

→ where @ $k=n$, It is Called Underfitting because, The Classifier is Underworking because its simply depends on majority points and not on the surface we draw or something

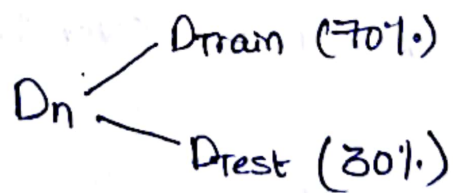
→ where as $k=5$ It is well-fitted for

Q.13 Need for Cross Validation

→ As we see depending upon the k -values we are Overfitting, under fitting for.

→ To determine 'K' or which k value will be suitable one idea is:

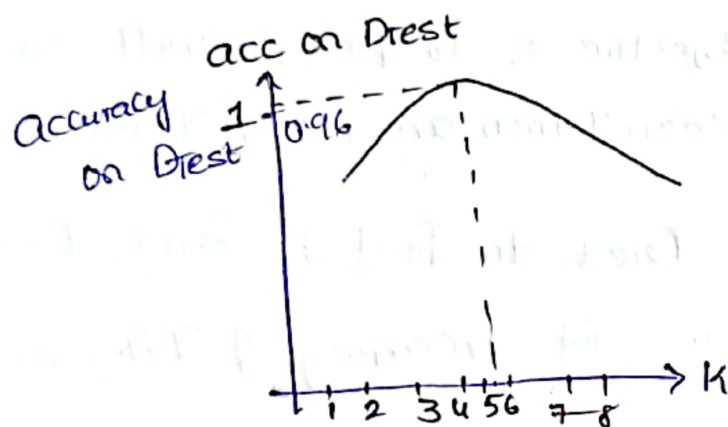
Consider D_n which is whole data set split this Dataset into D_{train} (70%) D_{test} (30%)



Now using D_{train} data by getting the accuracy on D_{test} for different k values.

	Train	accuracy on Dtest	No of Correctly classified Pts Total Points
$k=1$	D_{train}	0.78	
$k=2$	"	0.82	
$k=3$	"	0.85	

So we will check for different K values we will try to get the accuracy on D_{test} with the help of D_{train} . The K value on which the accuracy which is typically higher will be considered



let $K=6$ gives me the best accuracy on D_{test} when using D_{train} as Training data

we can say as

using D_{train} & 6-NN on Amazon reviews dataset

* I can get 96%. (if accuracy is 0.96)

* But there is a small problem.

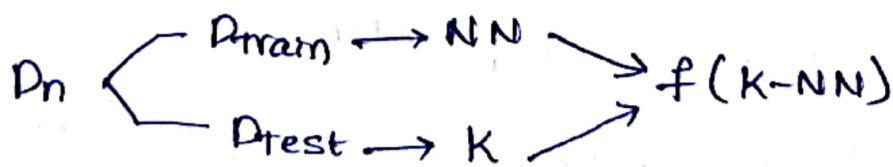
as in General our objective is if we have the Dataset D_n The whole dataset D_n can be divided into two subsets D_{train} & D_{test} using D_{train} & D_{test} we tends to create an algorithm or function

and for that function we try to get y_2 by giving x_2 . That means we try to apply the algorithm on future or unseen points. which is called "generalization."

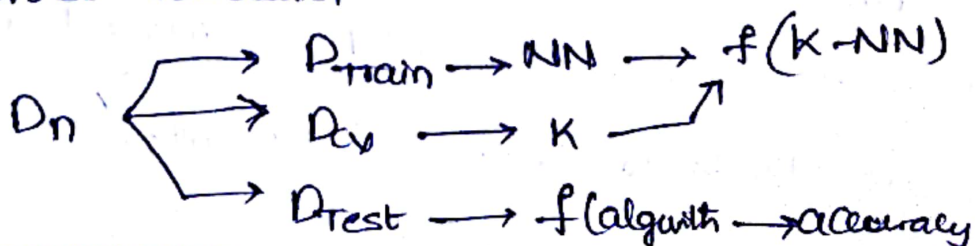
Therefore as our objective is to perform well on unseen pts. (unseen! which are not in Train & Test) and as we used D_{test} to find K and D_{train} to find NN as if we got accuracy of 96% on D_{test} by using "K"-NN we cannot say that my accuracy on future dataset or x_2 would also be 96% by using same K .

To overcome this there is concept called Cross-Validation (CV)

as in General case we used



In Cross Validation



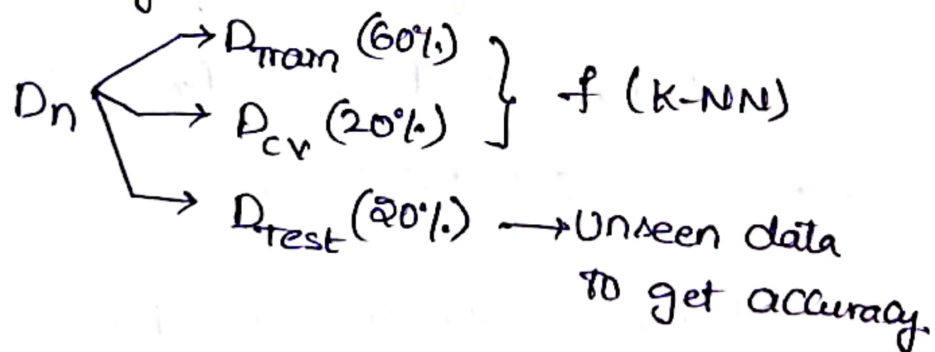
→ Now if we evaluate each point x_i in D_{test} as x_q and accuracy 93%. (as on unseen data)
 then we can say that 6-NN (k-NN) has an accuracy of 93% on unseen data, and this is called generalization accuracy.

and $100 - 93 = 7\%$ is the generalization error on ~~on~~ unseen data.

