

## Cheat Sheet: Building Supervised Learning Models

## Common supervised learning models

Process Name	Brief Description	Code Syntax
One vs One classifier (using logistic regression)	<p><b>Process:</b> This method trains one classifier for each pair of classes.</p> <p><b>Key hyperparameters:</b></p> <ul style="list-style-type: none"><li>- `estimator`: Base classifier (e.g., logistic regression)</li></ul> <p><b>Pros:</b> Can work well for small datasets.</p> <p><b>Cons:</b> Computationally expensive for large datasets.</p> <p><b>Common applications:</b> Multiclass classification problems where the number of classes is relatively small.</p>	<div><div>1</div><div>from sklearn.multiclass import OneVsOneClassifier</div><div>2</div><div>from sklearn.linear_model import LogisticRegression</div><div>3</div><div>model = OneVsOneClassifier(LogisticRegression())</div><div></div><div></div></div>
One vs All classifier (using logistic regression)	<p><b>Process:</b> Trains one classifier per class, where each classifier distinguishes between one class and the rest.</p> <p><b>Key hyperparameters:</b></p> <ul style="list-style-type: none"><li>- `estimator`: Base classifier (e.g., Logistic Regression)</li><li>- `multi_class`: Strategy to handle multiclass classification ('ovr')</li></ul> <p><b>Pros:</b> Simpler and more scalable than One vs One.</p> <p><b>Cons:</b> Less accurate for highly imbalanced classes.</p> <p><b>Common applications:</b> Common in multiclass classification problems such as image classification.</p>	<div><div>1</div><div>from sklearn.multiclass import OneVsRestClassifier</div><div>2</div><div>from sklearn.linear_model import LogisticRegression</div><div>3</div><div>model = OneVsRestClassifier(LogisticRegression())</div><div></div><div></div></div> <div>or</div> <div><div>1</div><div>from sklearn.linear_model import LogisticRegression</div><div>2</div><div>model_ova = LogisticRegression(multi_class='ovr')</div><div></div><div></div></div>
Decision tree classifier	<p><b>Process:</b> A tree-based classifier that splits data into smaller subsets based on feature values.</p> <p><b>Key hyperparameters:</b></p> <ul style="list-style-type: none"><li>- `max_depth`: Maximum depth of the tree</li></ul> <p><b>Pros:</b> Easy to interpret and visualize.</p> <p><b>Cons:</b> Prone to overfitting if not pruned properly.</p> <p><b>Common applications:</b> Classification tasks, such as credit risk assessment.</p>	<div><div>1</div><div>from sklearn.tree import DecisionTreeClassifier</div><div>2</div><div>model = DecisionTreeClassifier(max_depth=5)</div><div></div><div></div></div>
Decision tree regressor	<p><b>Process:</b> Similar to the decision tree classifier, but used for regression tasks to predict continuous values.</p> <p><b>Key hyperparameters:</b></p> <ul style="list-style-type: none"><li>- `max_depth`: Maximum depth of the tree</li></ul> <p><b>Pros:</b> Easy to interpret, handles nonlinear data.</p> <p><b>Cons:</b> Can overfit and perform poorly on noisy data.</p> <p><b>Common applications:</b> Regression tasks, such as predicting housing prices.</p>	<div><div>1</div><div>from sklearn.tree import DecisionTreeRegressor</div><div>2</div><div>model = DecisionTreeRegressor(max_depth=5)</div><div></div><div></div></div>
Linear SVM classifier	<p><b>Process:</b> A linear classifier that finds the optimal hyperplane separating classes with a maximum margin.</p> <p><b>Key hyperparameters:</b></p> <ul style="list-style-type: none"><li>- `C`: Regularization parameter</li><li>- `kernel`: Type of kernel function ('linear', 'poly', 'rbf', etc.)</li><li>- `gamma`: Kernel coefficient (only for 'rbf', 'poly', etc.)</li></ul> <p><b>Pros:</b> Effective for high-dimensional spaces.</p> <p><b>Cons:</b> Not ideal for nonlinear problems without kernel tricks.</p> <p><b>Common applications:</b> Text classification and image recognition.</p>	<div><div>1</div><div>from sklearn.svm import SVC</div><div>2</div><div>model = SVC(kernel='linear', C=1.0)</div><div></div><div></div></div>
K-nearest neighbors classifier	<p><b>Process:</b> Classifies data based on the majority class of its nearest neighbors.</p> <p><b>Key hyperparameters:</b></p> <ul style="list-style-type: none"><li>- `n_neighbors`: Number of neighbors to use</li><li>- `weights`: Weight function used in prediction ('uniform' or 'distance')</li><li>- `algorithm`: Algorithm used to compute the nearest neighbors ('auto', 'ball_tree', 'kd_tree', 'brute')</li></ul> <p><b>Pros:</b> Simple and effective for small datasets.</p> <p><b>Cons:</b> Computationally expensive as the dataset grows.</p> <p><b>Common applications:</b> Recommendation systems, image recognition.</p>	<div><div>1</div><div>from sklearn.neighbors import KNeighborsClassifier</div><div>2</div><div>model = KNeighborsClassifier(n_neighbors=5, weights='uniform')</div><div></div><div></div></div>
Random Forest regressor	<p><b>Process:</b> An ensemble method using multiple decision trees to improve accuracy and reduce overfitting.</p> <p><b>Key hyperparameters:</b></p> <ul style="list-style-type: none"><li>- `n_estimators`: Number of trees in the forest</li><li>- `max_depth`: Maximum depth of each tree</li></ul> <p><b>Pros:</b> Less prone to overfitting than individual decision trees.</p> <p><b>Cons:</b> Model complexity increases with the number of trees.</p> <p><b>Common applications:</b> Regression tasks such as predicting sales or stock prices.</p>	<div><div>1</div><div>from sklearn.ensemble import RandomForestRegressor</div><div>2</div><div>model = RandomForestRegressor(n_estimators=100, max_depth=5)</div><div></div><div></div></div>
XGBoost regressor	<p><b>Process:</b> A gradient boosting method that builds trees sequentially to correct errors from previous trees.</p> <p><b>Key hyperparameters:</b></p> <ul style="list-style-type: none"><li>- `n_estimators`: Number of boosting rounds</li><li>- `learning_rate`: Step size to improve accuracy</li><li>- `max_depth`: Maximum depth of each tree</li></ul> <p><b>Pros:</b> High accuracy and works well with large datasets.</p>	<div><div>1</div><div>import xgboost as xgb</div><div>2</div><div>model = xgb.XGBRegressor(n_estimators=100, learning_rate=0.1, max_depth=5)</div><div></div><div></div></div>

