

# Scikit-learn Cheatsheet

## Essential operations for machine learning and data science

This cheatsheet provides a quick reference to fundamental scikit-learn operations, algorithms, and best practices, ideal for both beginners and experienced data scientists for efficient machine learning workflows.

Data Preprocessing	Model Selection	Model Training
Prepare and transform data	Choose and configure algorithms	Fit models to data
Model Evaluation	Feature Engineering	
Assess model performance	Extract and select features	

## Installation & Imports

### Installation: `pip install scikit-learn`

Install scikit-learn and common dependencies.

```
# Install scikit-learn
pip install scikit-learn
# Install with additional packages
pip install scikit-learn pandas numpy matplotlib
# Upgrade to latest version
pip install scikit-learn --upgrade
```

### Essential Imports

Standard imports for scikit-learn workflows.

```
# Core imports
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score,
accuracy_classification_report

# Common algorithms
from sklearn.linear_model import LinearRegression,
LogisticRegression
from sklearn.ensemble import RandomForestClassifier,
GradientBoostingRegressor
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
```

## Data Preprocessing

Transform and prepare data for machine learning algorithms.

01	02	03
<b>Train-Test Split:</b> `train_test_split()`	<b>Feature Scaling:</b> `StandardScaler()` / `MinMaxScaler()`	<b>Encoding:</b> `LabelEncoder()` / `OneHotEncoder()`
Divide data into training and testing sets.	Normalize features to similar scales.	Convert categorical variables to numerical format.

## Supervised Learning - Classification

### Logistic Regression: `LogisticRegression()`

Linear model for binary and multiclass classification.

```
# Basic logistic regression
from sklearn.linear_model import LogisticRegression
log_reg = LogisticRegression(random_state=42)
log_reg.fit(X_train, y_train)
y_pred = log_reg.predict(X_test)

# With regularization
log_reg_l2 = LogisticRegression(C=0.1, penalty='l2')
log_reg_l1 = LogisticRegression(C=0.1, penalty='l1',
solver='liblinear')
```

### Decision Tree: `DecisionTreeClassifier()`

Tree-based model for classification tasks.

```
# Decision tree classifier
from sklearn.tree import DecisionTreeClassifier
tree_clf = DecisionTreeClassifier(max_depth=5,
random_state=42)
tree_clf.fit(X_train, y_train)
y_pred = tree_clf.predict(X_test)

# Feature importance
importances = tree_clf.feature_importances_
# Visualize tree
from sklearn.tree import plot_tree
plot_tree(tree_clf, max_depth=3, filled=True)
```

## Supervised Learning - Regression

### Linear Regression: `LinearRegression()`

Basic linear model for continuous target variables.

```
# Simple linear regression
from sklearn.linear_model import LinearRegression
lin_reg = LinearRegression()
lin_reg.fit(X_train, y_train)
y_pred = lin_reg.predict(X_test)
```

```
# Get coefficients and intercept
coefficients = lin_reg.coef_
intercept = lin_reg.intercept_
print(f'R^2 score: {lin_reg.score(X_test, y_test)}")
```

### Ridge Regression: `Ridge()`

Linear regression with L2 regularization.

```
# Ridge regression (L2 regularization)
from sklearn.linear_model import Ridge
ridge_reg = Ridge(alpha=1.0)
ridge_reg.fit(X_train, y_train)
y_pred = ridge_reg.predict(X_test)
```

```
# Cross-validation for alpha selection
from sklearn.linear_model import RidgeCV
ridge_cv = RidgeCV(alphas=[0.1, 1.0, 10.0])
ridge_cv.fit(X_train, y_train)
```

## Model Evaluation

Assess model performance using various metrics.

### Classification Metrics

Evaluate classification model performance.

```
# Basic accuracy
from sklearn.metrics import accuracy_score,
precision_score, recall_score, f1_score
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred,
average='weighted')
recall = recall_score(y_test, y_pred, average='weighted')
f1 = f1_score(y_test, y_pred, average='weighted')
```

```
# Detailed classification report
from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred))

# Confusion matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
```

### ROC Curve & AUC

Plot ROC curve and calculate Area Under Curve.

```
# ROC curve for binary classification
from sklearn.metrics import roc_curve, auc
fpr, tpr, thresholds = roc_curve(y_test, y_proba[:, 1])
roc_auc = auc(fpr, tpr)
```

```
# Plot ROC curve
import matplotlib.pyplot as plt
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, label=f'ROC Curve (AUC = {roc_auc:.2f})')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.legend()
```

## Performance & Debugging

### Learning Curves: `learning\_curve()`

Diagnose overfitting and underfitting.

```
# Plot learning curves
from sklearn.model_selection import learning_curve
train_sizes, train_scores, val_scores = learning_curve(
model, X, y, cv=5, train_sizes=np.linspace(0.1, 1.0, 10))
plt.figure(figsize=(10, 6))
```

```
plt.plot(train_sizes, np.mean(val_scores, axis=1), 'o-',
label='Validation Score')
plt.plot(train_sizes, np.mean(train_scores, axis=1), 'o-',
label='Training Score')
plt.xlabel('Number of Estimators')
plt.ylabel('Score')
plt.legend()
```

