**Tutorials** **DSA** **Data Science** **Web Tech** **Courses**

Python for Machine Learning    Machine Learning with R    Machine Learning Algorithms    EDA    Math for Machine

# Comprehensive Guide to Classification Models in Scikit-Learn

Last Updated : 17 Jun, 2024

Scikit-Learn, a powerful and user-friendly machine learning library in Python, has become a staple for data scientists and machine learning practitioners. It offers a wide array of tools for data mining and data analysis, making it accessible and reusable in various contexts.

*This article delves into the classification models available in Scikit-Learn, providing a technical overview and practical insights into their applications.*

## Table of Content

- [What is Classification?](#)
- [Scikit-Learn Classification Models](#)
  - [1. Logistic Regression](#)
  - [2. K-Nearest Neighbors (KNN)](#)
  - [3. Support Vector Machines (SVM)](#)
  - [4. Decision Trees](#)
  - [5. Random Forest](#)
  - [6. Naive Bayes](#)
  - [7. Gradient Boosting](#)
- [Model Evaluation Metrics](#)
- [Practical Tips for Using Scikit-Learn](#)

## What is Classification?

[Classification](#) is a [supervised learning](#) technique where the goal is to predict the categorical class labels of new instances based on past observations. It involves training a model on a labeled dataset, where the target variable is categorical. Common applications include spam detection, image recognition, and medical diagnosis.

### Key Concepts in Classification

Before diving into the models, it's essential to understand some key concepts:

- **Features**: The input variables used to make predictions.
- **Labels**: The output variable that the model is trying to predict.
- **Training Data**: The dataset used to train the model.
- **Test Data**: The dataset used to evaluate the model's performance.

## Scikit-Learn Classification Models

Scikit-Learn provides a variety of classification algorithms, each with its strengths and weaknesses. Here, we explore some of the most commonly used models.

**1. Logistic Regression**

[Logistic Regression](#) is a linear model used for binary classification problems. It models the probability that a given input belongs to a particular class.

**Advantages:**

- **Simplicity and Interpretability:** Logistic regression is easy to implement and interpret. It provides a clear probabilistic framework for binary classification problems.
- **Efficiency:** It is computationally efficient and works well with large datasets.
- **No Need for Feature Scaling:** Logistic regression does not require feature scaling, making it simpler to use.

**Disadvantages:**

- **Linear Decision Boundary:** It assumes a linear relationship between the independent variables and the log-odds of the dependent variable, which may not always be true.
- **Overfitting with High-Dimensional Data:** Logistic regression can overfit when the number of observations is less than the number of features.
- **Not Suitable for Non-Linear Problems:** It cannot handle non-linear relationships unless features are transformed

**Implementing Logistic Regression:**

It models the relationship between the input features and the probability of a binary outcome using a logistic function (sigmoid function). The model finds the best-fitting parameters to maximize the likelihood of the observed data.

Python

```python
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.datasets import load_wine
from sklearn.metrics import accuracy_score,
classification_report

X, y = load_wine(return_X_y=True)
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.33, random_state=42)
model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)

y_pred = model.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred))
print(classification_report(y_test, y_pred))
```

**Output:**

```
Accuracy: 0.95
              precision    recall  f1-score   support

           0       0.97      1.00      0.99        19
           1       0.95      0.95      0.95        20
           2       0.93      0.90      0.92        20

    accuracy                           0.95        59
   macro avg       0.95      0.95      0.95        59
weighted avg       0.95      0.95      0.95        59
```

**When to Use:**

- Logistic regression is ideal for binary classification problems where the relationship between the features and the target variable is approximately linear.
- It is also useful as a baseline model due to its simplicity and interpretability.

## 2. K-Nearest Neighbors (KNN)

KNN is a simple, instance-based learning algorithm. It classifies a data point based on the majority class among its k-nearest neighbors.

Advantages:

- **No Training Period:** KNN is a lazy learner, meaning it does not require a training phase, which makes it fast to implement.
- **Flexibility:** It can handle multi-class classification problems and is easy to understand and implement.
- **Adaptability:** New data can be added seamlessly without retraining the model.

