

Short Notes On Lasso & Ridge :-

- * Lasso Regulariser containing absolute value function, this values (absolute) are not differentiable and it contains more computational resources i.e. taking more time to compute.
- * Ridge Regularizers containing Square value function, these values are differentiable and it doesn't contain more time to compute because it is a simple resource.
- * We have an advantage in Lasso reg. is it contains feature selection. (i.e. whatever feature is not important it contains its coefficient as zero for unwanted features)
- * Ridge reg. has not this advantage of feature selection, both the regularizers perform same task to resolve the issue of "Overfitting".

General Terminologies :-

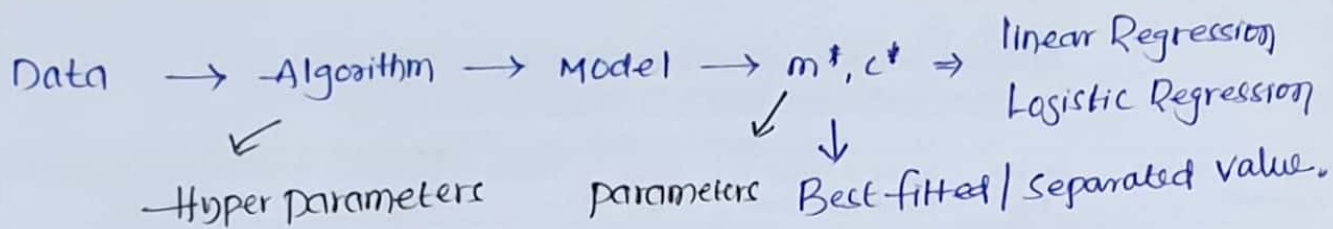
① In K-NN Regression/Classifier we have to choose 'k' value but how to know which 'k' we have to use, by experimentally ^{only} we can say that we have to choose this k value. In other terms we can't decide 'k' value.

② Eqn - $m_1 = m_0 - \eta \left[\frac{\partial f}{\partial m} \right]_{m_0}$ is update function of gradient-

descent. where (η) is learning rate, what ' η ' value we have to choose in eqn again we have to do experiments with different values. by experimental only we can decide the values.

from the above (k - in KNN, η - GD, λ - Lasso, Ridge) all are Hyper parameters which we tune the algorithms.

Difference b/w parameter & Hyperparameter :-



Parameter is nothing but output of our learning algorithm (ie model)
in the above case (m, c) are parameters.

Hyperparameter is nothing but which are tune the algorithm / model with
this we solve the issues of "overfitting". by this tuning the hyperparameters

we will achieve the Low bias and low Variance (not too low var).

Parameters comes from the model. (output), where Hyperparameters
are part of an Algorithm.

Linear Regression

$$(y_{act} - y_{pred})^2$$

$$\arg_{w, c} \left\{ \min \sum_{i=1}^n (y_i - (w^T x + c))^2 + \lambda \sum_{i=1}^d |w_i| \right\}; \arg_w \left\{ \min \sum_{i=1}^n (1 + \exp \{-y_i * (w^T x)\}) + \lambda \sum_{i=1}^d |w_i| \right\}$$

Lasso Ridge

$$\lambda * \sum_{i=1}^d |w_i| \quad \lambda * \sum_{i=1}^d (w_i)^2$$

Logistic Regression

$$(y_{act} * y_{pred})$$

Lasso Ridge

$$\lambda * \sum_{i=1}^d |w_i| \quad \lambda * \sum_{i=1}^d (w_i)^2$$

$d \rightarrow$ dimensionality

for Suppose if we consider Regularizer-term in either linear Regression
or logistic regression where we use either Lasso or Ridge Regularizer
if we consider λ value as zero (ie $\lambda = 0$)

Then The total term of regularizer will be zero. as shown
 in eqn., then we get only linear reg-term or logistic reg-term.
 ie our model having "overfitting" problem. if we have to improve
 lambda (λ) value then we get Underfitting problem. because there
 only regularization term value increases then no impact of linear and logistic

$\lambda \uparrow \rightarrow$ Under-fit $\lambda \downarrow \rightarrow$ Overfit

This means there is not maintain high lambda value and low lambda.
 Value. keep change the values of λ ie we have to tune every time with
 λ value in order to get best model.

