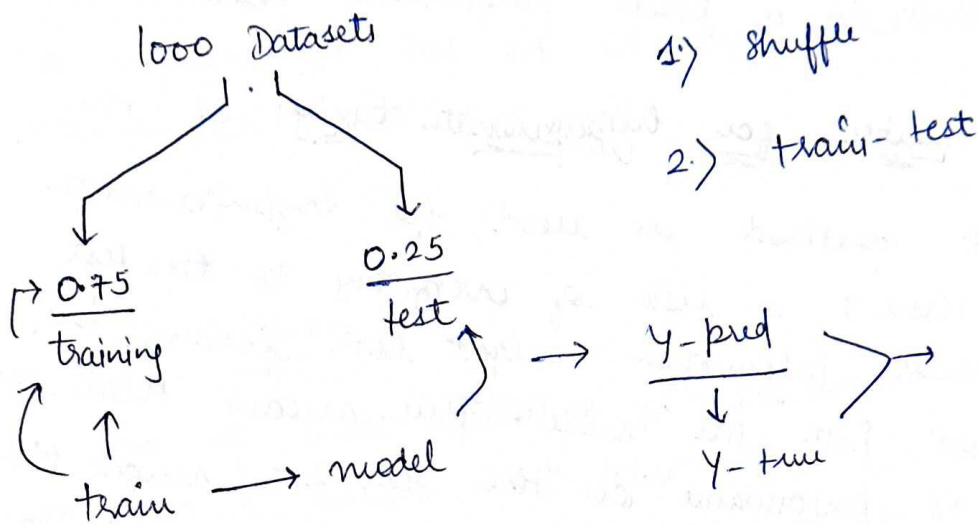


CROSS VALIDATION

The Hold-out Approach → train-test-split



Problem with Hold-out Approach

- Variability: The performance of the model can be very sensitive to how the data is divided into training and testing sets. If the split is unfortunate, the training set may not be representative of the overall distribution of data, or the test set might contain unusually easy or difficult examples. This leads to higher variance in the estimation of the model's performance.
- Data inefficiency: The holdout method only uses a portion of the data for training and a different portion for testing. This means that the model doesn't get to learn from all available data, which can be particularly problematic if the dataset is small.

3. Bias in performance estimation: If some classes or patterns are over- or under-represented in the training set or test due to the random split, it can lead to a biased performance estimation.

4. less reliable for hyperparameter tuning: If the holdout method is used for hyperparameter tuning, there's a risk of overfitting to the test set because information might leak from the test set into the model. This means that the model's performance on the test set might be overly optimistic and not representative of its performance on unseen data.

Why is hold-out approach used then?

1. Simplicity: The holdout method is straightforward and easy to understand. You simply divide your data into two sets: a training set and a test set. This simplicity makes it appealing especially for initial exploratory analysis or simple projects.

2. computation Efficiency: The holdout method is computationally less intensive than methods like k-fold cross-validation. In k-fold cross-validation, you need to train and test your model k times, which can be computationally expensive, especially for large datasets or complex models.

With the holdout method you only train the model once. (2)

3. large Datasets: For very large datasets, even a small proportion of the data may be sufficient to form a representative test set. In these cases, the holdout method can work quite well.

