

Validation data is used to evaluate the loss of a function h that is determined using the learning on the training data-set. If the loss on validation data is high for a given h , the hypothesis or model needs to be changed.

↳ helps improve the model by selecting the best configuration without overfitting

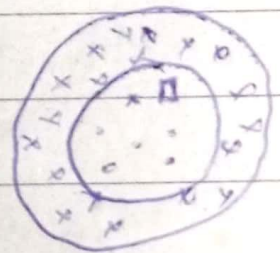
↳ used to tune hyperparameters (like value of K in K-NN algo)

"K - NEAREST NEIGHBOUR ALGORITHM"

Space & Time complexity $O(dn)$ ^{no. of feature (dimension)} _{no. of data points}

↳ used to assign labels to unknown data points

	x_1	x_2	class	Distance
training data	1	5	0	2
	0	8	0	6
	0	6	1	4
	1	2	1	1



Test (1,3) ^{calculated distances}

↳ we need to identify the class

Date _____

Day _____

K=1 means only 1 nearest neighbour check kserge

2 2 1 two 4 4 4 4

$$1 < \textcircled{K} \leq N$$

↳ Hyper parameter

↳ For binary classification $K = \text{odd}$

↳ For multiple classification $K = \text{even}$

→ Complexity of prediction increases with size of training data.

→ Two types of distance Metric

- Euclidean distance = $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

- Manhattan distance = $|x_2 - x_1| + |y_2 - y_1|$

- Minkowski distance = $\left[(|x_2 - x_1|)^p + (|y_2 - y_1|)^p \right]^{1/p}$

$p=1$
↳ if $p=1$ is Manhattan dist
 $p=2$ is Euclidean dist

Date _____

Day _____

→ Miss match in the values of data

This can be solved thru normalization.

• Choice of K :

▷ $K=1$

↳ if $K=1$, the model makes decision based on just one data point. If that point is incorrect, unusual or just slightly different, the entire prediction changes. This leads to overfitting

▷ $K=n$

↳ The algo assigns the majority class of the entire dataset

↳ It ignores pattern & predicts the "most frequent" class

↳ leads to underfitting

↳ where model is too general & performs poorly on both test / training data.

$$s = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n(\sum x^2) - (\sum x)^2][n(\sum y^2) - (\sum y)^2]}}$$

Date _____

Day _____

- (8) As 3 is smaller than our best distance i.e. 2 we don't visit (2,2)
- (9) Now checking for 9,6
- (10) Ab 6,7 & 7,2 ka 1 draw koege we get
1 (6-7) (7-2) level par hai (hain wo do k minus)
- (11) 0 < best distance ^{kra} we go to 9,6
- (12) (6,7) & (9,6) $d = 3.16$ (don't update current best)
- (13) (6,7) & (8,1) $d = 6$ id

so nearest is 4,7

Search bcm $O(\log n)$

▷ Curse of Dimensionality :

↳ Agr no. of dimensions zyada honge to data points dhooor dhooor honge.

↳ Data points ko gareeb rakhne k liye we will need huge amount of data

Solution:

- (1) Increase dataset
- (2) Reduce dimensions

→ KNN carries out prediction about the test point assuming we have data points near to the test point that are similar to test point.

↳ As we don't have neighbours in high dimensionality space, KNN bcms vulnerable.

Date _____

Day _____

→ All points tend to be almost equally far apart, making the concept of "nearest" unreliable.

▷ Parametric Algorithms :

↳ Assumes a fixed no. of parameters

↳ No matter how much data you throw it won't change its mind about how many parameters it needs.

↳ Like in linear regression

$$b_0 + b_1 \cdot x_1 + b_2 \cdot x_2 = 0$$

coefficient → ek fixed functional form hai.

↳ The values for these are learned thru data

▷ Non Parametric Algorithms :

↳ No fixed parameters, model structure

Date _____

Day _____

CONFUSION MATRIX

→ In skewed dataset we can't even tolerate 1% error rate.

→ Therefore to evaluate performance we use Confusion matrix

→ Detailed of 100 patients
 → 1% people have TB
 → 99% people have no TB

		Predicted	
		True	False
Actual	True	TP	FN -ve
	False	FP	TN

→ Type I error

Type II error

all diagonal values are TP

TP = True +ve = Correctly identified +ve

TN = True -ve = " " " -ve

FN = False -ve = incorrectly identified +ve (Actual True hai lekin humne False predict krna hai)

FP = False +ve = " " " -ve (Actual False hai lekin humne True predict krna hai)

Precision Vs Recall

→ Precision tells us how many of the correctly predicted cases actually turned out to be +ve.

→ This would determine whether our model is reliable or not.

→ Tells us accuracy of +ves.

Page No. _____

in Email Spam detection.

→ Out of all emails you marked as SPAM, how many were actually SPAM.

