

1. KNN Algorithm

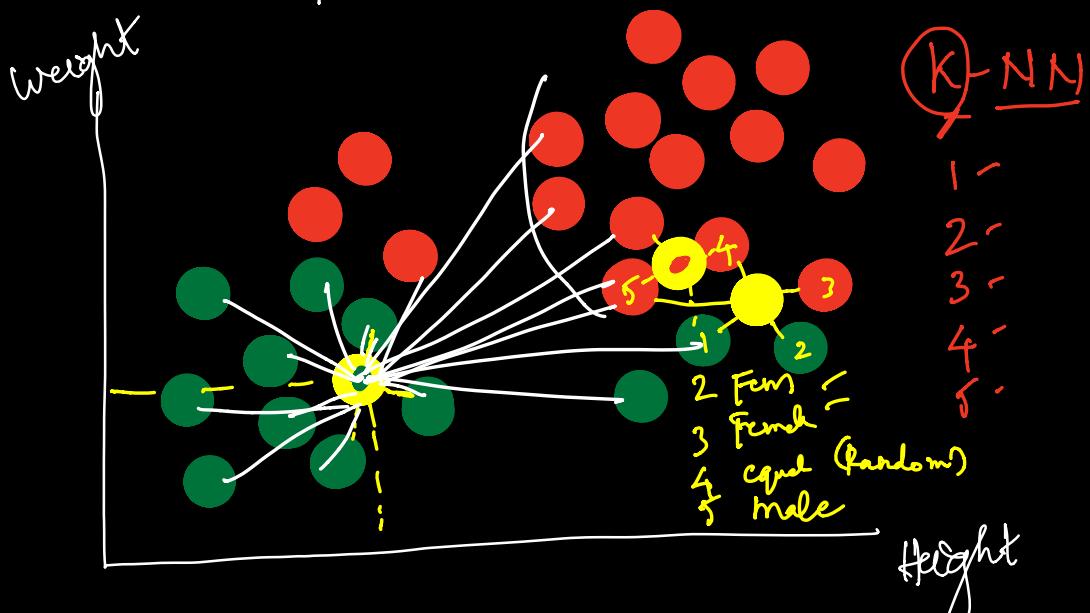
2. Naive Bayes Algorithm

1. KNN

Supervised learning

Regression, classification

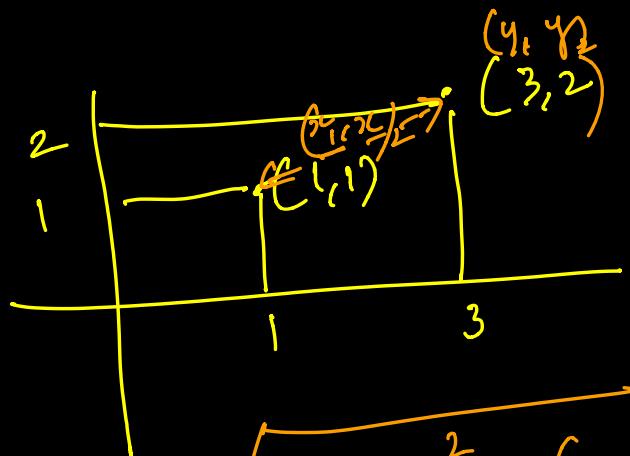
Simple, distance based Algorithms



Distance

Euclidean Distance ✓

$$\sqrt{\sum (x_i - y_i)^2} \quad \checkmark$$



$$\sqrt{(3-1)^2 + (2-1)^2} = \sqrt{2^2 + 1^2} = \sqrt{5} = 2\cdot\sqrt{5}$$

Manhattan

$$\sum |x_i - y_i|$$

Minkowski

$$\sum_{i=1}^k (|x_i - y_i|^q)^{\frac{1}{q}}$$

$$q=2$$

$$\mathbb{E}((\epsilon_i - y_i)^2)$$

Ecological

$K=3$

$$\mathbb{E} \sqrt{(\epsilon_i - y_i)^2}$$

Regression

Customer	Age	Income	No of Credit	Class
② George	35	35K	3	No
① Rachel	22	50K	2	Yes
④ Stere	63	200K	1	No
③ Tom	59	170K	1	Yes
Johm	37	50K	2	No?