Disadvantages:

- **Computationally Intensive:** KNN can be slow for large datasets as it requires computing the distance between the new point and all existing points.
- **Sensitive to Noise and Outliers:** KNN is sensitive to noisy data and outliers, which can affect its performance.
- **Feature Scaling Required:** It requires feature scaling to ensure that all features contribute equally to the distance calculations

**Implementing K-Nearest Neighbors (KNN):**

The algorithm calculates the distance between the new data point and all existing points, then assigns the class of the majority of the k-nearest neighbors.

Python

```
from sklearn.neighbors import KNeighborsClassifier

knn = KNeighborsClassifier(n_neighbors=3)
knn.fit(X_train, y_train)
y_pred_knn = knn.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred_knn))
print(classification_report(y_test, y_pred_knn))
```

Output:

```
Accuracy: 0.95
              precision    recall  f1-score   support
```

|              |      |      |      |    |
|--------------|------|------|------|----|
| 0            | 1.00 | 1.00 | 1.00 | 19 |
| 1            | 0.95 | 0.95 | 0.95 | 20 |
| 2            | 0.90 | 0.90 | 0.90 | 20 |
|              |      |      |      |    |
| accuracy     |      |      | 0.95 | 59 |
| macro avg    | 0.95 | 0.95 | 0.95 | 59 |
| weighted avg | 0.95 | 0.95 | 0.95 | 59 |

**When to Use:**

- KNN is suitable for small to medium-sized datasets where the computational cost is manageable.
- It is also useful for problems where the decision boundary is not linear and can be complex.

### 3. Support Vector Machines (SVM)

SVM is a powerful classifier that works by finding the hyperplane that best separates the classes in the feature space.

**Advantages:**

- **Effective in High-Dimensional Spaces:** SVMs work well with high-dimensional data and are effective when the number of dimensions exceeds the number of samples.
- **Robust to Overfitting:** SVMs have good generalization performance and are less prone to overfitting, especially with the use of regularization.
- **Versatility:** They can handle both linear and non-linear classification problems using the kernel trick.

**Disadvantages:**

- **Computationally Expensive:** SVMs can be slow to train, especially with large datasets.
- **Choice of Kernel:** Selecting the appropriate kernel and tuning hyperparameters can be challenging.
- **Interpretability:** SVMs are less interpretable compared to simpler models like logistic regression

**Implementing Support Vector Machines (SVM):**

SVM finds the hyperplane that maximizes the margin between the closest points of different classes. For non-linear data, it uses kernel functions to transform the data into higher dimensions.

Python

```python
from sklearn.svm import SVC

svm = SVC(kernel='linear')
svm.fit(X_train, y_train)
y_pred_svm = svm.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred_svm))
print(classification_report(y_test, y_pred_svm))
```

Output:

```
Accuracy: 0.98
              precision    recall  f1-score   support

           0       1.00      1.00      1.00        19
           1       0.95      1.00      0.97        20
           2       1.00      0.95      0.97        20

    accuracy                           0.98        59
   macro avg       0.98      0.98      0.98        59
weighted avg       0.98      0.98      0.98        59
```

When to Use:

- SVMs are ideal for complex classification problems with high-dimensional data and when the decision boundary is non-linear.
- They are also useful when the dataset is small but has many features.

4. Decision Trees

Decision Trees are non-parametric models that split the data into subsets based on the value of input features. They are easy to interpret and visualize.

Advantages:

- **Interpretability:** Decision trees are easy to understand and interpret, making them useful for explaining model predictions to non-technical stakeholders.

- **Handling Different Data Types:** They can handle both numerical and categorical data without requiring feature scaling or encoding.
- **Non-Linear Relationships:** Decision trees can capture non-linear relationships between features and the target variable.

**Disadvantages:**

- **Prone to Overfitting:** Decision trees can easily overfit the training data, especially if they are deep and complex.
- **High Variance:** Small changes in the data can result in significantly different trees, making them unstable.
- **Bias Toward Dominant Classes:** They can be biased toward features with many levels or dominant classes in imbalanced datasets

**Implementing Decision Trees:**

The algorithm splits the data into subsets based on the value of input features, creating branches until it reaches a decision node (leaf).