Q - what if Overfit and Underfit a problem?

A - when ever we have a Overfit or Underfit issue, we have to
 use Regularizer term, this introduce the ' λ ', which tune our entire
 algorithm, ie it helps to our model from Overfitting/Underfitting issue.

Hyper parameter tuning in SVM :-

In SVM our task is to minimize the Margin,

$$\text{S.V.M eqn} - \min \left\{ \frac{1}{2} \times \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \xi_i \right\}$$

In the above eqn $\rightarrow c$ is hyperparameter, if " $c=0$ " then

$\min \left\{ \frac{1}{2} \times \|w\|^2 \right\} \Rightarrow$ which simply means we maximize the

margins. then $\max \left\{ \frac{2}{\|w\|^2} \right\} \Rightarrow$ So this eqn is nothing but "hard

Margin" S.V.M equation. the main problem with this eqn

Is to the data is Linearly Separable. without any misclassifications (ie there is misclassification but we assumed it).

So this hard margin SVM eqn gives us Underfit Issue because Initially we consider there is no misclassifier point but In real world we didn't get like this data, Some data points overlapping with other points.

Now we consider $C=0$ then the cost term will be 0 and we left with only Hard margin eqn and that gives us Underfitting Issue.

$C \downarrow$ - Underfitting $C \uparrow$ - Overfitting.

from this we can say if our 'C' value decreases we face Underfitting issues, if we increase C value we get Overfitting Issues.

In SVM we have kernel's also ~~we~~ if we consider R.B.F there is kernel value will be -

$$\exp\left\{-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right\}$$

So here ' σ ' is hyperparameter, So we can represent it with

$$\gamma \text{ also } \Rightarrow \exp\left\{-\gamma \|x_1 - x_2\|^2\right\}$$

$$\text{where } \gamma = \frac{1}{2\sigma^2}$$

Note: ① if we consider linear SVM eqn \rightarrow we have only one hyperparameter that is 'C'

② if we consider kernel SVM with R.B.F \rightarrow we get two hyperparameters are - C, γ So In this case we have to tune two values where In simple SVM eqn we have to tune only 'C' value.

Hyperparameters with different algorithms :-

① Linear or logistic Regression $\rightarrow \lambda \downarrow$ - Overfit problem
 $\lambda \uparrow$ - Underfit problem

② Support vector Machine $\rightarrow C \downarrow$ - Underfitting ; $\gamma \downarrow$ - Underfitting
 $C \uparrow$ - Overfitting ; $\gamma \uparrow$ - Overfitting

③ k-Nearest Neighbour $\rightarrow k \downarrow$ - Overfitting
 $k \uparrow$ - Underfitting

In kNN, we have two hyper parameters are - k, p .

where 'p' is a distance matrix when $p=1 \Rightarrow$ Manhattan distance
 $p=2 \Rightarrow$ Euclidean distance
 $p=3 \Rightarrow$ Minkowski distance

In Order to get best model In kNN, we have to tune both (k, p) Parameters by experiments like cross validation methods.

④ Decision Tree \rightarrow Depth \uparrow - Overfitting ; No of points in leaf
Depth \downarrow - Underfitting ;

In D.T we are memorize the entire data, where we mention each and every condition that gets more complex and increases depth of the decisions i.e. deeply going to conditions. So that is considered as a Overfitting Issue.

