to data
6 # Unsupervised learning
7 reg = k_means.fit(X_train)
8 pca_model = pca.fit_transform(X_train) # Fit to the data, and transform it
```

4. Prediction

```
1 # Prediction
2 # Supervised learning
3 y_pred = svc.predict(X_test) # Predict labels
4 y_pred = lr.predict(X_test) # Predict labels
5 y_pred = knc.predict(X_test) # Predict labels
6 # Unsupervised learning
7 y_pred = k_means.predict(X_test) # Predict labels clustering
```

5. Performance Evaluation

Classification Metrics

```
1 # Classification Metrics
2 # 1. Accuracy score
3 from sklearn.metrics import accuracy_score
4 knc.score(X_test, y_test)
5 accuracy_score(y_test, y_pred)
6 # 2. Classification Report
7 from sklearn.metrics import classification_report
8 classification_report(y_test, y_pred)
9 # 3. Confusion matrix
10 from sklearn.metrics import confusion_matrix
11 confusion_matrix(y_test, y_pred)
```

Cross Validation

```
1 # Cross-validation
2 from sklearn.model_selection import cross_val_score
3 clf = svc.fit(X_train, y_train)
4 scores = cross_val_score(clf, X, y, cv=5, scoring='f1_macro')
```

Regression Metrics

```
1 # Regression Metrics
2 # 1. Mean Absolute Error
3 from sklearn.metrics import mean_absolute_error
4 mean_absolute_error(y_test, y_pred)
5 # 2. Mean Squared Error
6 from sklearn.metrics import mean_squared_error
7 mean_squared_error(y_test, y_pred)
8 # 3. R^2 score
9 from sklearn.metrics import r2_score
10 r2_score(y_test, y_pred)
```

Clustering Metrics

```
1 # Clustering Metrics
2 # Adjusted Rand Index
3 from sklearn.metrics import adjusted_rand_score
4 adjusted_rand_score(y, y_pred)
5 # Homogeneity
6 from sklearn.metrics import homogeneity_score
7 homogeneity_score(y_pred, y_pred)
8 # V-measure
9 from sklearn.metrics import v_measure_score
10 v_measure_score(y, y_pred)
```

Model Tuning

Grid Search

```
1 # Grid search
2 from sklearn.model_selection import GridSearchCV
3 parameters = {'kernel':('linear', 'rbf'), 'C':[1, 10]}
4 clf = GridSearchCV(svc, parameters)
5 clf.fit(X_train, y_train)
6 # To check the results
7 clf.cv_results_
```

Randomized Search

```
1 # Randomized Parameter Optimization
2 from sklearn.model_selection import RandomizedSearchCV
3 params = {'n_neighbors': [2,3,4], 'weights':['uniform','distance']}
4 clf = RandomizedSearchCV(estimator=knc,
5                           param_distributions=params,
6                           cv=4,
7                           n_iter=8,
8                           random_state=5)
9 clf.fit(X_train, y_train)
10 clf.cv_results_
```

Pipeline

```
1 from sklearn import datasets
2 from sklearn.model_selection import train_test_split
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.decomposition import PCA
5 from sklearn.tree import DecisionTreeClassifier
6 from sklearn.model_selection import GridSearchCV
7 from sklearn.pipeline import Pipeline
8 # Load data
9 iris = datasets.load_iris()
10 X, y = iris.data, iris.target
11 # Splitting data
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25)
13 # Making the Pipeline: PCA -> Scaling the data -> Classification
14 pipe = Pipeline([('pca', PCA()),
15                  ('scaler', StandardScaler()),
16                  ('classifier', DecisionTreeClassifier())])
17 # Fitting the model
18 parameters = {'pca__n_components': [2, 3, 4],
19               'classifier__max_depth': [5, 10, 20]}
20 grid = GridSearchCV(pipe, parameters).fit(X_train, y_train)
21 # Stores the optimum model in best_pipe
22 best_pipe = grid.best_estimator_
23 print(best_pipe)
24 print('Test set score: ' + str(best_pipe.score(X_test, y_test)))
```

