# Title: Day 12 - K-fold cross validation

## Introduction:

K-fold cross-validation is a method in machine learning for testing a model's performance. It involves dividing the data into K subsets, training the model on K-1 subsets, and testing it on the remaining one. This process is repeated K times, providing a reliable assessment of the model's ability to generalize to new data. It's a crucial technique for model evaluation and selection.

## Tasks and Operations on Implementing K fold:

1. Importing Libraries:
   We imported necessary libraries, including Pandas for data manipulation and scikit-learn for machine learning tools.
2. Loading the Dataset:
   We loaded the dataset from a CSV file called 'FAKEBILL_third_day.csv' using Pandas. This dataset is assumed to contain information about bills, including features and a target variable.
3. Extracting Features and Target Variable:
   We separated the dataset into features (X) and the target variable (y). The features are a set of characteristics used to make predictions, and the target variable is what we want to predict.
4. Creating a K-fold Cross-Validation Iterator:
   We created a K-fold cross-validation iterator (kf) with five splits (K=5), which will be used to split the data into training and testing sets in a controlled way. We also enabled shuffling for randomization and set a random seed (random_state=42) for reproducibility.

5. Choosing a Machine Learning Model:

   We selected the Random Forest Classifier (clf) as the machine learning model. This model is used for classification tasks and is known for its effectiveness in various applications.

6. Performing K-fold Cross-Validation:

   We used the cross_val_score function to perform K-fold cross-validation. The model (clf), the feature data (X), and the target data (y) are provided as input. The cv parameter is set to the cross-validation iterator (kf), and we specified that we want to measure accuracy (scoring='accuracy').

7. Calculating and Printing Accuracy:

   We calculated the accuracy for each fold and printed it. The accuracy scores for each fold represent how well the model performs on each subset of the data.

   We also calculated and printed the mean accuracy across all the folds. The mean accuracy is a measure of the model's overall performance in a robust manner, considering different subsets of the data.


## Benefit of using K-fold cross validation:

K-fold cross-validation reduces overfitting risk, facilitates effective model selection, and ensures reproducible and reliable performance metrics by evaluating a model across different data subsets.


## Conclusion:

In conclusion, K-fold cross-validation is a valuable technique in machine learning, offering important benefits that include mitigating overfitting, aiding in model selection, and providing robust and reproducible performance metrics. It is a crucial tool for building reliable and generalizable models in diverse applications.