### Validation Curves: `validation\_curve()`

Analyze the effect of hyperparameters.

```
# Validation curve for single hyperparameter
from sklearn.model_selection import validation_curve
param_grid = {'n_estimators': [10, 50, 100, 200, 500]}
train_scores, val_scores = validation_curve(
RandomForestClassifier(random_state=42), X, y,
param_name='n_estimators',
param_grid=param_grid, cv=5)
```

```
plt.figure(figsize=(10, 6))
plt.plot(param_grid, np.mean(val_scores, axis=1), 'o-',
label='Validation')
plt.plot(param_grid, np.mean(train_scores, axis=1), 'o-',
label='Training')
plt.xlabel('Number of Estimators')
plt.ylabel('Score')
plt.legend()
```

## Configuration & Best Practices

Optimize workflow and maintain reproducible results.

### Random State & Reproducibility

Ensure consistent results across runs.

```
# Set random state for
reproducibility
import numpy as np
np.random.seed(42)
# Set random_state in all
sklearn components
model = RandomForestClassifier(random_state=42)
train_test_split(X, y, random_state=42)
# For cross-validation
cv = StratifiedKFold(n_splits=5,
shuffle=True, random_state=42)
```

```
# For cross-validation
cv = StratifiedKFold(n_splits=5,
shuffle=True, random_state=42)
# For cross-validation
cv = StratifiedKFold(n_splits=5,
shuffle=True, random_state=42)
```

### Grid Search: `GridSearchCV()`

Exhaustive search over parameter grid.

```
# Grid search for hyperparameter tuning
from sklearn.model_selection import GridSearchCV
param_grid = {
'n_estimators': [100, 200, 300],
'max_depth': [3, 5, 7, None],
'min_samples_split': [2, 5, 10]
}
grid_search = GridSearchCV(model,
param_grid, cv=5, scoring='accuracy',
n_jobs=-1)
```

```
grid_search.fit(X_train, y_train)
best_params = grid_search.best_params_
best_params
```

### Random Search: `RandomizedSearchCV()`

Random sampling from parameter distributions.

```
# Randomized search (faster for large parameter spaces)
from sklearn.model_selection import RandomizedSearchCV
param_distributions = {
'n_estimators': randint(100, 500),
'max_depth': [3, 5, 7, None],
'min_samples_split': randint(2, 10)
}
random_search = RandomizedSearchCV(
RandomForestClassifier(random_state=42),
param_distributions, n_iter=50, cv=5)
```

```
random_search.fit(X_train, y_train)
best_params = random_search.best_params_
best_params
```

### Principal Component Analysis: `PCA()`

Dimensionality reduction technique.

```
# PCA for dimensionality reduction
from sklearn.decomposition import PCA
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X)
explained_variance = pca.explained_variance_ratio_
```

```
# Find optimal number of components
pca_full = PCA()
pca_full.fit(X)
cumsum =
np.cumsum(pca_full.explained_variance_ratio_)
# Find components for 95% variance
n_components = np.argmax(cumsum >= 0.95) + 1
```

### Ensemble Methods: `VotingClassifier()` / `BaggingClassifier()`

Combine multiple models for better performance.

```
# Voting classifier (ensemble of different algorithms)
from sklearn.ensemble import VotingClassifier
voting_clf = VotingClassifier(estimators=[
('lr', LogisticRegression(random_state=42)),
('rf', RandomForestClassifier(random_state=42)),
('svc', SVC(probability=True, random_state=42)])
voting_clf = VotingClassifier(estimators=[('lr', LogisticRegression(random_state=42)),
('rf', RandomForestClassifier(random_state=42)),
('svc', SVC(probability=True, random_state=42))], voting='soft')
```

```
voting_clf.fit(X_train, y_train)
y_pred = voting_clf.predict(X_test)

# Bagging classifier
from sklearn.ensemble import BaggingClassifier
bagging_clf = BaggingClassifier(DecisionTreeClassifier(),
n_estimators=100, random_state=42)
bagging_clf.fit(X_train, y_train)
```

### Gradient Boosting: `GradientBoostingClassifier()`

Sequential ensemble method with error correction.

```
# Gradient boosting classifier
from sklearn.ensemble import GradientBoostingClassifier
gb_clf = GradientBoostingClassifier(n_estimators=100,
learning_rate=0.1, random_state=42)
gb_clf.fit(X_train, y_train)
y_pred = gb_clf.predict(X_test)
```

```
# Feature importance
importances = gb_clf.feature_importances_
# Learning curve
from sklearn.model_selection import learning_curve
train_sizes, train_scores, val_scores = learning_curve(
gb_clf, X, y, cv=5)
```

### Model Selection & Hyperparameter Tuning

Find the best model and optimize hyperparameters.