⑤ Random forest \rightarrow Depth + No of D.T + No of points in leaf

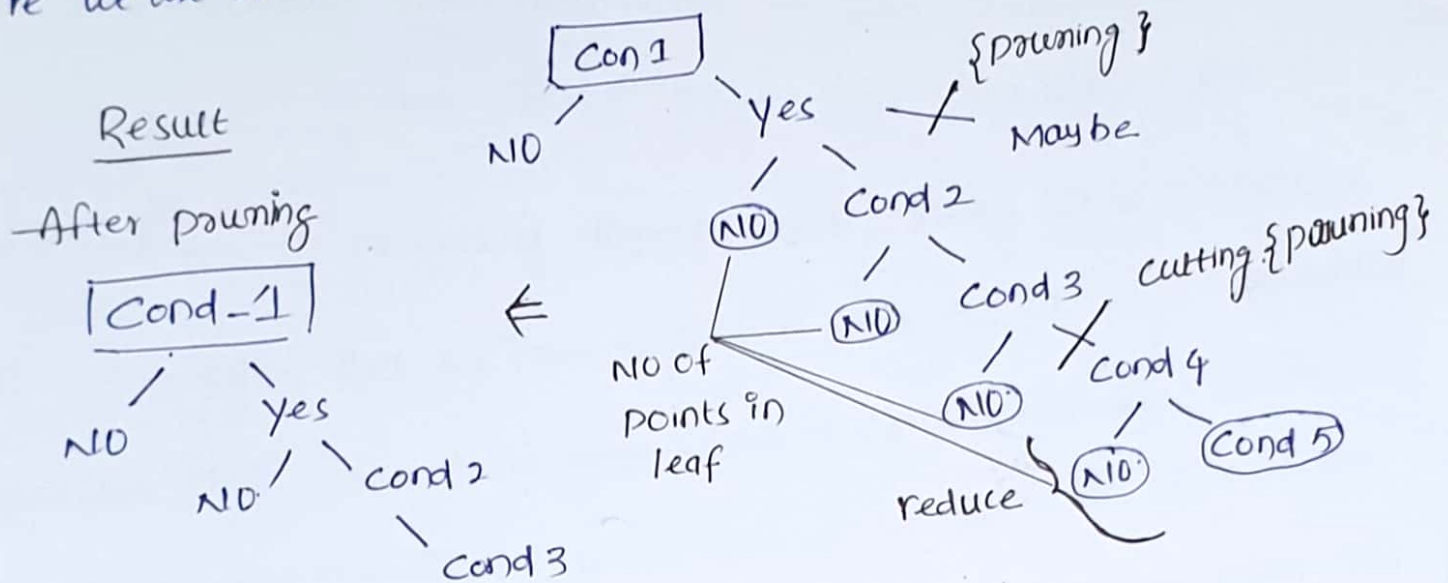
In Random forest also something in Decision tree because in it also we have decision making while training data.

In Random forest \rightarrow Base learners (D.T) + Row & Column Sampling + Aggregators
 \hookrightarrow Overfitting \rightarrow Depth (\uparrow)

that gives best Model.

Q - How Do you treat Overfitting in DT?

A - when we create a Decision tree for Suppose it has issue of overfitting (ie we have high depth) it contains more decisions making after getting deep tree, we have a technique called as pruning. by this we cut down the depth of DT, then the remaining part will be lost- ie we are reduce the depth of decision tree.



Q - How do you treat the Overfit Issue in DT?

In Decision Tree we have to use Regularizer which are

- ① Pruning
- ② changing the hyperparameter \rightarrow H.P for DT is \rightarrow Depth $\{$ # of points in leaf.
- ③ Random forest

Note :- NO of points in leaf's means - terminal end points in DT

In case of linear Reg - we use $\lambda \rightarrow \lambda$

logistic Reg - λ

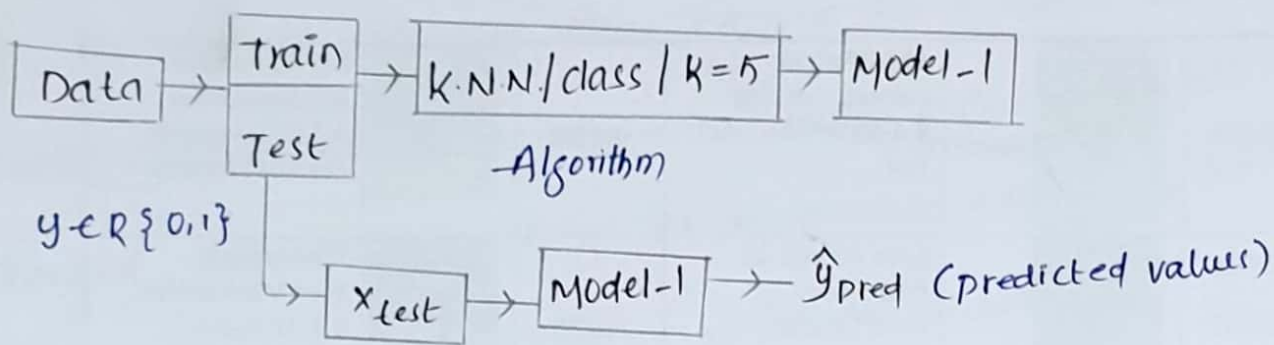
SVM - C, γ (where RBF kernel)

