

Last class (September 1)

- 1) Blinkit Problem Statement
↳ issue with logistic regression
- 2) Geometric Intuition of kNN
- 3) kNN Algorithm
- 4) kNN scratch code
- 5) Assumptions of kNN
- 6) sklearn's kNN implementation
- 7) Bias-Variance trade-off
- 8) Train & Test time complexity of kNN

Today's class

- 0) Recap- Quizzes
- 1) kNN with categorical features
- 2) Different distance metrics
- 3) LSH for kNN
- 4) Missing data with kNN — imputation
- 5) Employee Attrition Problem Statement
- 6) Decision Tree Intuition
- 7) How to split the nodes
- 8) Purity and Impurity of the nodes

Target Encoding / output

x_2	x_1	x_2	y
0.75	0.8	RED	1
0.5	0.1	GREEN	0
0.75	2.2	RED	1
0.75	-1.5	RED	1
0.5	-0.9	GREEN	1
0.75	3.2	RED	0

$P(y=1 | x_2 = \text{GREEN}) = \frac{1}{2}$
 $= \frac{P(y=1 \cap x_2 = \text{Green})}{P(x_2 = \text{Green})} = 0.5$

$P(y=1 | x_2 = \text{RED}) = \frac{3}{4} = 0.75$
 $= \frac{P(y=1 \cap x_2 = \text{RED})}{P(x_2 = \text{RED})}$

OHE

3 values \rightarrow RED, GREEN, YELLOW



$3-1 = 2 \rightarrow$ sufficient

\rightarrow Yellow

	is_Green	is_Red	is_Yellow
Not Red	0	0	1
Not Green	0	1	0
Red	1	0	0

N unique cat values \Rightarrow OHE needs

num feat

$\rightarrow n-f \leq 6 \rightarrow$ OHE

$N-1$ dimensions

$> 6 \rightarrow$ Target Encoding

Categorical values: A, B, C
 3 classes: $0, 1, 2$

$$P(y=0 | x=A/B/C) = \begin{matrix} 0.5 & 0.25 & 0.25 \\ A & B & C \end{matrix}$$

$$P(y=1 | x=A/B/C) = \begin{matrix} 0.3 & 0.6 & 0.1 \\ A & B & C \end{matrix} \quad \checkmark$$

$$P(y=2 | x=A/B/C) = \begin{matrix} 0.4 & 0.45 & 0.15 \\ A & B & C \end{matrix}$$