Python

```python
from sklearn.tree import DecisionTreeClassifier

tree = DecisionTreeClassifier(max_depth=5)
tree.fit(X_train, y_train)
y_pred_tree = tree.predict(X_test)

print("Accuracy:", accuracy_score(y_test, y_pred_tree))
print(classification_report(y_test, y_pred_tree))
```

Output:

```
Accuracy: 0.93
              precision    recall  f1-score   support

           0       1.00      0.95      0.97        19
           1       0.91      0.95      0.93        20
           2       0.89      0.90      0.89        20

    accuracy                           0.93        59
   macro avg       0.93      0.93      0.93        59
weighted avg       0.93      0.93      0.93        59
```

**When to Use:**

- Decision trees are suitable for problems where interpretability is crucial, and the relationships between features and the target variable are non-linear.
- They are also useful for quick data exploration and feature selection.

### 5. Random Forest

Random Forest is an ensemble method that combines multiple decision trees to improve the model's accuracy and robustness.

**Advantages:**

- **High Accuracy:** Random forests generally provide high accuracy by averaging the predictions of multiple decision trees, reducing the variance.
- **Robustness to Noise:** They are resilient to noisy data and outliers, making them suitable for real-world datasets.
- **Feature Importance:** Random forests can provide insights into feature importance, aiding in feature selection and understanding the dataset.

**Disadvantages:**

- **Computationally Intensive:** Building multiple decision trees can be computationally expensive and require more resources.
- **Interpretability:** The ensemble nature makes it challenging to interpret the reasoning behind individual predictions compared to a single decision tree.
- **Bias in Imbalanced Datasets:** Random forests may be biased toward the majority class in imbalanced datasets, affecting the performance for minority classes

**Implementing Random Forest:**

The algorithm creates multiple decision trees using random subsets of the data and features, then averages their predictions.

Python

```
from sklearn.ensemble import RandomForestClassifier

forest = RandomForestClassifier(n_estimators=100, max_depth=5,
random_state=42)
forest.fit(X_train, y_train)
```

```
y_pred_forest = forest.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred_forest))
print(classification_report(y_test, y_pred_forest))
```

Output:

```
Accuracy: 0.98
              precision    recall  f1-score   support

           0       1.00      1.00      1.00        19
           1       0.95      1.00      0.97        20
           2       1.00      0.95      0.97        20

    accuracy                           0.98        59
   macro avg       0.98      0.98      0.98        59
weighted avg       0.98      0.98      0.98        59
```

**When to Use:**

- Random forests are ideal for complex classification problems with large and high-dimensional datasets.
- They are also useful when robustness to noise and feature importance insights are required.

**6. Naive Bayes**

Naive Bayes classifiers are based on Bayes' theorem and assume independence between features. They are particularly effective for text classification.

**Advantages:**

1. **Works Well with High-Dimensional Data:** Naive Bayes performs well with high-dimensional data, such as text classification tasks.
2. **Handles Multi-Class Problems:** It is effective for multi-class prediction problems.
3. **Less Training Data Required:** Naive Bayes can perform well with less training data if the independence assumption holds.

**Disadvantages:**

1. **Independence Assumption:** The algorithm assumes that all features are independent, which is rarely true in real-world scenarios. This can limit its

accuracy when features are correlated.

2. **Zero-Frequency Problem:** If a categorical variable in the test data set has a category that wasn't present in the training data, the model will assign it zero probability. This can be mitigated using smoothing techniques.

3. **Poor Estimation:** Naive Bayes is known to be a poor estimator, so the probability outputs from `predict_proba` should not be taken too seriously.

4. **Sensitivity to Irrelevant Features:** The algorithm can be sensitive to the presence of irrelevant features, which can affect its performance

**Implementing Naive Bayes:**

The algorithm applies Bayes' theorem with the assumption of independence between features. It calculates the probability of each class given the input features and selects the class with the highest probability.