Gradient Descent - η
Hyperparameter

In this all algorithms after tuning we have to again tune the values with Gradient descent ' η '

Hyper parameter tuning process in k.N.N. Algorithm :-

17



Accuracy \rightarrow Train data $\rightarrow y_{trn}, \hat{y}_{trn-pred} - 90\%$
 \hookrightarrow Test data $\rightarrow y_{tst}, \hat{y}_{tst-pred} - 60\%$ } Overfitting

Step-1 :- In process of hyperparameter tuning let us consider a dataset after that we have to split it into two parts as train data, test data. then take train data which is classification type i.e. the 'y' target variable is categorical type so we have to do classification task. So here we are applying k.N.N. Classification algorithm to train data and consider hyperparameter "k" as 5 and that gives a model which we consider as "model-1", after making a model, we take x_{test} data value and put it in model that gives us predicted values of "y".

Step-2 :- Now we have to use Accuracy matrices to calculate the accuracy of a test data as well as train data. for suppose consider train data $(y_{trn}, \hat{y}_{trn-pred})$ that gives us 90% accuracy. Similarly consider test data $(y_{tst}, \hat{y}_{tst-pred})$ that gives 60%. then we can understand our model has Overfitting Issue.

19
Accuracy-3 \rightarrow Train-data \rightarrow 50% } Underfitting.
 \hookrightarrow Test-data \rightarrow 60%

Now we get train data accuracy as 50% and test data as 60%.
In this test data has more accurate values but in train data we didn't have good accuracy this gives us result as Underfitting.

If we observe accuracy result of train and test, test has more accuracy as compare train but with 60% result we can't predict well. So we have to improve our accuracy in order to get 90% result in test prediction. then we can say that is good model.

So again we have to do trials in order to get best model, let consider $k=13$ and do step-1 we get model-4. if it shows values as -

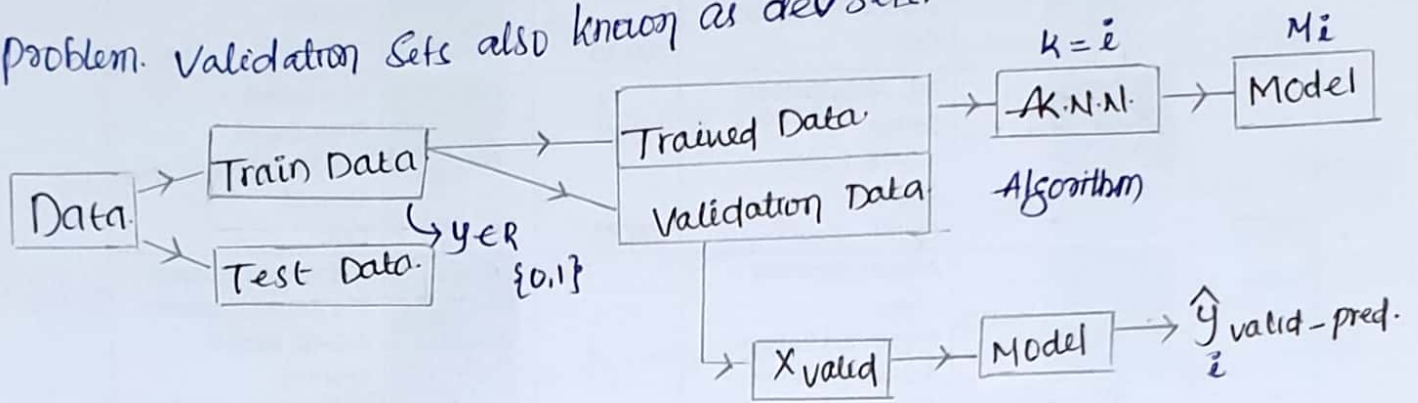
Accuracy-4 \rightarrow Train-data \rightarrow 90% } Best model
 \hookrightarrow Test-data \rightarrow 91%

So here " $k=13$ " where model-4 gives values of accuracy is predicted well in train data as well test data. So we can say that our model-4 is "Best model" which has no issues of overfitting and Underfitting.