$$\sqrt{(37-59)^2 + (50K-170K)^2 + (2-1)^2}$$

122

k is too small \rightarrow sensitive to the
(graph) noise points

k is too large \rightarrow this might be
more generic
(independent)

$k < \log\sqrt{k}(n) \rightarrow n \Rightarrow$ number of
sample

Age	Loan	$k=3$	Default	Distance
25	40,000	N.	N.	102,000
35	60,000	N.	③	82,000
45	80,000	N.	N.	62,000
20	20,000	N.	①	1,221,000
35	120,000	Y	N.	22,000
52	181,000	Y	②	1,24,000
23	95,000	0.62	?	47,000
48	142,000	~	④	~
25	470,000	~	~	~

θ^{-1} Min-Max Scaling (Normalization)

$$\begin{array}{ccc} \text{Max} & \rightarrow & 1 \\ \text{min} & \rightarrow & 0 \end{array} \quad \begin{array}{l} \text{Max} \rightarrow 142000 \\ \text{Min} \rightarrow 18000 \end{array}$$

$$X = \frac{X - \text{min}}{\text{max} - \text{min}} \quad \begin{array}{c} 142000 \\ 18000 \\ \hline 124000 \end{array}$$

$$\frac{142000 - 18000}{124000}$$

$$\frac{124000}{124000} = 1$$

$$\Rightarrow \frac{95000 - 18000}{124000} \quad \begin{array}{c} 95 \\ 18 \\ \hline 77 \end{array}$$

$$= \frac{77000}{124000} = 0.620$$

Max = 52

Min = 20

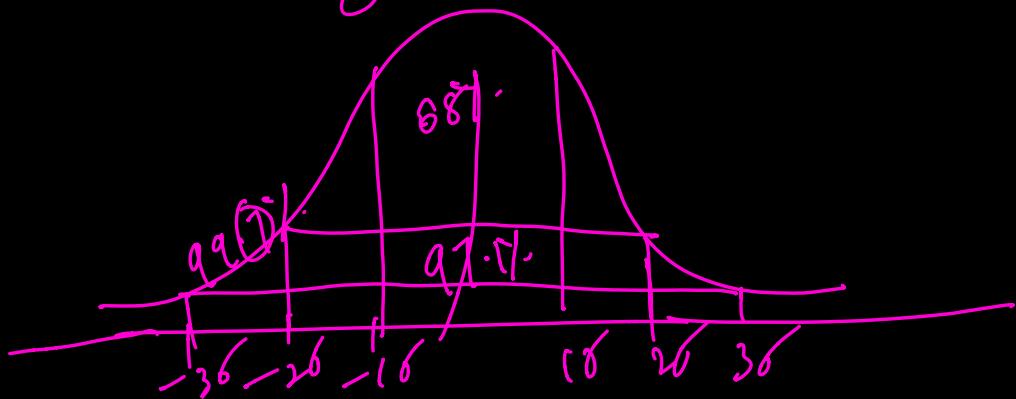
$$\frac{48 - 20}{52 - 20} = \frac{28}{32} = 0.875$$

Age	Loan	Distance	Default
0.1562	0.177	0.71	N
0.466	0.33	③ 0.35	N /
0.78	0.5	① 0.09	N ✓
0	0.01	0.87	N
0.46	0.82	0.40	N
1	0	② 0.12	Y ✓
0.093	0.62	0.78	Y
<hr/>			
0.87	1 ↘	?	N

