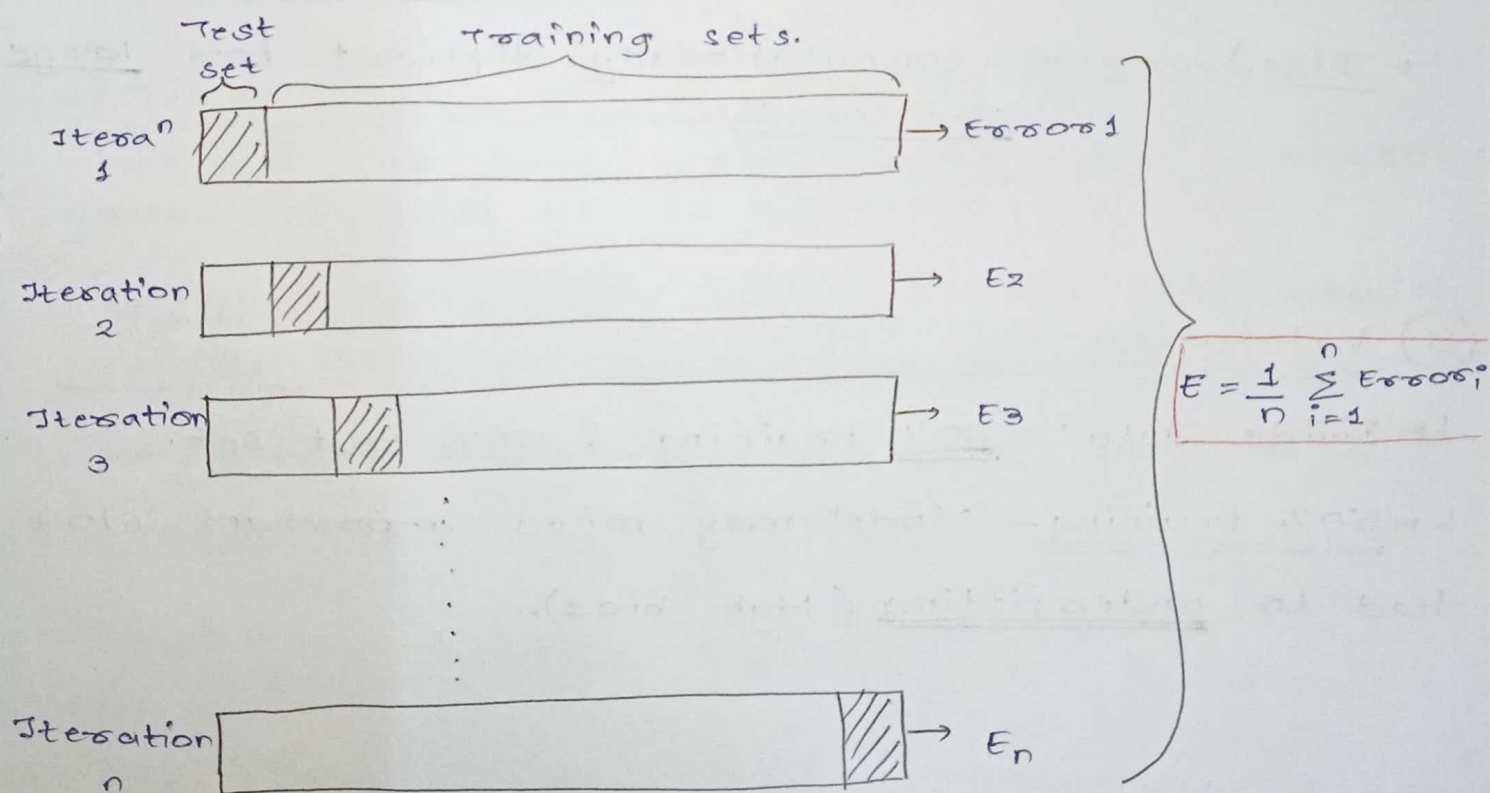# *Cross-Validation*

↳ **Randomly Sampling Data:-**

- Select <u>random records</u> for train & test.
- May result in <u>overfitting</u> (Model train well on train set, but perform poor on test set).

↳ **Cross-Validation:-**

- Divide data <u>multiple sets.</u>
- <u>1 set reserve for test</u> & use others to train.
- Repeat <u>n-times</u>, each time diff test set.
- All <u>results averaged</u> ⇒ performance.
- <u>Prevent overfitting.</u>
- Ensure <u>model robust</u> & generalize well on new data ( for testing).



$$E = \frac{1}{n} \sum_{i=1}^{n} Errors_i$$

# ⊙ Types of Cross-Validation Techniques:—

## ① Leave-One-Out Cross-Validation (LOOCV) —

↳ Take only 1 data point for testing & use other for training.

↳ Advantage — Bcz train on all data ⇒ low bias.

↳ Disadvantage — • Repeat for len(dataset), require high computational time.
  • If test on outlier → High variation in result.

## ② Leave-P-Out Cross-Validation (LPOCV) —

↳ Leave P records for test & others for training.

↳ Can control test set size.

↳ Repeat for all samples & average error.

↳ Disadvantage — computationally difficult for large size of P.

## ③ Validation-Set Approach —

↳ Divide data: 50% training & 50% test set.

↳ 50% training — Model may miss important info & lead to underfitting (high bias).

④ K-Fold Cross-Validation (KFCV) —

↳ Divide data in K-equal sized subsets (folds).

↳ Train & test model K-times, with different test set each time.

↳ Average the errors for performance evaluation.

↳ Easy to understand & output less biased than other techniques.

⑤ Stratified K-Fold Cross-Validation —

↳ Similar to RFCV, but works on stratification concept.

↳ stratification: Rearrange data to ensure each fold/group represent all classes from complete dataset.