$$V_1 = [-5, 5]$$

$$V_2 = [-3, 3]$$

Similarity $\propto \frac{1}{\text{Distance}}$

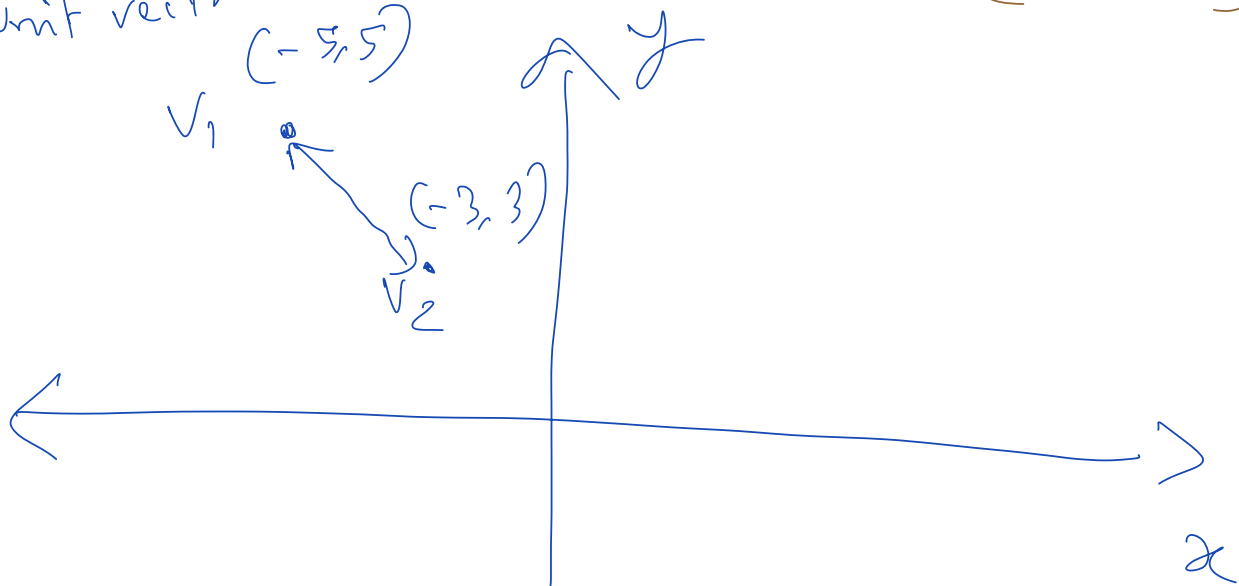
$$\cos(V_1, V_2) =$$

$$L2\text{-norm of } V_1 = \sqrt{(-5)^2 + 5^2} = \sqrt{50}$$

$$\begin{aligned} \downarrow \\ \text{unit vector } V_1' &= \left[-\frac{5}{\sqrt{50}}, \frac{5}{\sqrt{50}} \right] \quad \checkmark \\ L2\text{-norm of } V_2 &= \sqrt{(-3)^2 + 3^2} = \sqrt{18} = 3\sqrt{2} \end{aligned}$$

$$v_2' = \left[-\frac{3}{3\sqrt{2}}, \frac{3}{3\sqrt{2}} \right] = \left[-\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right]$$

unit vector



$$\cos(v_1, v_2) = v_1' \cdot v_2'$$

$$= \left(-\frac{1}{\sqrt{2}} \right) \left(-\frac{1}{\sqrt{2}} \right) + \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{2} + \frac{1}{4} = \frac{3}{4}$$

Angle between v_1 & $v_2 \rightarrow$ small,

\Rightarrow cosine similarity $\rightarrow 1$
very high

$\cos(\theta) = -1$ to $+1$

cosine similarity $\propto \frac{1}{\text{Euclidean Distance}}$

cosine distance = $1 - \text{cosine similarity}$

Test time for kNN: $O(nd + n \log n)$

assuming k is very small

$k \ll n$

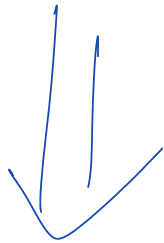
$1 \text{ B} \Rightarrow 10^9$
 $100 \rightarrow 10^2$

$\downarrow 10^7 \text{ times}$

	f_1	f_2	f_3	f_4	f_5	z
x_1	0.5	1.1	-1.1	10	7	
x_2	10.0	18	NAN	15	14	
x_3	13.0	7	8	9	13	
			\vdots			
x_{100}						

$\rightarrow 1.5$

Drop f_3



$k \text{ NN}$
 \downarrow

get $k \text{ NN}$ of x_2

	f_1	f_2	f_4	f_5	z
x_1	0.5	1.1	10	7	
x_2	10.0	18	15	14	
x_3	13.0	7	9	13	
	\vdots	\vdots	\vdots	\vdots	\vdots
x_{100}					