Python

```python
from sklearn.naive_bayes import GaussianNB
nb = GaussianNB()
nb.fit(X_train, y_train)
y_pred_nb = nb.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred_nb))
print(classification_report(y_test, y_pred_nb))
```

Output:

```
Accuracy: 0.95
              precision    recall  f1-score   support

           0       1.00      0.95      0.97        19
           1       0.91      1.00      0.95        20
           2       0.95      0.90      0.92        20

    accuracy                           0.95        59
   macro avg       0.95      0.95      0.95        59
weighted avg       0.95      0.95      0.95        59
```

**When to Use:**

- Due to its speed, it is suitable for applications requiring real-time predictions.
- It is useful for problems involving multiple classes.

## 7. Gradient Boosting

Gradient Boosting is an ensemble technique that builds models sequentially, each correcting the errors of its predecessor.

Advantages:

1. **High Predictive Accuracy:** Gradient Boosting often provides predictive accuracy that is difficult to surpass.
2. **Flexibility:** It can optimize different loss functions and offers several hyperparameter tuning options, making it highly flexible.
3. **No Data Pre-Processing Required:** It often works well with both categorical and numerical values without requiring extensive data pre-processing.
4. **Handles Missing Data:** Gradient Boosting can handle missing data without the need for imputation.

Disadvantages:

1. **Overfitting:** Gradient Boosting models can overfit the training data, especially if not properly regularized. This can be mitigated using techniques like penalized learning, tree constraints, randomized sampling, and shrinkage.
2. **Computationally Expensive:** The algorithm is computationally intensive and often requires many trees (sometimes more than 1000), which can be time and memory exhaustive.
3. **Complex Hyperparameter Tuning:** The high flexibility results in many parameters that interact and influence the model's behavior, requiring extensive grid search during tuning.

### Implementing Gradient Boosting:

The algorithm combines weak learners (typically decision trees) in a sequential manner, where each new model focuses on the residual errors of the previous models.

Python

```
from sklearn.ensemble import GradientBoostingClassifier

gb = GradientBoostingClassifier(n_estimators=100,
learning_rate=1.0, max_depth=1, random_state=42)
```

```
gb.fit(X_train, y_train)
y_pred_gb = gb.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred_gb))
print(classification_report(y_test, y_pred_gb))
```

Output:

```
Accuracy: 0.97
              precision    recall  f1-score   support

           0       1.00      1.00      1.00        19
           1       0.95      1.00      0.97        20
           2       1.00      0.95      0.97        20

    accuracy                           0.97        59
   macro avg       0.98      0.98      0.98        59
weighted avg       0.97      0.97      0.97        59
```

**When to Use:**

- When the dataset has missing values, Gradient Boosting can be a good choice as it handles missing data natively.
- It is suitable for problems where the relationships between features and the target variable are complex and non-linear

*The outputs above are based on the `load_wine` dataset and a specific random state for splitting the data.*

## Model Evaluation Metrics

Evaluating the performance of classification models is crucial. Scikit-Learn provides several metrics to assess model performance:

- **Accuracy**: The ratio of correctly predicted instances to the total instances.
- **Precision**: The ratio of correctly predicted positive observations to the total predicted positives.
- **Recall (Sensitivity)**: The ratio of correctly predicted positive observations to all observations in the actual class.
- **F1-Score**: The weighted average of Precision and Recall.
- **Confusion Matrix**: A table used to describe the performance of a classification model.

Python

```python
from sklearn.metrics import confusion_matrix, precision_score,
recall_score, f1_score
cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:\n", cm)

precision = precision_score(y_test, y_pred, average='weighted')
print("Precision:", precision)

recall = recall_score(y_test, y_pred, average='weighted')
print("Recall:", recall)

f1 = f1_score(y_test, y_pred, average='weighted')
print("F1-Score:", f1)
```

Output:

```
Confusion Matrix:
 [[19  0  0]
  [ 0 19  1]
  [ 0  2 18]]

Precision: 0.95
Recall: 0.95
F1-Score: 0.95
```

## Practical Tips for Using Scikit-Learn

1. **Data Preprocessing**: Always preprocess your data by handling missing values, scaling features, and encoding categorical variables.
2. **Feature Selection**: Use techniques like Recursive Feature Elimination (RFE) to select the most relevant features.
3. **Cross-Validation**: Use cross-validation to get a more accurate estimate of your model's performance.
4. **Model Interpretation**: Use tools like SHAP and LIME to interpret your model's predictions.