↳ Useful for imbalanced dataset.

⑥ Time Series Cross-Validation —

↳ Used for time-series data, when temporary order of data points essential.

↳ Sequentially split train & test sets, with training data preceding testing data.

↳ Useful — Time-series forecast, stock market prediction, weather forecast.

- **Purpose of Cross validation:**

① Estimate how well model perform on the unseen (testing) data.

② Detect & prevent overfitting (occur when model train well on train set & perform poor on test set.

- **Advantages of cross validation:**

① Robust model — More robust than single train-test split. Bcz averages evaluation.

② Maximize Data Utilization — Use each data for both training & testing, imp for limited data.

③ Bias Reduction — Reduce bias from single split.

④ Overfitting Detection — when model perform well on train set, but poor on test set.

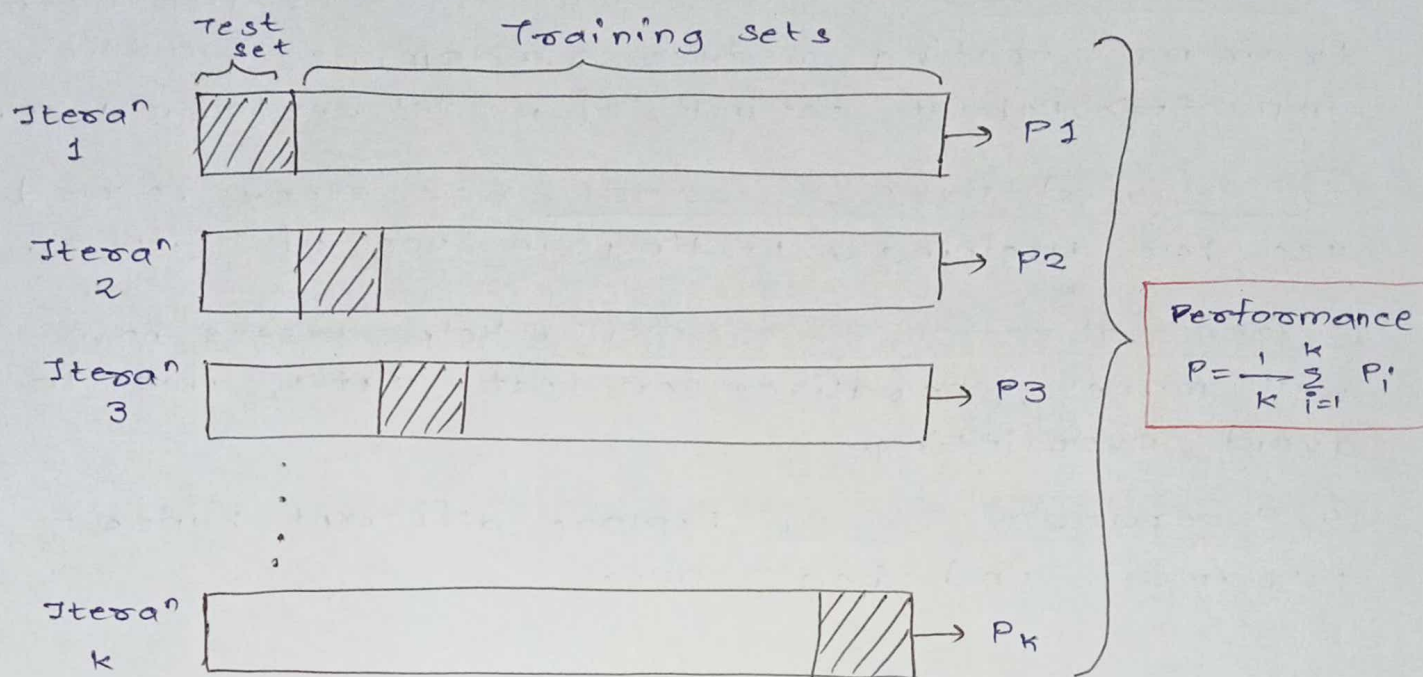⑤ Hyperparam Tuning — Help explore different hyperparams & choose one with best result.

- **Limitations of cross validation:**

① computational cost — n iterations, large datasets.

② Not suitable for time-series data — order is imp. Should consider time-based data splits.

③ Randomness — Effectiveness depend on randomness.

④ Imbalance Datasets — For highly imbalanced data (1 class dominate all others) ⇒ should use stratified KFCV.

⑤ Data Leakage — when not handle properly lead data leakages.

# 1) K-Fold Cross-Validation (KFCV):-

↳ Used to evaluate performance of predictive model.



$$P = \frac{1}{K} \sum_{i=1}^{K} P_i$$

## ◉ Steps in KFCV :-

① **Dataset Splitting:** original dataset divided in K equally-sized, non-overlapping subsets/folds.

② **Training & Testing:** Model train & evaluate k times. In each itration, use 1 for testing & others trainin Repeat process k-times.

③ **Performance Evaluation:** In each iteration, evaluation metrics (accuracy, precision, f1-score) recorded on each test set.

④ **Aggregate Metrics** — After k iterations, all performance metrics averaged to find single performance estimate.

## ⊙ Advantages of KFCV :—

① <u>Robustness</u> : since use multiple sets for both training & testing, reduce randomness & outlier impact ⟹ Robust estimat of model performance.

② <u>Maximized Data Utilization</u> : Bcz every record use for training & testing in one of K iteraⁿ.

③ <u>Generalization assessment</u> : Help assess how well model generalizes on test data, imp to avoid overfitting.

④ <u>Hyperparam Tuning</u> : Explore different hyper-params & find best one.

## ⊙ Choosing K Value :

↳ Generally between 5 to 10.
↳ Smaller — less time, but more variations.
↳ Larger — Reduce variations, but more time.