# Intro to Data Science
## Assignment # 04

Name:
    Muhammad Nadeem

Roll:
    201980050

Q: K-fold cross-validation with it's types
Solution:
K-fold cross-validation is a technique used to assess the performance of a machine learning model and to estimate its generalization ability on new, unseen data. It involves splitting the dataset into multiple subsets (folds), training and evaluating the model multiple times while rotating which fold is used for testing. This provides a more robust estimate of the model's performance compared to a single train-test split.

Here's a **step-by-step** explanation of K-fold cross-validation:

## Data Preparation:

Begin by preparing your dataset. This includes cleaning, preprocessing, and any feature engineering that might be necessary.

## Choose K and Shuffle:

Choose the number of folds, denoted as K. Common values are 5 or 10. The higher the value of K, the more accurate the estimate of model performance, but it also requires more computation. Shuffle the dataset randomly to ensure that the data points are not ordered in a specific way that might bias the results.

## Split Data into Folds:

Divide the dataset into K equal-sized (or nearly equal-sized) subsets, also known as folds. Each fold will be used as a validation set once, and the rest of the folds will be used for training.

**For each fold, do:**
**a. Model Training:**

Train your machine learning model on K-1 folds. These folds are used as the training set. This means that you'll train the model K times, each time using a different fold as the validation set.

**b. Model Validation:**

Use the remaining fold as the validation set to evaluate the model's performance. Calculate the relevant performance metrics (e.g., accuracy, F1-score, etc.) on this fold.

**Performance Metrics Aggregation:**

After performing K iterations (one for each fold), you will have K sets of performance metrics. Calculate the average and standard deviation of these metrics to assess the overall performance of your

model.

## Model Selection and Tuning (Optional):

If you're comparing multiple models, different hyperparameters, or doing feature selection, you can use the average performance metrics to make informed decisions about which model or configuration performs better.

## Final Model Training:

Once you've chosen your model and configuration based on cross-validation results, you can train your final model using the entire dataset (without splitting into folds).

## Model Evaluation:

Finally, evaluate the performance of your chosen model on a completely separate, unseen test dataset. This gives you an indication of how well your model is likely to perform on new, real-world data.

## Advantages of K-fold cross-validation:

* Utilizes the entire dataset for both training and validation, reducing the risk of bias.
* Provides a more reliable estimate of model performance compared to a single train-test split.
* Useful for small datasets where a single train-test split might lead to unreliable results.
* Remember that cross-validation can be computationally expensive, especially if you have a large dataset or complex models. However, it's a crucial step in ensuring that your model's performance estimates are robust and generalizable to new data.