## Conclusion

Scikit-Learn offers a comprehensive suite of tools for building and evaluating classification models. By understanding the strengths and weaknesses of each algorithm, you can choose the most appropriate model for your specific

problem. Remember to preprocess your data, tune hyperparameters, and evaluate your models using appropriate metrics to achieve the best results.

R   rs987…

Previous Article                                     Next Article

Interpreting Random Forest                 Is DevOps Useful for Machine Learning?
Classification Results

## Similar Reads

### Building a Custom Estimator for Scikit-learn: A Comprehensive Guide

Scikit-learn is a powerful machine learning library in Python that offers a wide range of tools for data analysis and modeling. One of its best features is the...

5 min read

### Save and Load Machine Learning Models in Python with scikit-learn

In this article, let's learn how to save and load your machine learning model in Python with scikit-learn in this tutorial. Once we create a machine learning...

4 min read

### Gaussian Mixture Models (GMM) Covariances in Scikit Learn

Gaussian Mixture Models (GMMs) are a type of probabilistic model used for clustering and density estimation. They assist in representing data as a...

6 min read

### How to Identify Overfitting Machine Learning Models in Scikit-Learn

Identifying overfitting in machine learning models is crucial to ensuring their performance generalizes well to unseen data. In this article, we'll explore ho...

7 min read

### Hidden Markov Models with Scikit-Learn

Hidden Markov Models (HMMs) are statistical models that represent systems that transition between a series of states over time. They are specially used i...

6 min read

View More Articles

**Article Tags :**          AI-ML-DS With Python          Data Science Blogathon 2024          Sklearn          AI-ML-DS

+3 More

**Practice Tags :**          Machine Learning

---

## GeeksforGeeks
Sanchhaya Education Private Limited

A-143, 9th Floor, Sovereign Corporate Tower, Sector-136, Noida, Uttar Pradesh - 201305

GET IT ON Google Play          Download on the App Store

### Company
About Us
Legal
In Media
Contact Us
Advertise with us
GFG Corporate Solution
Placement Training Program
GeeksforGeeks Community

### Languages
Python
Java
C++
PHP
GoLang
SQL
R Language
Android Tutorial
Tutorials Archive

### DSA
Data Structures
Algorithms
DSA for Beginners
Basic DSA Problems
DSA Roadmap
Top 100 DSA Interview Problems
DSA Roadmap by Sandeep Jain
All Cheat Sheets

### Data Science & ML
Data Science With Python
Data Science For Beginner
Machine Learning Tutorial
ML Maths
Data Visualisation Tutorial
Pandas Tutorial
NumPy Tutorial
NLP Tutorial
Deep Learning Tutorial

### Web Technologies
HTML
CSS
JavaScript

### Python Tutorial
Python Programming Examples
Python Projects
Python Tkinter

TypeScript

ReactJS

NextJS

Bootstrap

Web Design

Web Scraping

OpenCV Tutorial

Python Interview Question

Django

## Computer Science

Operating Systems

Computer Network

Database Management System

Software Engineering

Digital Logic Design

Engineering Maths

Software Development

Software Testing

## DevOps

Git

Linux

AWS

Docker

Kubernetes

Azure

GCP

DevOps Roadmap

## System Design

High Level Design

Low Level Design

UML Diagrams

Interview Guide

Design Patterns

OOAD

System Design Bootcamp

Interview Questions

## Inteview Preparation

Competitive Programming

Top DS or Algo for CP

Company-Wise Recruitment Process

Company-Wise Preparation

Aptitude Preparation

Puzzles

## School Subjects

Mathematics

Physics

Chemistry

Biology

Social Science

English Grammar

Commerce

World GK

## GeeksforGeeks Videos

DSA

Python

Java

C++

Web Development

Data Science

CS Subjects