# K-Fold Cross Validation In Machine Learning

| ⌄ Type | Web development |
| --- | --- |

**K-Fold Cross Validation** is a technique used in machine learning to evaluate the performance of a model. It helps to reduce the variability in model evaluation by using different subsets of the data. This method is particularly useful when working with limited data.

In K-fold cross-validation, the dataset is divided into $K$ smaller subsets (folds). The model is trained and evaluated $K$ times, each time using a different fold as the test set and the remaining folds as the training set.

## How it Works:

1. **Split the Data**: Divide the dataset into $K$ equally sized folds (subsets).

2. **Train and Test**: For each fold:

   - Treat the current fold as the test set.

   - Combine the remaining $K-1$ folds as the training set.

   - Train the model on the training set and evaluate it on the test set.

3. **Average the Results**: Once all $K$ iterations are complete, the model's overall performance is averaged (e.g., accuracy, precision, recall) to provide a more reliable estimate of its effectiveness.

## Key Points:

- **K**: The number of subsets (folds) in which to split the data. Common choices are 5 or 10, but you can choose any number.

- **No Overlap**: Every instance of the data is used exactly once as a part of the test set and multiple times as a part of the training set.

- **Performance Evaluation**: The performance scores from each fold are averaged to produce a final score.

# Advantages of K-Fold Cross Validation:

1. **Reduces Bias**: By testing the model on different subsets of the data, K-Fold reduces the variance of the performance estimate.

2. **Utilizes All Data**: All data points are used for both training and testing, providing a more robust performance measure.

3. **Efficient**: It's computationally more efficient compared to a simple train-test split, especially with small datasets.

# Disadvantages of K-Fold Cross Validation:

1. **Computationally Expensive**: For large datasets or complex models, K-Fold cross-validation can be computationally expensive because the model needs to be trained `K` times.

2. **Data Leakage Risk**: If data is not properly split, there can be the risk of data leakage between training and test sets, leading to overly optimistic results.

# Implementation in Python (Code Example):

To demonstrate K-Fold cross-validation, we will use the `KFold` class from `sklearn.model_selection` along with a simple classifier (e.g., `LogisticRegression`), but the same principles apply to other models as well.

## Step-by-Step Code Example:

```
# Import necessary libraries
import numpy as np
from sklearn.model_selection import KFold
from sklearn.linear_model import LogisticRegression
from sklearn.datasets import load_iris
```

```python
from sklearn.metrics import accuracy_score

# Load dataset (Iris dataset as an example)
data = load_iris()
X = data.data  # Features
y = data.target  # Labels

# Number of folds
k = 5

# Initialize KFold
kf = KFold(n_splits=k, shuffle=True, random_state=42)

# List to store the accuracy for each fold
accuracies = []

# Loop over each fold
for train_index, test_index in kf.split(X):
    # Split data into train and test sets for this fold
    X_train, X_test = X[train_index], X[test_index]
    y_train, y_test = y[train_index], y[test_index]

    # Initialize the model
    model = LogisticRegression(max_iter=200)

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

    # Predict on the test set
    y_pred = model.predict(X_test)

    # Calculate accuracy
    accuracy = accuracy_score(y_test, y_pred)
    accuracies.append(accuracy)

    print(f'Fold Accuracy: {accuracy:.4f}')
```

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

## Explanation of the Code:

1. **Data**: We use the Iris dataset, which is available in `sklearn.datasets`.

2. **KFold**: We create an instance of `KFold` with `n_splits=5`, meaning the data will be split into 5 folds. We set `shuffle=True` to shuffle the data before splitting to ensure random distribution of data.

3. **Model**: We use a logistic regression model for classification, but this can be replaced by any machine learning model (e.g., decision trees, SVM, etc.).

4. **Training and Testing**: For each fold, we train the model on the training data and evaluate it on the test data.

5. **Accuracy**: We calculate the accuracy for each fold and store it in the list `accuracies`.

6. **Average Performance**: After all folds are completed, the average accuracy is computed to get the final model performance.

## Output Example:

```
Fold Accuracy: 0.9667
Fold Accuracy: 0.9667
Fold Accuracy: 0.9667
Fold Accuracy: 0.9333
Fold Accuracy: 0.9667

Average Accuracy: 0.9667
```

Here, the model's performance is averaged over 5 folds to give an overall measure of how well the model is likely to perform on unseen data.

# Tuning the Number of Folds (K):

- **Small** `K` : With small values of `K` (e.g., 2 or 3), the variance in the performance metric may increase. Also, some portions of the data may not be used for training in some folds.

- **Large** `K` : With larger `K` (e.g., 10), the model is evaluated more times, leading to a more reliable estimate, but it becomes more computationally expensive.

---

# Using Cross-Validation with Scikit-learn's `cross_val_score` :

`scikit-learn` provides an easier way to perform cross-validation using the `cross_val_score` function, which automatically splits the data into folds and evaluates the model.

```python
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression

# Initialize the model
model = LogisticRegression(max_iter=200)

# Perform 5-fold cross-validation and get accuracy scores
scores = cross_val_score(model, X, y, cv=5, scoring='accuracy')

# Print the results
print(f'Accuracy per fold: {scores}')
print(f'Average accuracy: {scores.mean():.4f}')
```

## Explanation:

- `cross_val_score` : This function performs K-fold cross-validation in one line. The `cv` parameter specifies the number of folds, and the `scoring` parameter is set to `'accuracy'` for classification tasks.

- This method simplifies the process and is highly recommended for quick experimentation.

## Summary of Key Concepts:

- **K-Fold Cross Validation** splits the data into `K` folds, trains the model on `K-1` folds, and tests it on the remaining fold. This process is repeated for each fold.

- **Shuffling** the data before splitting ensures randomness in data selection and prevents biases.

- The performance across all folds is averaged to provide a more reliable estimate of the model's performance.

- **Scikit-learn's** `cross_val_score` simplifies the implementation of cross-validation.

## Key Takeaways:

- K-Fold Cross Validation is a robust method for assessing model performance.

- It is particularly useful in cases where the dataset is small or when you want a more reliable evaluation of a model.

- The number of folds `K` should be chosen based on the dataset size, model complexity, and computational resources.