Standardization

$$Z = \frac{X - \mu}{\sigma}$$

Z-score

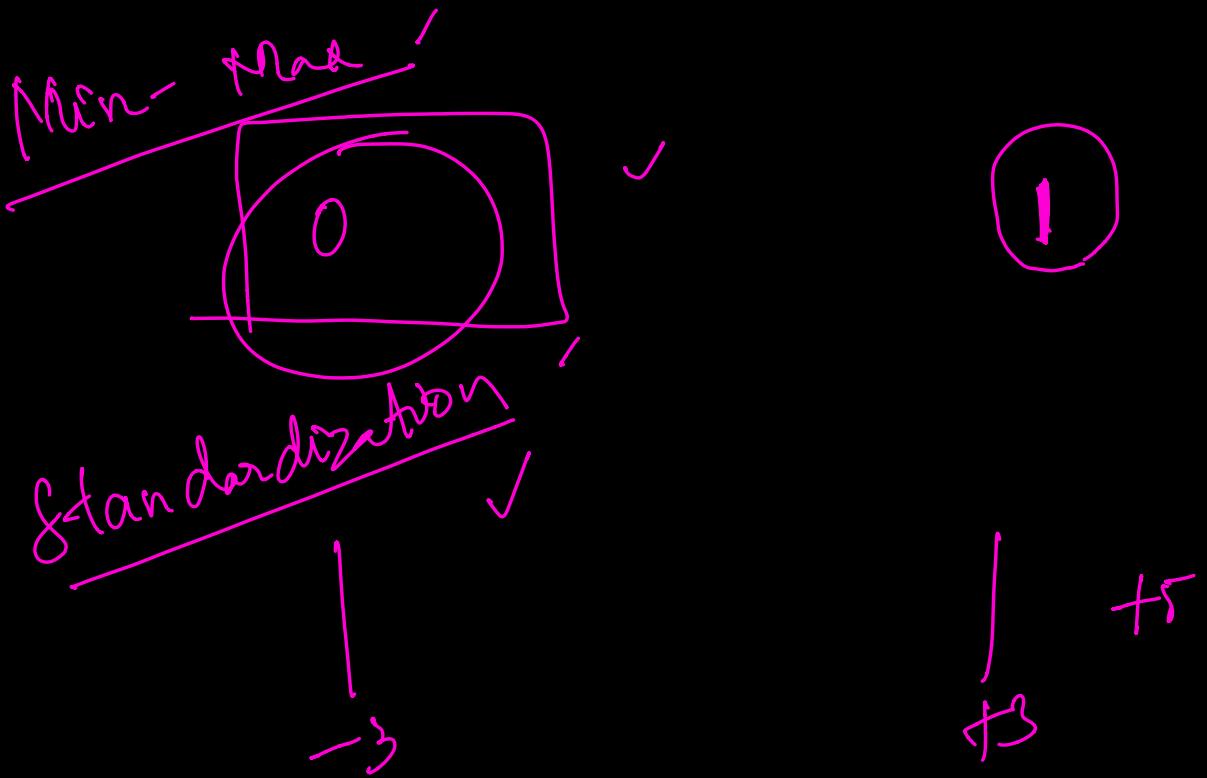


99.7% \rightarrow $-3 \quad +3$

$-5 \quad +5$

Box plot
Z-score

-5 $\quad +5$

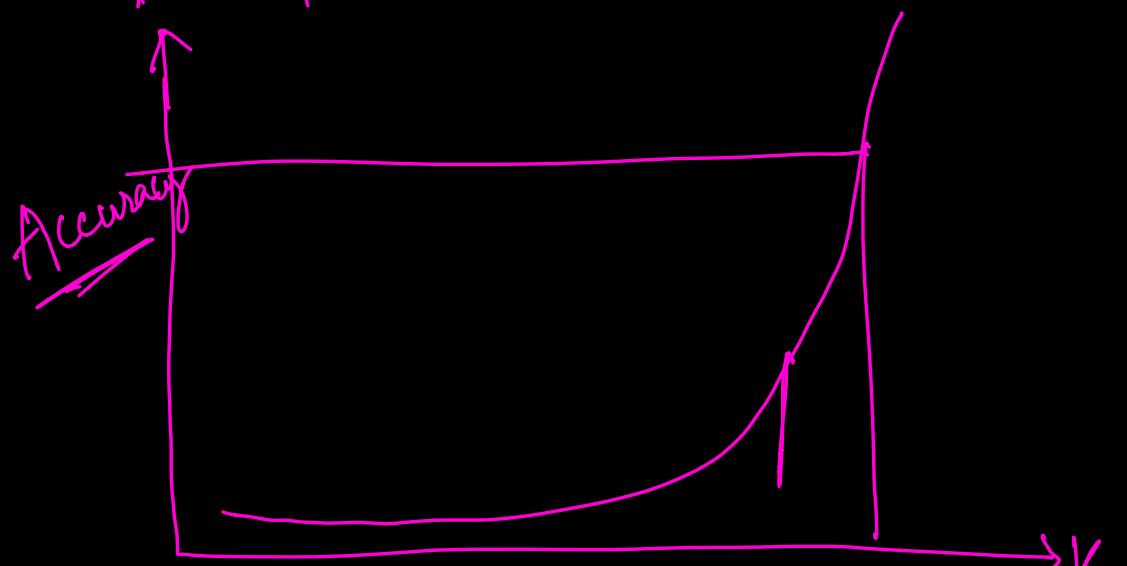
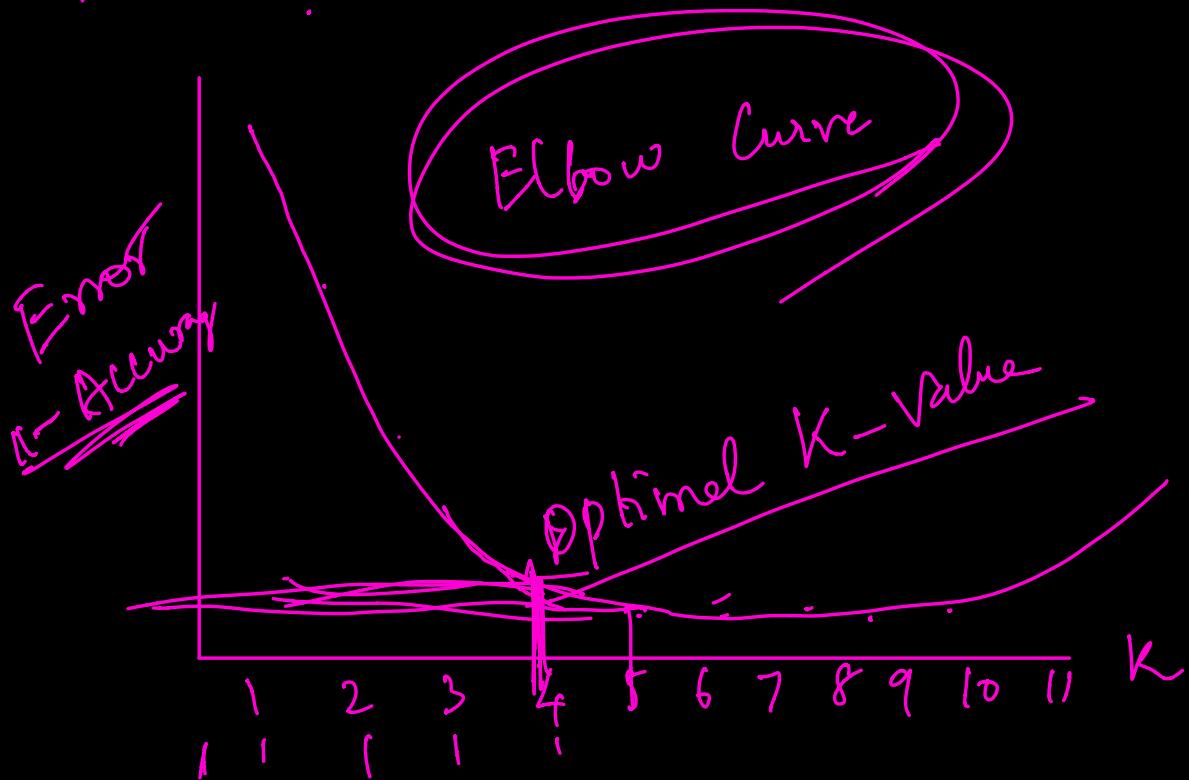


Linear Regression
model $\Rightarrow mx + b \rightarrow$ predict

KNN

model \Rightarrow No model
Lazy model

$K = ?$



Disadvantages:
more time to classify ✓

Need to calculate lot of distance

Choosing right k is tricky

Need large Number of samples

KD Tree can be used to reduce
the complexity of time

Naive Bayes Algorithm

Supervised Algorithm

Classification

Naive Bayes classifier works on the principle of conditional probability as given by Bayes' theorem