Q.14 K'-fold Cross Validation:-

K-fold Cross Validation is where we combine D_{train} and D_{cv}

→ Generally we use whole data set D_n

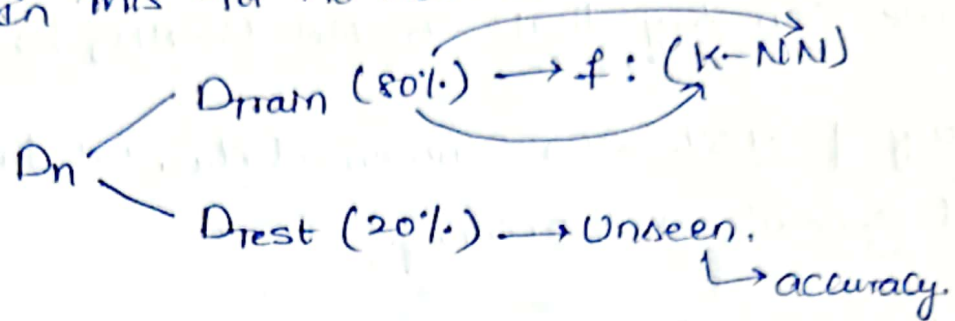


→ But in machine one of the key point is more the training data the better is your algorithm.

→ In k-fold cross validation we combine D_{train} and D_{cv} and total 80% data will be used to compute NN

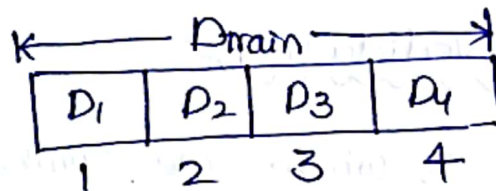
→ The k in k -fold CV and k in KNN both are different

Step 1:- So in this for the Dataset D_n



In Step 2:-

we will try to break the whole D_{train} into 4 equal sized parts



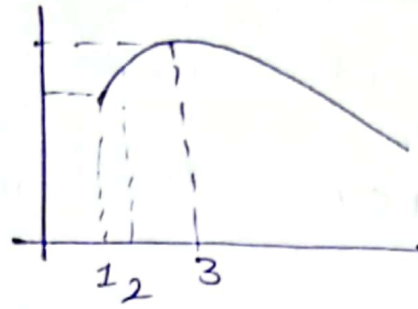
Step 3:-

we will find the accuracy on CV data set.
as we only divided D_n into D_{train} & D_{test} then how could we get CV data set?

So as we divided D_{train} into 4 equal parts

for different k values we will find the avg of accuracies and which ever ^{'k'} gives the highest avg accuracy value that k needs to be considered

→ let us consider our accuracy on data set is low



		Train	CV	acc on CV
4 times	$K=1$	$D_1 D_2 D_3$	D_4	a_4
	$K=1$	$D_1 D_2 D_4$	D_3	a_3
	$K=1$	$D_1 D_3 D_4$	D_2	a_2
	$K=1$	$D_2 D_3 D_4$	D_1	a_1
				$\left. \begin{matrix} a_4 \\ a_3 \\ a_2 \\ a_1 \end{matrix} \right\} \text{avg} (a_1, a_2, a_3, a_4) = a_{K=1}$
by we find for $K=2$	$K=2$	$D_1 D_2 D_3$	D_4	a'_4
	$K=2$.	.	.
	$K=2$.	.	.
	$K=2$.	.	.
				$\rightarrow \text{avg } a_{K=2}$

let the avg value we got is higher at $K=3$

→ as we are repeating for every K value four times so in K -fold $K'=4$ here, so it is 4 fold C.V.

→ let $K=3$ has highest accuracy so we decided $K=3$ as best classifier so we get function as

$f: 3\text{-NN}$

→ So $f: 3\text{-NN}$
 \uparrow $\rightarrow D_{\text{train}}$
 avg
 accuracy on
 4-fold CV

→ So if we given a review x_q we will use whole of Training data to compute NN and we use 3 NN to decide Class labels.

$x_q \xrightarrow{\text{NN}} D_{\text{train}}$
 \rightarrow 3 nearest neighbours

→ To calculate accuracy of the model.

we use $D_{\text{test}} \rightarrow$

$D_{\text{test}} \rightarrow$ Accuracy of 3-NN on D_{test}

let accuracy we got is 93%. on unseen data

Note ↓

So avg-accuracy and accuracy of model is dif.

By using avg accuracy we are find best k values depending upon which ever avg-accuracy value is h and will be calculated with D_{train} .

→ where as accuracy of the model will be calculated by Dtest

→ Note:-

Mostly used k' -fold CV is 10-fold CV. typically.

→ Note:-

Time it takes to compute the optimal/best k in KNN increases by k' times if we use k' -fold CV.

2.15 Visualizing, Train, CV & test datasets



when we are considering the whole data set and we are randomly sampling the data set into