The result "Best model" whatever we get by change of k value ie here every time we tune the model by hyperparameter ' k '.
So here we get the best model in order to change ' k ' value by every time we observe train and test data which is unseen but in real time there is no use of test data (ie we have to check train data accuracy values) only.

There we use the Concept of Validation. In this we Split²⁰ train data in to two parts as validation data and train data.

Validation Set :- it is a Set of data used to train model (AI) with the goal of finding and optimizing the best model to solve a given problem. Validation Sets also known as dev sets.



Step-1 :- first of all we have to choose a data and Split the data in to two parts as train & test data after that keep test data as it is. Consider trained data. Split this trained in to two parts as trained and validation data. after this we have to choose algorithm in order to do this we have to check our target variable if it is a numerical value contain do regression task, for Suppose if we consider categorical then do classification task. So here I choose algorithm. K.N.N. and apply Our data to this algo. we get model after it we have to take "X-Validated-data", apply to a model that gives ' \hat{y} ' pred- values.

Step-2 :- now here we consider Accuracy matrices to get Accuracy of a model in order to do it we consider both datas and apply

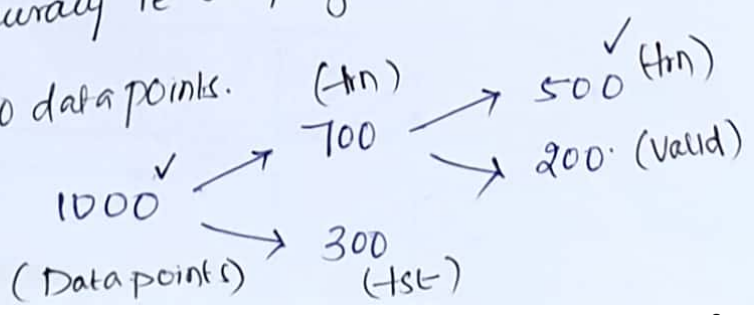
matrices to that for Suppose initially i considered $k = 5$

Accuracy \rightarrow Train data $(y_{\text{trn}}, \hat{y}_{\text{trn}}) \rightarrow 95\%$
 \hookrightarrow Validated data $(y_{\text{val}}, \hat{y}_{\text{val}}) \rightarrow 85\%$ } Overfit problem.

Now here we calculate accuracy of both train and validated data. the get good result but as compare to train, valid data has low accuracy this will get the issue of Overfit.

Step - 3 :- In Order to reduce it we have to do hyperparameter tuning. So in this algorithm 'k' is a hyperparameter So every time we have to change the value of 'k' ie we have to keep trials on this data. Until unless we get "good model" from $(k = 1, 2, 3, \dots, i)$ we get models as (Model-1, Model-2, ..., model-i) & Accuracies are also respectively. So from this models we pick easily "best model". as similar as hyperparameter tuning on train and test data.

But here One biggest drawback of validation process is we take train data and further divided as validation, train data. So here we decrease the train samples from train data that will be effect on result like accuracy ie Overfitting (or) Underfitting issues majorly. for Suppose we have 1000 data points.