→ If you marked 10 as Spam & 3 are actually SPAM then 30% precision

Date _____

Day _____

For this to be max i.e. = 1
then FP must be = 0

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{\text{Correctly identified +ves}}{\text{Total +ve Predicted}}$$

→ Recall tells us how many of the actual +ve cases we were able to predict correctly

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{\text{Correctly identified +ves}}{\text{Total Actual +ves}}$$

→ Out of all actual spam emails, how many did you correctly predict?

if 20 spam emails were there & we correctly deleted 8 then recall is 40%.

$$F1\text{-Score} = \frac{2}{\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}} \rightarrow \text{it is max when precision = recall}$$

Micro average =

Date _____

Day _____

Doctor testing for a disease

Positive \rightarrow has ~~dis~~ TB

Neg \rightarrow Donot have TB

Sick person diagnosed
as sick
Healthy person
diagnosed as Sick

	Has TB	Donot have TB
predicted TB (True)	TP	FP
donot predicted TB (False)	FN	TN

A sick person diagnosed
healthy

A ~~sick person~~ healthy person
diagnosed as healthy

Recall is usefull metric in cases where

False Neg is higher concern than PP.

Is important in medical cases. Q k agr ek Sick ko healthy
diagnosed kr dia to marla hojayege.

Imbalance data :

Where one class has significantly more samples than
the other in a classification problem.

Is in a fraud detection dataset \rightarrow 95% train are normal
 \rightarrow 5% are fraudulent

Date _____

Day _____

		Actual		
Predicted		Dog	Cat	Rabbit
	Dog	23 _{TP}	12 _{FP}	7 _{FP}
	Cat	11	29 _{FP}	13
	Rabbit	4	10	24 _{TP}

		Actual	
Predicted		Dog	Cat & Rabbit
	Dog	TP 23	FP 12+7 19
	Cat & Rabbit	FN 11+4 15	TN 29+13+10+24 76

		Actual	
		Cat	Dog & Rabbit
	Cat	TP 29	FP 11+7 18
	Dog & Rabbit	FN 12+10 22	TN 23+13+4+24 58

		Actual	
Predicted		Rabbit	Dog & Cat
	Rabbit	TP 24	FP 4+11 15
	Dog & Cat	FN 7+13 20	TN 23+12+11+29 75

	TP	TN	FP	FN
Dog	23	76	19	15
Cat	29	58	24	22
Rabbit	24	75	14	20

Macro is along along
micro is the south

50 4 2
5 45 5
3 6 40

Date _____

Day _____

Predicted		Actual		
		A	B	C
Predicted	A	50	4	2
	B	5	45	5
	C	3	6	40

FP → row

FN → col

Predicted		Actual	
		A	B & C
Predicted	A	TP 50	FP 4+2=6
	B & C	FN 5+3=8	TN 45+5+6+40=96

Predicted		Actual	
		B	A & C
Predicted	B	TP 45	FP 5+5=10
	A & C	FN 4+6=10	TN 50+2+3+40=95

Predicted		Actual	
		C	A & B
Predicted	C	TP 40	FP 3+6=9
	A & B	FN 2+5=7	TN 50+4+5+45=104

$$\text{Accuracy} = \frac{\text{TP}_A + \text{TP}_B + \text{TP}_C}{\text{Total Samples}} =$$

Formulas for Distance
 $|x_2 - x_1| + |y_2 - y_1| \rightarrow \text{manhattan}$
 $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \rightarrow \text{euclidean}$

$$1 < K < N$$

Date: _____

K-Nearest Neighbour (KNN)

Can be used for classification as well as regression

	x_1	x_2	class	Distance
Training data	1	5	0	2
	0	8	0	6
	0	6	1	4
	1	2	1	1

we need to find class of this test point



Test: (1, 3)

distance formula se tamam training data se distance nikal liya

if $K=1$ means we need to find one nearest neighbour.

we choose odd values (if K is even then we choose odd values because if K is even then we have two classes and we cannot choose one class label as we have two classes. So we choose odd values.)

"K" ~~can~~ \neq no. of data points in training data.

Height	Weight	Label	Distance	K=3 label for Test is 170
170	70	Normal	16	
160	60	Normal	36	
180	80	Over	4	
175	75	over	6	
182	82	overweight	8	
169	66	normal	21	
Test (170, 70)				

→ Machine Learning of ~~the following~~ ~~process~~.
Date: _____

If $K=1$, then of K a model overfit hojayege
we sirf ek hi point ko learn karege as
 $K=1$. Training accuracy to achi hogi (ek
testing accuracy achi nhi aayega.
↳ overlearn

If $K=N$, then models is underfit.
Pattern nhi learn kr payega.

Testing & Training both accuracy is low.
↳ no learning

→ if Binary classification then choose odd
value

→ if Multiple or the even -

- Document A (Sports): "Team wins championship match."
- Document B (Politics): "Election results were announced yesterday."
- New Doc (?): "The match results were announced"
↳ we need to identify label for this