Cross Validation

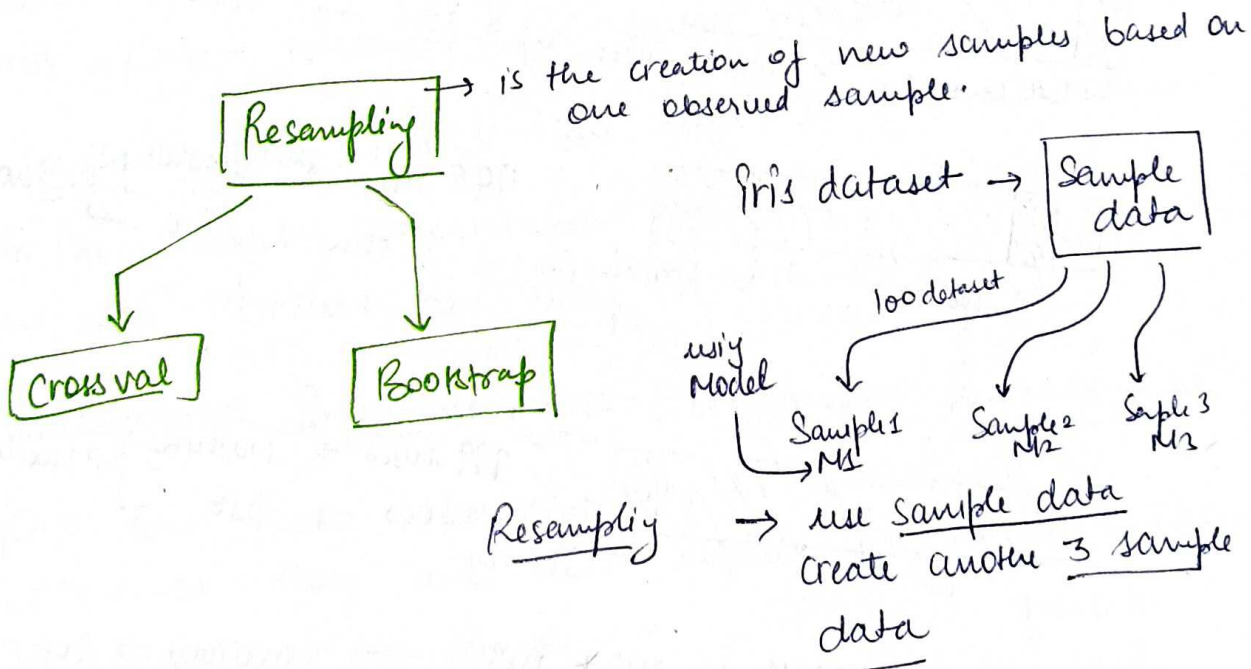
Cross Validation

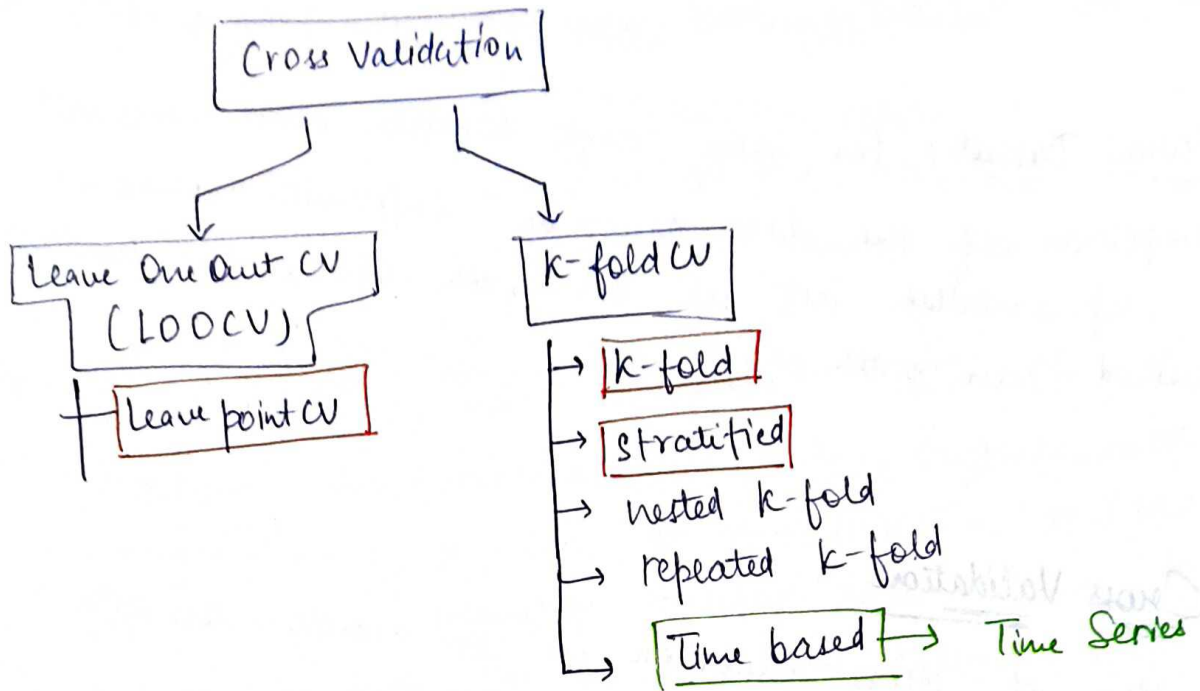
The idea of cross-validation is to divide the data into several subsets or "folds". The model is then trained on some of these subsets and tested on the remaining ones. This process is repeated multiple times, with different subsets used for training and validation each time. The results from each round are usually averaged to estimate the model's overall performance.

