→ CV is model validation for assessing how the result of a statistical analysis will ~~generate~~ generalize to an independent data set.
→ used in order to limits problems like overfitting

# Cross Validation :-

→ Cross Validation is a technique which involves reserving a particular sample of a data set on which we don't train the ~~data~~ Model. Later we test the Model on this sample before finalizing the Model.
→ Steps in Cross validation :-

① We reserve a sample data set
② Train the Model using the remaining part of data set.
③ Use the reserve sample of the data set test (validation) set. This will help us to know the effectiveness of Model performance. If our Model deliver +ve result on validation set go ahead with current model. It rocks!

⇒ Common Method used for Cross-Validation :-

① Validation Set Approach :-

→ In this, we divided data set    50% for validation
                                   50% for training

→ Dis :-
   We train Model only on 50% data whereas may be possible we are leaving some interesting inf^n i.e. high bias.

②

② **Leave one out cross Validation (LOOCV) :-**

→ In this approach, we reserve only one data point and train Model on rest data points. This process iterate for each data points.

→ **Adv & Dis :-**

- We make use of all data points i.e low bias
- We repeat the cross Validation process iterate $n$ times means high execution time.
- This approach leads to higher variation in testing Model :: we testing against one data point. If the data point turns out to be an outlier - It can lead to higher variation i.e high variance.

③ **K-fold Cross Validation :-**

- from above 2 methods — i) high bias
                            ii) high variance
So, it will take Care of both.

**Steps :-**

①  Randomly split our entire data set into K "folds"
②  for each K folds in our data set, build our model on (K-1) folds of data set. Then test the Model to check effectiveness of Kth fold.
③  Record the error we see on each predictions.
④  Repeat this until each of K folds has served as the test set.
⑤  The avg of our K recorded errors is called the cross validation error & will serve as our performance metric for the Model

**How to choose right value of K :**

→ lower value of K ⇒ more bias ⇒ undesirable
    large    "    "    " ⇒ more variability.

small value of K leads to validation set approach
higher   "    "    "    "    "      LOOCV    approach

So, always suggested K = 10

**How to measure Model's bias-variance ?**

→ After K-fold Cross validation, we will get k- different
model estimation errors ($e_1, e_2, ... e_k$). In ideal scenario,
these error values should add to zero.

→ For Model's bias, take avg of all errors,
     Lower avg value, better Model.

→ For Model's variance, take standard deviation of all error.
     Lower value of std, our model doesn't vary a lot with
             different subset of training data.

→ So, our focus on achieving balance b/w bias & variance
   i.e, for better predictive Model.
         reduce variance, controlling bias to an extent

→ This trade-off usually leads to building less complex predictive
models.

Model is too simple ⇒ suffer from underfitting
Model is too Complex ⇒ suffer from overfitting