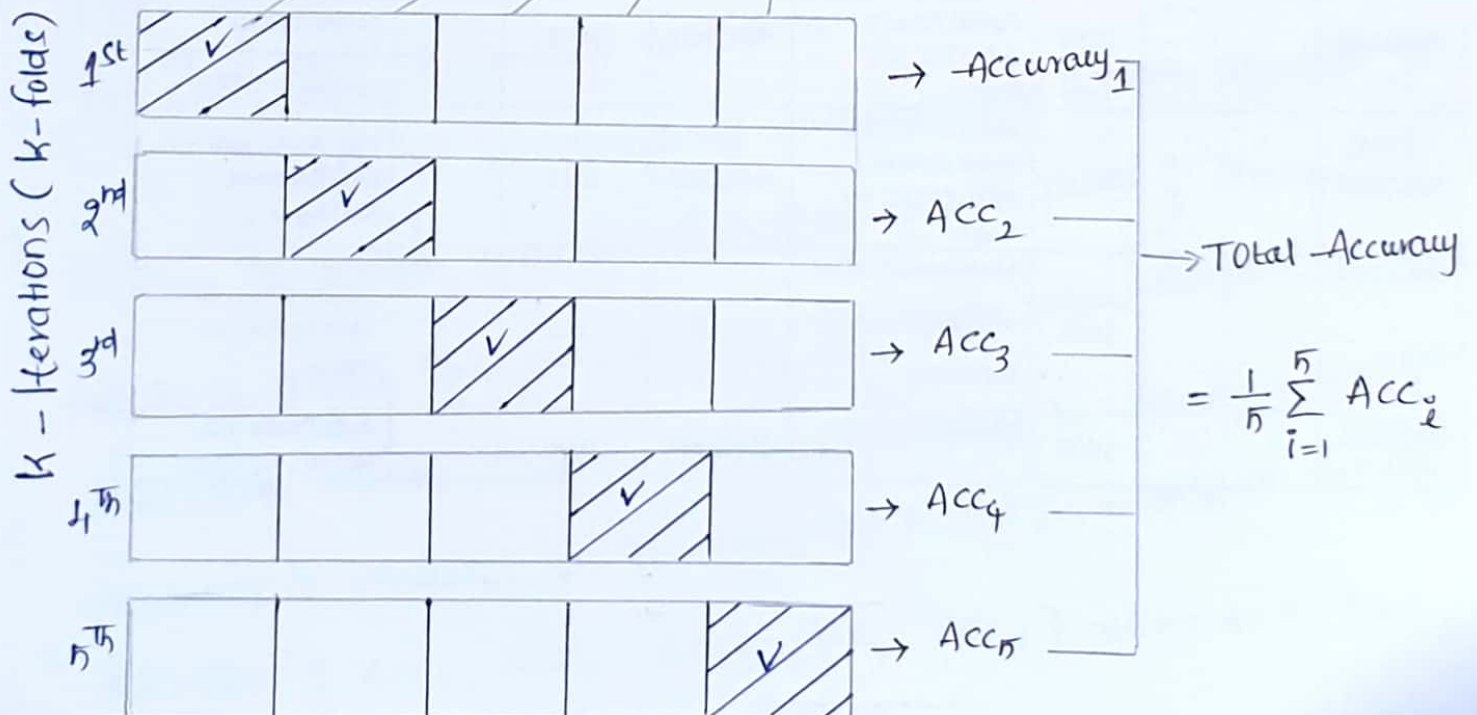
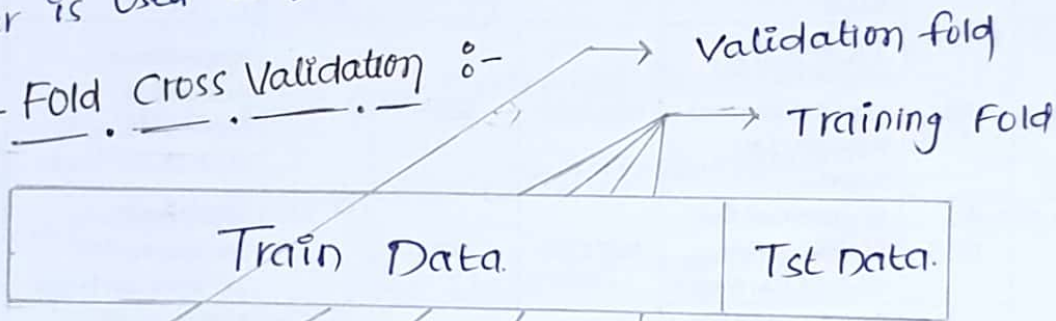


So if we look in to train data points initially they are 1000 data points after validation split we are left with 500 data points only is half of the data points, but we know that In order to get a best model we need to train more data points, then only we get best performances like Accuracy, the resultant tells us whether the model is best or not.

