

Stratified K-Fold Cross-Validation in Machine Learning

▼ Type

@datasciencebrain

Introduction

Stratified K-Fold Cross-Validation is an advanced technique used in machine learning for evaluating the performance of a model. It is a modification of the traditional K-Fold Cross-Validation, which ensures that each fold has a similar distribution of the target variable (class labels). This technique is especially useful when dealing with imbalanced datasets, as it ensures that each class is well-represented in both the training and validation sets.

In standard K-Fold Cross-Validation, the dataset is split into K subsets, or "folds," where one fold is used as the validation set and the remaining K-1 folds are used for training. This process is repeated K times, with each fold being used as the validation set once. Stratified K-Fold Cross-Validation takes this further by preserving the percentage of samples for each class in every fold, making it particularly effective when the target variable is imbalanced.

When to Use Stratified K-Fold Cross-Validation

1. **Imbalanced Datasets:** When the dataset contains a disproportionate number of examples in one class compared to others, Stratified K-Fold ensures that each fold has a proportional representation of each class.
2. **Classification Problems:** It is mainly used for classification problems, especially when there are multiple classes or class imbalances.
3. **Reliable Model Evaluation:** It provides a more reliable and robust estimate of model performance than simple train-test splits, especially in situations where

you want to ensure that each class is adequately represented in the training and validation sets.

How Stratified K-Fold Cross-Validation Works

1. **Data Splitting:** The dataset is divided into K equal-sized folds.
2. **Class Distribution:** In each fold, the proportion of samples for each class in the training and validation sets is similar to the entire dataset.
3. **Training and Validation:** The model is trained on K-1 folds and tested on the remaining fold, repeated for each of the K folds. This ensures that the model is evaluated on different data subsets and the performance is averaged.

Advantages of Stratified K-Fold Cross-Validation

1. **Preservation of Class Distribution:** The technique ensures that the distribution of the target classes is preserved in each fold, leading to more reliable estimates of model performance.
2. **Better Performance on Imbalanced Data:** For imbalanced datasets, Stratified K-Fold Cross-Validation avoids situations where some folds may not contain samples from the minority class.
3. **Efficient Use of Data:** Each data point is used for both training and validation, improving the efficiency of the model evaluation process.

Disadvantages of Stratified K-Fold Cross-Validation

1. **Computational Complexity:** Stratified K-Fold Cross-Validation can be computationally expensive, especially for large datasets, as the model is trained K times.

2. **Not Suitable for Regression:** It is primarily used for classification tasks and may not be ideal for regression tasks, where K-Fold is more appropriate.

Implementing Stratified K-Fold Cross-Validation in Python

The `scikit-learn` library provides an easy-to-use implementation of Stratified K-Fold Cross-Validation using the `StratifiedKFold` class. Below is a detailed explanation along with code examples.

1. Importing Required Libraries

```
import numpy as np
import pandas as pd
from sklearn.model_selection import StratifiedKFold
from sklearn.datasets import make_classification
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
```

2. Creating a Sample Dataset

For the sake of this example, we will create a synthetic classification dataset using `make_classification`.

```
# Create a synthetic classification dataset
X, y = make_classification(n_samples=1000, n_features=20, n_informative=2,
                          n_classes=3, weights=[0.2, 0.3, 0.5], random_state=42)
```

3. Setting Up Stratified K-Fold Cross-Validation

The `StratifiedKFold` class can be used to split the dataset into K folds. We specify the number of splits (K) and whether to shuffle the data before splitting.

```
# Initialize Stratified K-Fold with 5 splits
skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
```

4. Running the Cross-Validation

For each fold, we will split the data into training and validation sets, train a Logistic Regression model, and calculate the accuracy.

```
# Initialize a list to store accuracy scores for each fold
accuracies = []

# Loop through the Stratified K-Folds
for train_index, val_index in skf.split(X, y):
    # Split data into training and validation sets
    X_train, X_val = X[train_index], X[val_index]
    y_train, y_val = y[train_index], y[val_index]

    # Initialize the model (Logistic Regression)
    model = LogisticRegression(max_iter=1000)

    # Train the model
    model.fit(X_train, y_train)

    # Predict on the validation set
    y_pred = model.predict(X_val)

    # Calculate accuracy for the current fold
    accuracy = accuracy_score(y_val, y_pred)
    accuracies.append(accuracy)

# Calculate the average accuracy
average_accuracy = np.mean(accuracies)
print(f"Average Accuracy across all folds: {average_accuracy:.4f}")
```

5. Output

The output will display the accuracy of the model on each fold and the average accuracy across all folds.

Average Accuracy across all folds: 0.92

Code Explanation

1. **Dataset Creation:** We use `make_classification` to generate a synthetic classification dataset with 1000 samples, 20 features, and 3 classes.
2. **Stratified K-Fold Setup:** We initialize `StratifiedKFold` with 5 splits, setting `shuffle=True` to ensure the data is shuffled before splitting.
3. **Model Training and Evaluation:** In each fold, we split the dataset into training and validation sets. We use Logistic Regression for model training and evaluate the performance using accuracy. This process is repeated for each fold.
4. **Average Performance:** Finally, we calculate and print the average accuracy across all folds, which provides a reliable estimate of model performance.

Key Points to Remember

- **Stratification:** The main benefit of Stratified K-Fold Cross-Validation is that it ensures each fold has a similar distribution of the target variable.
- **Shuffling:** Shuffling the data before splitting can improve generalization, but it is optional depending on the nature of your data.
- **K-Fold Variations:** Stratified K-Fold Cross-Validation is specifically used for classification tasks. For regression problems, you should use K-Fold Cross-Validation without stratification.