### Grid Search: `GridSearchCV()`

Chain preprocessing and modeling pipeline.

```
# Create preprocessing and modeling pipeline
from sklearn.pipeline import Pipeline
pipeline = Pipeline([
('scaler', StandardScaler()),
('classifier', RandomForestClassifier(random_state=42))]
)
pipeline.fit(X_train, y_train)
y_pred = pipeline.predict(X_test)
```

### Random Search: `RandomizedSearchCV()`

Chain preprocessing and modeling pipeline.

```
# Create preprocessing and modeling pipeline
from sklearn.pipeline import Pipeline
pipeline = Pipeline([
('scaler', StandardScaler()),
('classifier', RandomForestClassifier(random_state=42))]
)
pipeline.fit(X_train, y_train)
y_pred = pipeline.predict(X_test)
```

### Feature Selection: `SelectKBest()` / `RFE()`

Select the most informative features.

```
# Uivariate feature selection
from sklearn.feature_selection import SelectKBest,
f_classif
selector = SelectKBest(score_func=f_classif, k=10)
X_selected = selector.fit_transform(X_train, y_train)
```

```
# Recursive Feature Elimination
from sklearn.feature_selection import RFE
rfe = RFE(RandomForestClassifier(random_state=42),
n_features_to_select=10)
X_rfe = rfe.fit_transform(X_train, y_train)
```

### Model Selection & Hyperparameter Tuning

Address class imbalance in datasets.

```
# Install imbalanced-learn: pip install imbalanced-learn
from imblearn.over_sampling import SMOTE
smote = SMOTE(random_state=42)
X_resampled, y_resampled = smote.fit_resample(X_train, y_train)
```

```
# Using class weights
rf_balanced =
RandomForestClassifier(class_weight='balanced',
random_state=42)
rf_balanced.fit(X_train, y_train)
```

```
# Manual class weights
from sklearn.utils.class_weight import
compute_class_weight
class_weights = compute_class_weight('balanced',
classes=np.unique(y_train), y=y_train)
weight_dict = dict(zip(np.unique(y_train), class_weights))
```

### Handling Imbalanced Data: `SMOTE` / `Class Weights`

Handle common issues and debug models.

```
# DBSCAN clustering
from sklearn.cluster import DBSCAN
dbscan = DBSCAN(eps=0.5, min_samples=5)
cluster_labels = dbscan.fit_predict(X)
n_clusters = len(set(cluster_labels)) - (1 if -1 in
cluster_labels else 0)
n_noise = list(cluster_labels).count(-1)
```

```
print(f'Number of clusters: {n_clusters}')
print(f'Number of noise points: {n_noise}')
```

### Hierarchical Clustering: `AgglomerativeClustering()`

Build hierarchy of clusters.

```
# Agglomerative clustering
from sklearn.cluster import AgglomerativeClustering
agg_clustering = AgglomerativeClustering(n_clusters=3,
linkage='ward')
cluster_labels = agg_clustering.fit_predict(X)
```

```
# Dendrogram visualization
from scipy.cluster.hierarchy import dendrogram, linkage
plt.figure(figsize=(12, 8))
dendrogram(linkage)
```

### Pipeline: `Pipeline()`

Chain preprocessing and modeling steps.

```
# Create preprocessing and modeling pipeline
from sklearn.pipeline import Pipeline
pipeline = Pipeline([
('scaler', StandardScaler()),
('classifier', RandomForestClassifier(random_state=42))]
)
pipeline.fit(X_train, y_train)
y_pred = pipeline.predict(X_test)
```

### Feature Selection: `SelectKBest()` / `RFE()`

Select the most informative features.

```
# Uivariate feature selection
from sklearn.feature_selection import SelectKBest,
f_classif
selector = SelectKBest(score_func=f_classif, k=10)
X_selected = selector.fit_transform(X_train, y_train)
```

```
# Recursive Feature Elimination
from sklearn.feature_selection import RFE
rfe = RFE(RandomForestClassifier(random_state=42),
n_features_to_select=10)
X_rfe = rfe.fit_transform(X_train, y_train)
```

### Feature Importance Visualization

Understand which features drive model predictions.

```
# Plot feature importance
importances = model.feature_importances_
indices = np.argsort(importances)[::-1]
plt.figure(figsize=(12, 8))
plt.title("Feature Importance")
plt.bar(range(X.shape[1]), importances[indices])
plt.xticks(range(X.shape[1]), X.columns[indices], rotation=90)
```

```
# SHAP values for model interpretability
# pip install shap
explainer = shap.TreeExplainer(model)
shap_values = explainer.shap_values[X_test]
shap.summary_plot(shap_values, X_test)
```

### Model Comparison

Compare multiple algorithms systematically.

```
# Compare multiple models
from sklearn.model_selection import cross_val_score
models = [
'LogisticRegression(random_state=42),
RandomForestClassifier(random_state=42),
'SVM': SVC(random_state=42),
'GradientBoosting':
GradientBoostingClassifier(random_state=42)]
results = {}
for name, model in models.items():
scores =
```