bob (observation) \rightarrow estimate

Sampling: - Population Data $\xrightarrow{\text{pop}}$ Observation $n \rightarrow$ estimate
 ↑
 Sample

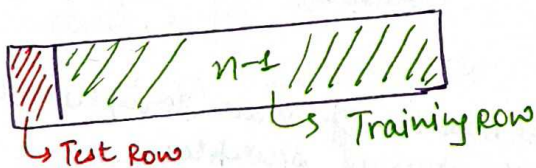
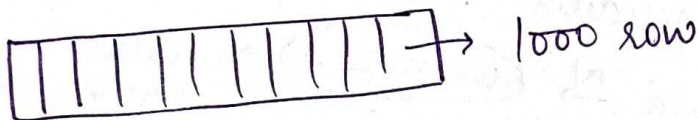
(eg:- Avg Salary of Eng. in India)



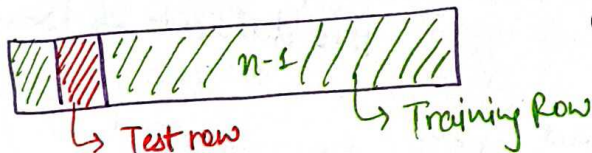


Leave One Out Cross-Validation (LOOCV)

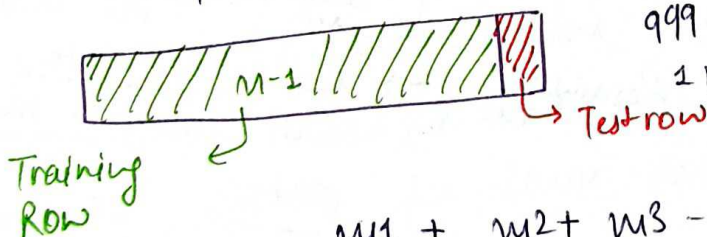
data \rightarrow 1000 rows \rightarrow $\frac{\text{model}}{(1000)}$ \rightarrow model evaluation



999 row \rightarrow Training } model(1)
1 row \rightarrow Test



999 row \rightarrow Training } model(2)
1 row \rightarrow Test



999 row \rightarrow Training } model(1000)
1 row \rightarrow Test

$m_1 + m_2 + m_3 + \dots + m(1000) \rightarrow \text{avg final accuracy score}$

Advantages

(3)

1. Use of data: LOOCV uses almost all of the data for training, which can be beneficial in situation where the dataset is small and every data is valuable.
2. Less Bias: Since each iteration of validation is performed on just one data point. LOOCV is less biased than other method, such as k-fold cross validation. The validation process is less dependent on the random partitioning of data.
3. No Randomness: There's no randomness in the train/test split, so the evaluation is stable, without variation in the results due to different random split.

Disadvantages

1. Computational Expense: LOOCV requires fitting the model N times, which can be computationally expensive and time-consuming for large datasets.
2. High Variance: LOOCV can lead to higher variance in the model performance since the training sets in all iterations are very close to each other.
3. Inappropriate Performance Metric: Performance metrics like R^2 are not appropriate to be used with LOOCV as they are not defined when the validation set only has one sample.