	<b>Cons:</b> Computationally intensive, complex to tune. <b>Common applications:</b> Predictive modeling, especially in Kaggle competitions.	
--	---	--

Associated functions used

Method Name	Brief Description	Code Syntax
OneHotEncoder	Transforms categorical features into a one-hot encoded matrix.	<div><div>1</div><div>from sklearn.preprocessing import OneHotEncoder</div><div>2</div><div>encoder = OneHotEncoder(sparse=False)</div><div>3</div><div>encoded_data = encoder.fit_transform(categorical_data)</div><div><div></div><div></div></div></div>
accuracy_score	Computes the accuracy of a classifier by comparing predicted and true labels.	<div><div>1</div><div>from sklearn.metrics import accuracy_score</div><div>2</div><div>accuracy = accuracy_score(y_true, y_pred)</div><div><div></div><div></div></div></div>
LabelEncoder	Encodes labels (target variable) into numeric format.	<div><div>1</div><div>from sklearn.preprocessing import LabelEncoder</div><div>2</div><div>encoder = LabelEncoder()</div><div>3</div><div>encoded_labels = encoder.fit_transform(labels)</div><div><div></div><div></div></div></div>
plot_tree	Plots a decision tree model for visualization.	<div><div>1</div><div>from sklearn.tree import plot_tree</div><div>2</div><div>plot_tree(model, max_depth=3, filled=True)</div><div><div></div><div></div></div></div>
normalize	Scales each feature to have zero mean and unit variance (standardization).	<div><div>1</div><div>from sklearn.preprocessing import normalize</div><div>2</div><div>normalized_data = normalize(data, norm='l2')</div><div><div></div><div></div></div></div>
compute_sample_weight	Computes sample weights for imbalanced datasets.	<div><div>1</div><div>from sklearn.utils.class_weight import compute_sample_weight</div><div>2</div><div>weights = compute_sample_weight(class_weight='balanced', y=y)</div><div><div></div><div></div></div></div>
roc_auc_score	Computes the Area Under the Receiver Operating Characteristic Curve (AUC-ROC) for binary classification models.	<div><div>1</div><div>from sklearn.metrics import roc_auc_score</div><div>2</div><div>auc = roc_auc_score(y_true, y_score)</div><div><div></div><div></div></div></div>

Author

Jeff Grossman  
Abhishek Gagneja