"In Order to prevent this problem we have to use a technique called as Cross Validation."

It is a Statistical method of evaluating and comparing learning algos. by deviding data into two segments One is used to learn or train model Other is used to validate the model.

K - Fold Cross Validation :-



In this validation process firstly we take a dataset and divide it in to two parts as train data and test data now consider train data and Split it further k -folds in this we have to take k -value as $\{5, 7, 10, \dots\}$ any great value but in generally Consider " $k=10$ " here this k is different from KNN, so this ' k ' will be divided that part in to such many kinds, here we considered $k=5$ now we will get " 5 -folds" ie train data Split in to 5 many times with 5 parts in train data. So far Suppose if we divided Split Initially data as 80% train and 20% test, this 80% divided equally to all train parts as shown if in-figure ie 16% equally divided each part.

So if we looking in to the diagram in 1st iteration of train data. the first part will Considered as Validated data and remaining all are Considered as train data, Similarly in 2nd iteration the 2nd part will Considered as Validated remaining all are train data. lets do it Untill the 5th iteration train data.

After it we Considered train data and applying algorithm to it then it makes a model, then after Consider values of x_{train} and put in to model that gives us $\hat{y}_{pred-valid}$ ~~observed~~ value. Compare it with y_{val} then finally applying performance matrix ie Accuracy we get Accuracy of model-1 (or) iteration 1st result, Similarly do for all we get $\{Acc_1, Acc_2, Acc_3, Acc_4, Acc_5\}$ where $k=5, p=1$

Then finally we do summation of all these values we get

24

total Accuracy at $\{k=5; p=1\}$

So if we observe here we consider $k=5$ (which means KNN points) and $p=1$ (Manhattan distance), which are hyperparameters, and ' k ' is a bold in cross validation. here $k=5$ we decided so we get '5' models with 5 iterations. So if we change the values of $(k=5, 7, 9; kNN)$ and $p=(1, 2)$ Initially consider $\{k=5; p=1\}$. Now we take $\{k=5; p=2\}$ again we get 5 iterations with 5 models. So now look in to how many pairs we have there are — $k = \{5, 7, 9\}$ } \Rightarrow Totally we got '6' pairs.

$p = \{1, 2\}$

for each pair we have '5' models because of our k -fold is '5'

from this we got total '30' models internally when we apply " k -fold".

In practice it is very difficult to train the '30' models, In order to overcome it we have to use some techniques in k -fold cross validation. those are —

- ① Grid Search Cross Validation
- ② Random Search Cross Validation.

Grid Search C.V :- In this process we have to grid all the parameters

what ever we have $\{k, p\}$

$\Rightarrow k = \{5, 7, 9\}$

$p = \{1, 2\}$



So in Grid Search Cross Validation firstly we have to make a Structure of Grid Containing all values of hyperparameters (k, p) In case of k -NN Algorithm as shown in figure, then we have to train model as- $k=5$ and $p=1$, then train another model where $\{k=7, p=1\}$ Similarly we have to do it till $\{9, 2\}$ as shown in figure, where we train total data by taking different hyper parameters, every time this will be shuffle data where you mention random state or not.

Randomized Search Cross Validation :-

R.S.C.V is better than the Grid Search C.V. in terms of Computational Complexity. In Grid Search C.V. it computes each and every pair that means in above example we have '30' models all this models compute by Grid Search C.V. but in Randomized Search C.V. its don't do all models computation. just randomly select the models to train then it will pick '5' randomly selected models, in case-II select 10 models, in case-III select 8 models for suppose, all this selection is done by randomly. and do cross validation

\therefore Both the techniques do same operation but in Grid do all models validation (ie train), in case of R.S.C.V it pick randomly the model as compare to Grid, random.S.C.V has less Computational Complexity, So from this we can say R.S.C.V is better than G.C.V., Generally if we have low models pick. G.S.C.V

	k		
	5	7	9
$p = 1$		✓	
2			✓

$$k = \{5, 7, 9\} ; p = \{1, 2\}$$

out of this pick Randomly '2' models only.