Toss two coins :

$$SS = \{ \overbrace{HH, HT, \cancel{TT}, TH}^{\text{Two Heads}} \}$$
$$P(\text{Getting Two Heads}) = \frac{1}{4}$$



$$P(\text{Second coin being head given first coin is tail})$$

= $\left\{ TT, TH \right\}$

1
2

$$P(\text{Getting two heads given first coin is head})$$

= $\left\{ HH, HT \right\}$

$\frac{1}{2}$

$$P(A|B) = \frac{P(A) P(B|A)}{P(B)}$$

$P(\text{Second coin being head given first coin is tail})$

$$P(A/B) \quad \{ HH, HT, TH, TT \}$$

$A = \text{Second coin being head}$

$B = \text{first coin is tail}$

$$P(A/B) = \frac{P(A) \cdot P(B/A)}{P(B)}$$

$$P(A) = \frac{2}{4}$$

$$P(B) = \frac{2}{4} \quad \{ HH, HT, TH, TT \}$$

$$P(B/A) = \frac{\text{first coin is tail given second coin being head}}{2}$$

$$\left\{ \begin{matrix} \text{H}\bar{\text{H}}, \bar{\text{T}}\text{H} \end{matrix} \right\}^{\vee} \\ = \mathbb{V}_2$$

$$\frac{\frac{2}{4} * \frac{1}{2}}{\frac{2}{4}} = \frac{\frac{1}{2} * \frac{1}{2}}{\frac{1}{2}}$$

$$\frac{\frac{1}{2} * \frac{2}{1}}{\frac{1}{2}}$$


$$P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)}$$

Class Prior Probability

Likelihood

↓

Posterior Probability

↓

Predictor Prior Probability

$A =$ Target
 $B =$ Independant Data

L.

Day = Holiday
 Discount = Yes
 Free Delivery = Yes

$$P(\text{Buy} = \text{Yes})$$

$$P(\text{Buy} = \text{No})$$

$$P(\text{Buy} = \text{Yes})$$

Day = Holiday
 Discount = Yes
 Free Delivery = Yes

$$\begin{aligned}
 P\left[\text{Buy} = \text{Yes}\right] &= P(\text{Day} = \text{Holiday} \mid \text{Buy} = \text{Yes}) \\
 &+ P(\text{Discount} = \text{Yes} \mid \text{Buy} = \text{Yes}) \\
 &+ P(\text{Free Delivery} = \text{Yes} \mid \text{Buy} = \text{Yes})
 \end{aligned}$$

$$\frac{P(\text{Day} = \text{Holiday}) \times P(\text{Discount} = \text{Yes}) \times P(\text{Free Delivery} = \text{Yes})}{P(\text{Buy} = \text{Yes})}$$

Likelihood Tables

Likelihood Table		Buy		11/30
		Yes	No	
Day	Weekday	9/24	2/6	11/30
	Weekend	7/24	1/6	
	Holiday	8/24	3/6	11/30
		24/30	6/30	

Frequency Table		Buy		20/30
		Yes	No	
Discount	Yes	19/24	1/6	20/30
	No	5/24	5/6	
		24/30	6/30	

Frequency Table		Buy		23/30
		Yes	No	
Free Delivery	Yes	21/24	2/6	23/30
	No	3/24	4/6	
		24/30	6/30	

Calculating Conditional Probability of purchase on the following combination of day, discount and free delivery:

Where B equals:

- Day = Holiday
- Discount = Yes
- Free Delivery = Yes

Let A = Buy $P = \text{No Buy}$

$$P(A|B) = P(\text{Yes Buy} | \text{Discount} = \text{Yes}, \text{Free Delivery} = \text{Yes}, \text{Day} = \text{Holiday})$$

$$= \frac{P(\text{Discount} = \text{Yes} | \text{Yes}) * P(\text{Free Delivery} = \text{Yes} | \text{Yes}) * P(\text{Day} = \text{Holiday} | \text{Yes})}{P(\text{Discount} = \text{Yes}) * P(\text{Free Delivery} = \text{Yes}) * P(\text{Day} = \text{Holiday})}$$

$$= \frac{(19/24) * (21/24) * (8/24) * (24/30)}{(20/30) * (23/30) * (11/30)}$$

$$= 0.986$$

$$\frac{\frac{24}{30} * \frac{8}{24} * \frac{19}{24} * \frac{21}{24}}{\frac{11}{30} * \frac{20}{30} * \frac{23}{30}}$$

$$\Rightarrow 0.986 \approx$$

$$P(\text{NoBuy} \mid \text{---})$$

$$= 0.178$$

0.

$$0.629$$

Play = yes
~~Play~~

$$0.3107$$

$$0.628$$

$$0.62156$$

$$P(\text{No}) = 0.212$$