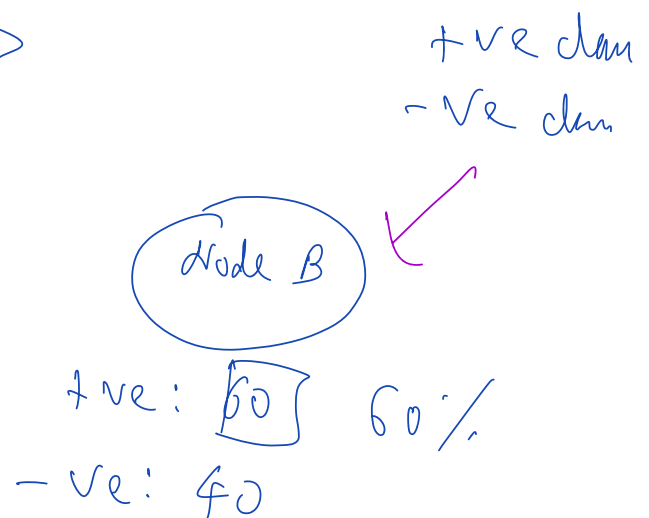
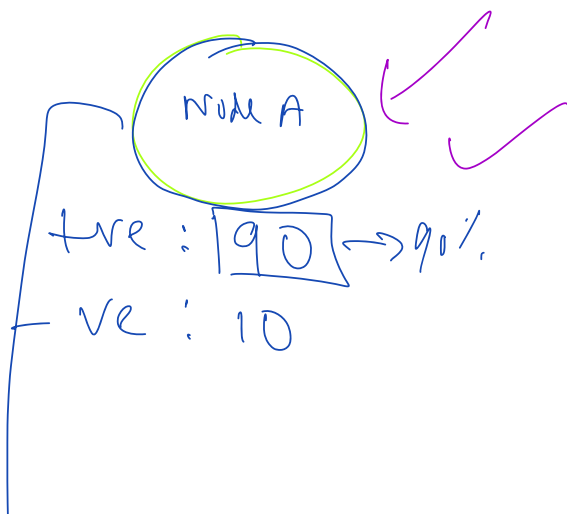
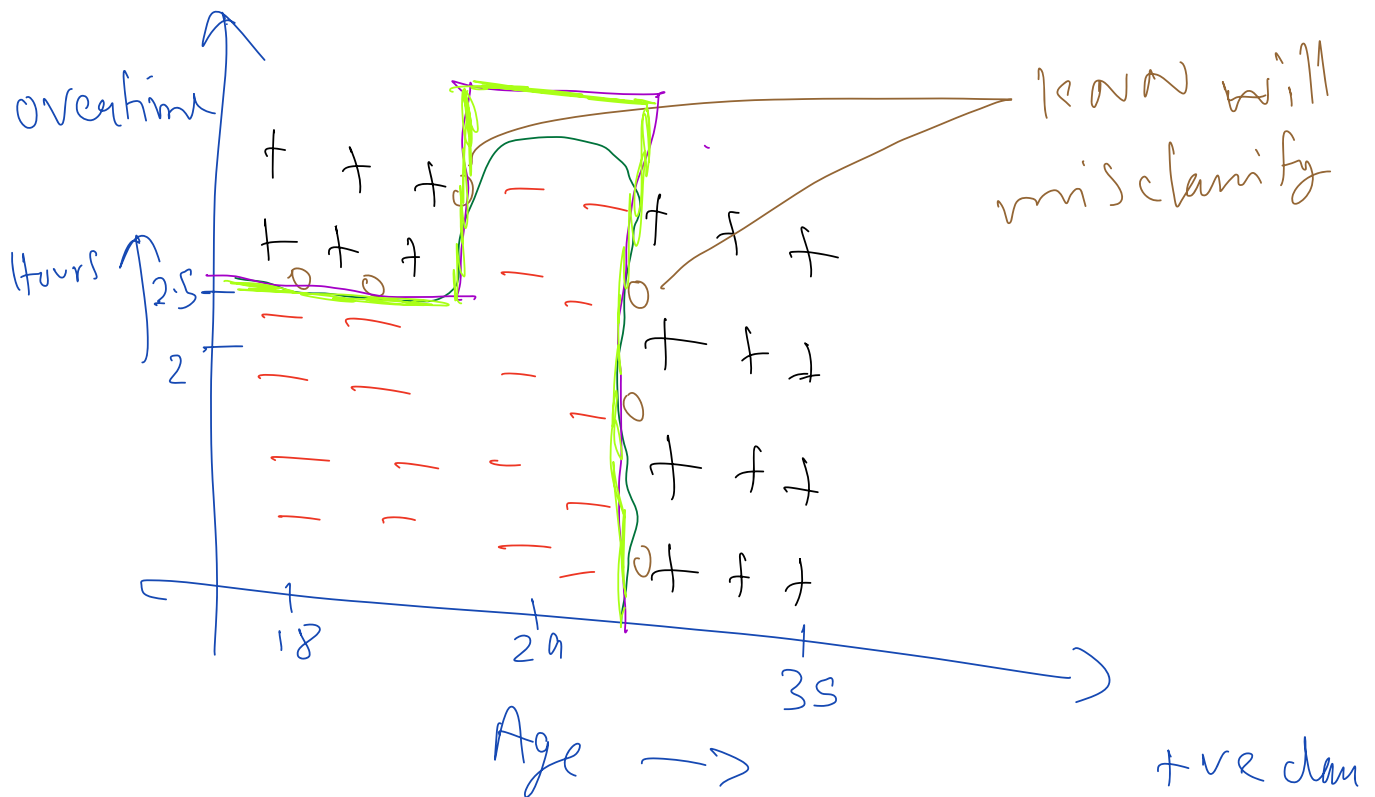
x_1
 x_5
 x_{10}
 x_{100}