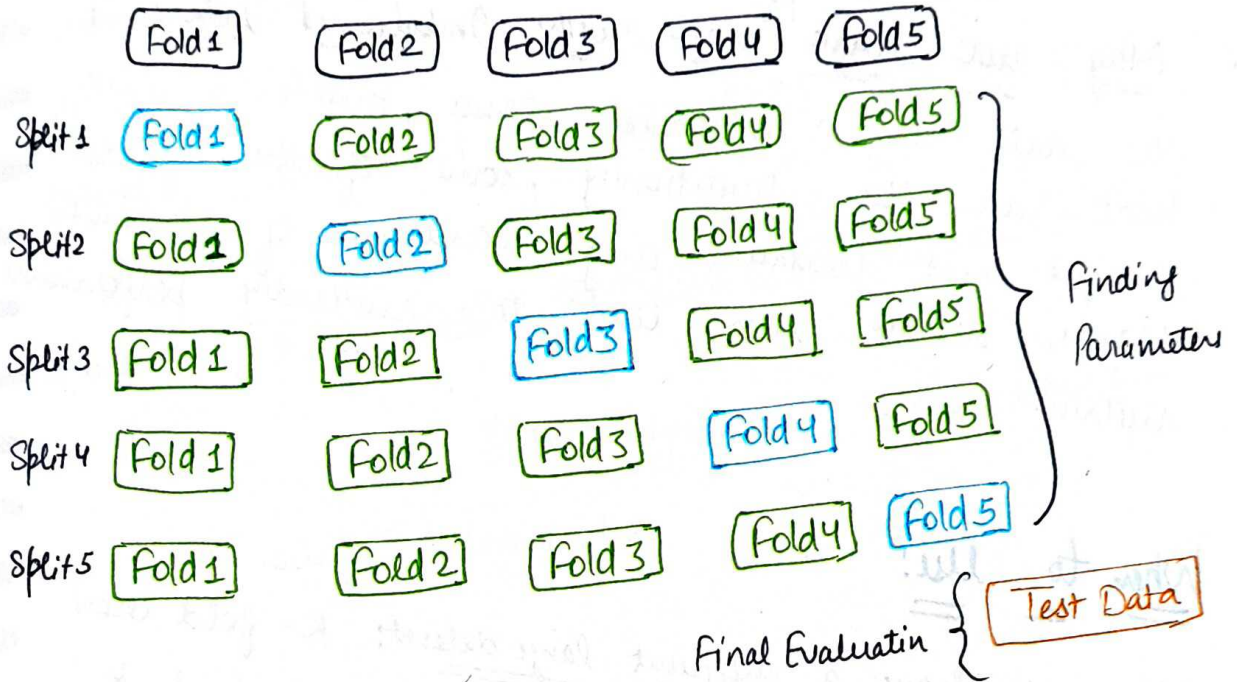
4. Not ideal for Imbalanced Data: In classification problems. If you have imbalanced classes, LOOCV may not provide a reliable estimate of model performance because the single validation sample in each iteration may not be representative of the overall class distribution.

When to use:

1. Small datasets: LOOCV is most beneficial when you have a limited amount of data. With small datasets, you want to use as much data as possible for training to get a reliable model, which is exactly what LOOCV offers by using all but one data point for training.
2. Balanced datasets: LOOCV might not perform well on imbalanced datasets especially in classification problems, because the training set might end up missing some classes. Thus, it's more appropriate to use LOOCV when you have a balanced dataset.
3. Need for less biased performance estimate: Since LOOCV uses nearly all the data for training, it gives a less biased estimate of model performance compared to other methods like k-fold cross-validation.

K-Fold Cross Validation

9



Advantage

1. Reduction of variance: By averaging over k different partitions, the variance of the performance estimation is reduced. This is beneficial because it means that the performance estimate is less sensitive to the particular random partitioning of the data.
2. Computationally Inexpensive: Take less time and space in comparison to LOOCV.