\downarrow nearest neighbors

$$x_{23} = 2$$

$$\text{avg}(x_{1,3}, x_{5,3}, x_{10,3}, x_{100,3})$$

$$x_{23} = 1.5$$



↓ more homogenous → more pure

Node A

class 0 : 60%
1 : 15%
2 : 25%

Node B

class 0 : 30%
1 : 40%
2 : 30%

more pure

Node 1

class 0 : 30%
" 1 : 33%
" 2 : 37%

Node 2

class 0 : 40%
1 : 50%
2 : 10%

class 2

class 1

f_1	f_2	$f_3 \dots$	f_{10}
τ_{11}	$\tau_{2,1}$		
τ_{12}	τ_{22}		
\vdots	\vdots		
$\tau_{1,10}$	$\tau_{2,10}$		

100 + ve

100 - ve

100 0

...

200 \rightarrow x_i 's

$f_1 > \tau_{11} <$

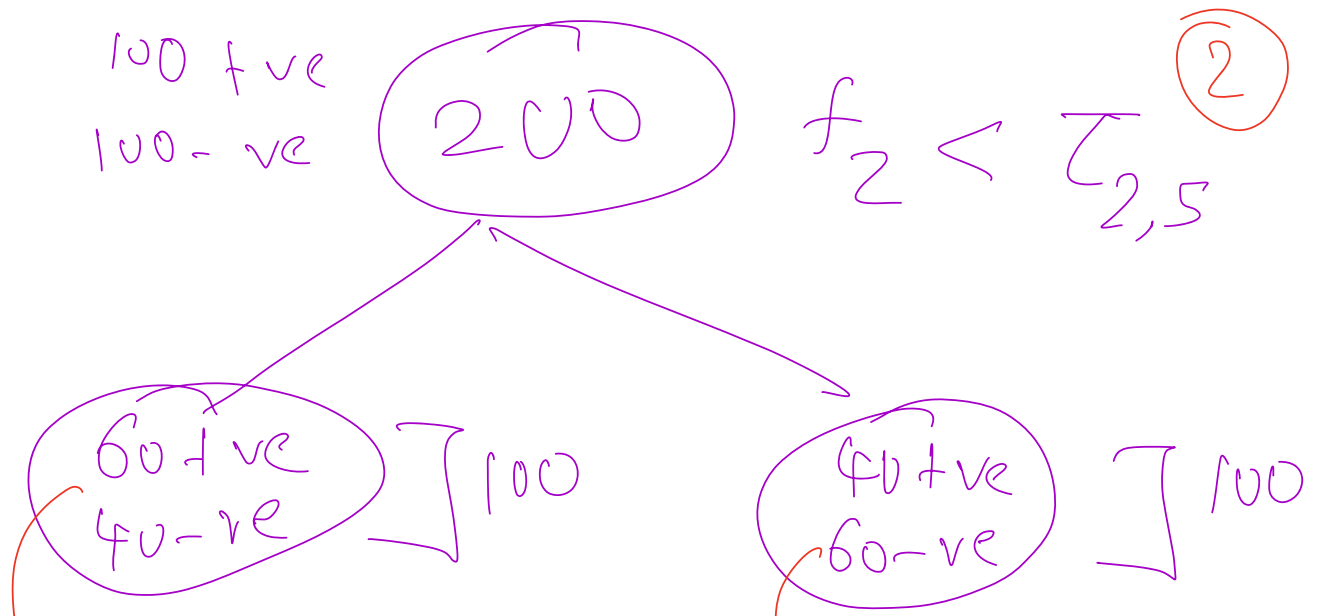
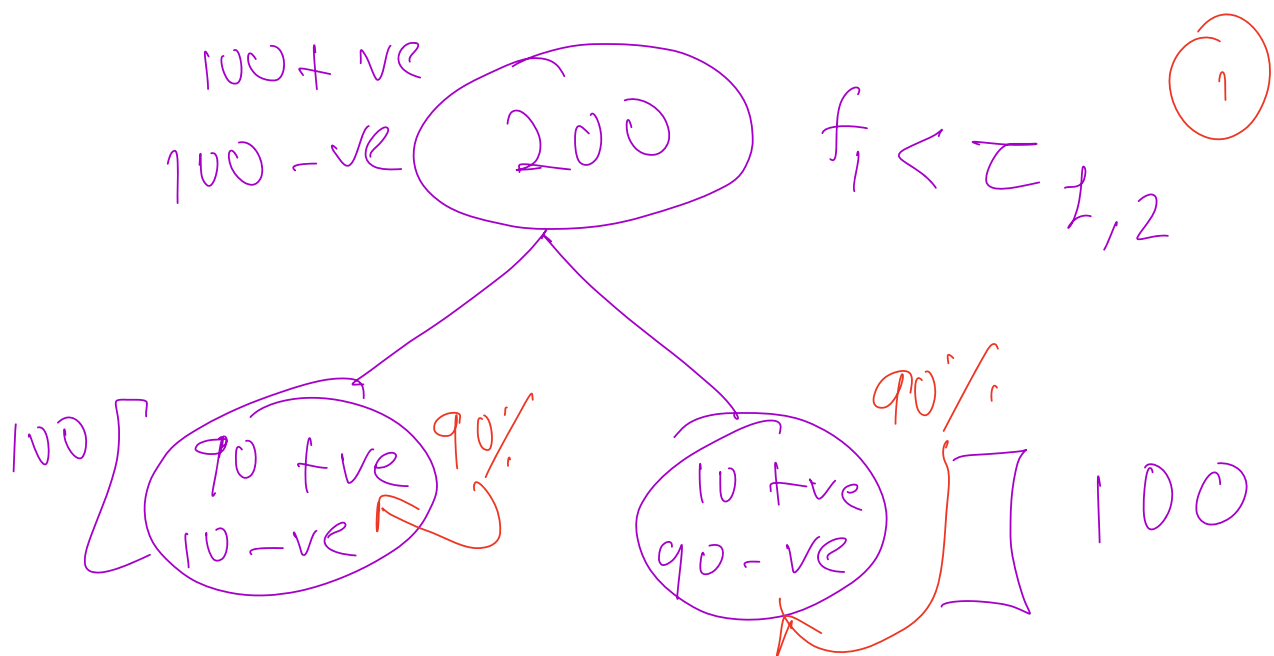
$f_1 > \tau_{1,10} <$

\vdots

possibilities

$$f_2 > \tau_{2,1} \leftarrow$$

$$f_3 < \tau_{3,10} \leftarrow$$



$\hookrightarrow 60\% +ve$

$\hookrightarrow 60\% -ve$