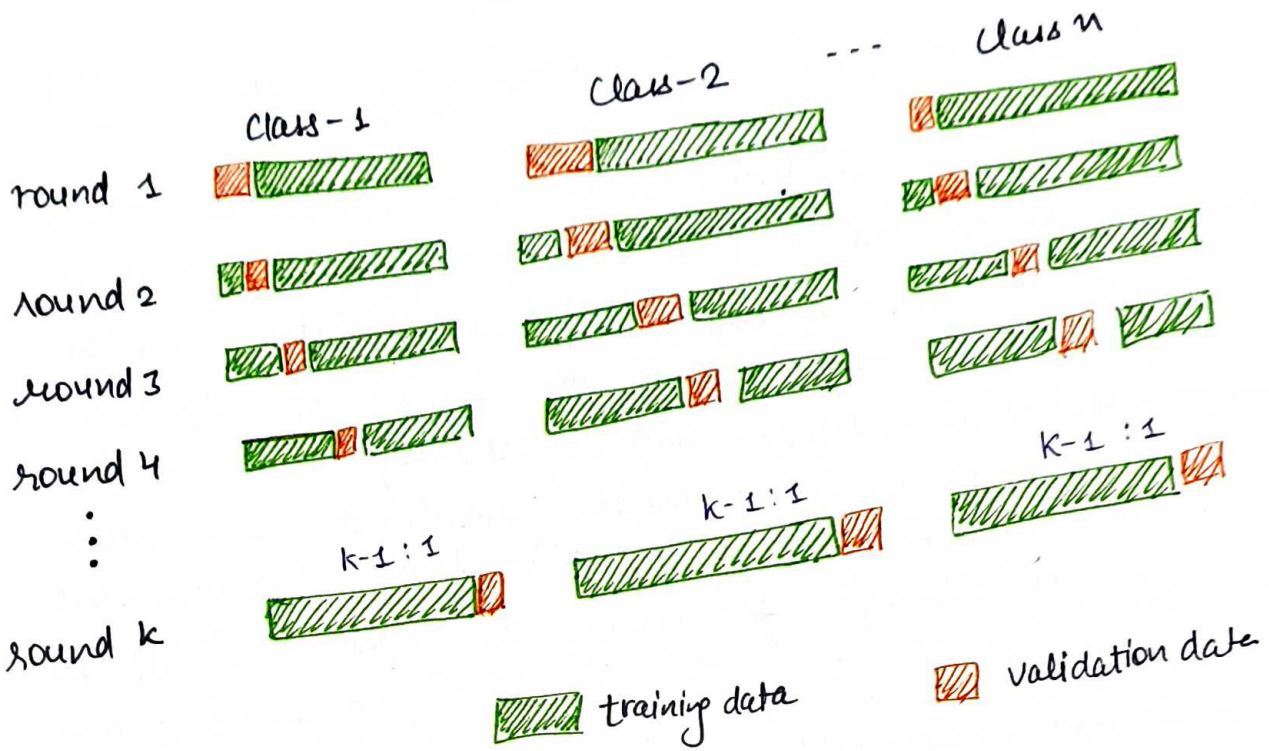
Disadvantage

1. Potential for High Bias: If k is too small, there could be higher bias if the test set is not representative of the overall population.
2. May not work well with Imbalanced data: If the data has imbalanced classes, there's a risk that in the partitioning, some of the folds might not contain any samples of minority classes, which can lead to misleadingly performance metrics.

When to use:

1. When you have a sufficient large dataset: k -fold cross validation requires the model to be trained k times, so, it can be computationally intensive. However, if your dataset is large enough, this increased computational cost can be justified by the more reliable estimate of model performance.
2. When your data is evenly distributed: k -fold cross validation works best when the data is evenly distributed.

Stratified K-fold CV



* Work on Imbalanced data

Validation Set

