

Previous class - 22 June 2023

- 1) Dataset & Problem Statement (Blinkit)
- 2) Visualising data & Motivation for KNN
- 3) KNN algorithm & assumptions
- 4) Code: scratch implementation
- 5) Train and Test-time Complexity for KNN

Today's agenda

- ✓ 1) Review of code: scratch implementation ✓
- ✓ 2) sklearn's KNN implementation ✓
- ✓ 3) Different distance metrics ✓
- ✓ 4) Weighted KNN overview ✓
- ✓ 5) Bias-Variance Trade off in KNN ✓
- ✓ 6) Impact of outliers ✓
- ✓ 7) How to handle categorical features ✓
- ✓ 8) Applications: Google Image Search ✓
- ✓ 9) KNN based imputation ✓
- ✓ 10) Overview of Employee Attrition Dataset
- ✓ 11) AMA Session (11:15 pm)

Different Distance Metrics

- 1) Euclidean $\rightarrow \sqrt{\sum_{j=1}^d (x_{ij} - x_{qj})^2}$
- 2) Manhattan $\rightarrow \left[\sum_{j=1}^d |x_{ij} - x_{qj}| \right]^1$
- 3) Minkowski $\rightarrow \left[\sum_{j=1}^d |x_{ij} - x_{qj}|^p \right]^{1/p}$
- 4) Cosine

$p = 1 \Rightarrow$ Manhattan

$p = 2 \Rightarrow$ Euclidean

$p = 1.5, 2.5, 4, 6, \dots$

Cosine: $(x_{i1}, x_{i2}, \dots, x_{id}) \rightarrow x_i$
 $(x_{q1}, x_{q2}, \dots, x_{qd}) \rightarrow x_q$

$\overline{x_i} = L2 \text{ norm of } x_i = \sqrt{x_{i1}^2 + x_{i2}^2 + \dots + x_{id}^2}$

$\hat{x_i} = \frac{x_i}{\overline{x_i}} = \left(\frac{x_{i1}}{\overline{x_i}}, \frac{x_{i2}}{\overline{x_i}}, \frac{x_{i3}}{\overline{x_i}}, \dots, \frac{x_{id}}{\overline{x_i}} \right)$

\downarrow
unit vector

$\overline{x_q} = L2 \text{ norm of } x_q = \sqrt{x_{q1}^2 + x_{q2}^2 + \dots + x_{qd}^2}$

$\hat{x_q} = \frac{x_q}{\overline{x_q}} = \left(\frac{x_{q1}}{\overline{x_q}}, \frac{x_{q2}}{\overline{x_q}}, \dots, \frac{x_{qd}}{\overline{x_q}} \right)$

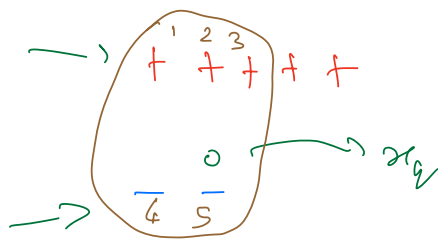
$$\text{cosine}(\hat{x}_q, \hat{x}_i) = \hat{x}_{i1} \cdot \hat{x}_{q1} + \hat{x}_{i2} \cdot \hat{x}_{q2} + \dots + \hat{x}_{id} \cdot \hat{x}_{qd}$$

↪ 0 to 1

As a general practice :

- 1) Lower dimensions → Euclidean Distance
- 2) Higher Dimensions → Cosine Similarity

Weighted kNN



$$W_i = \frac{1}{d_i}$$

K = 5

↓
-ve

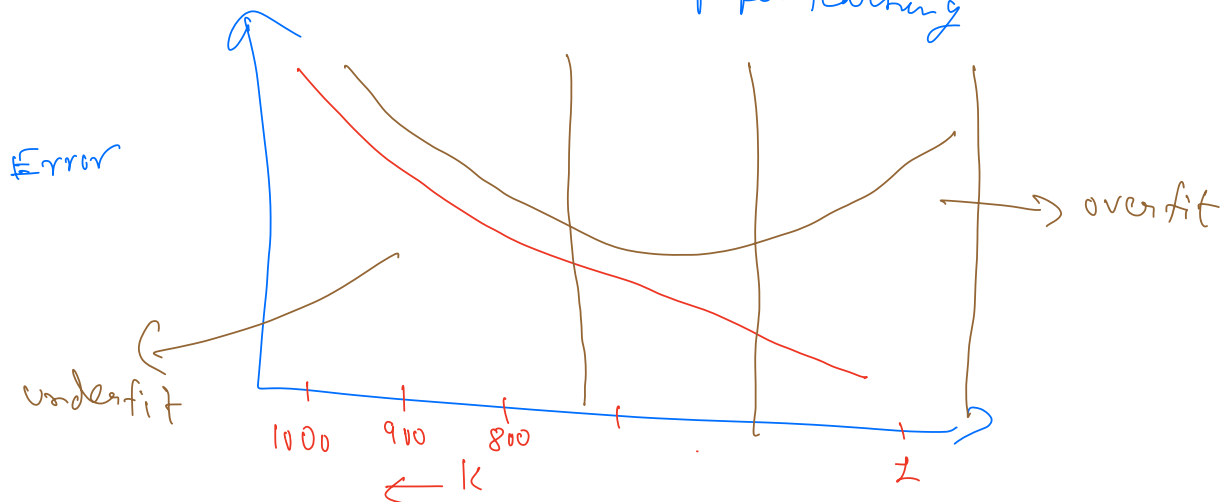
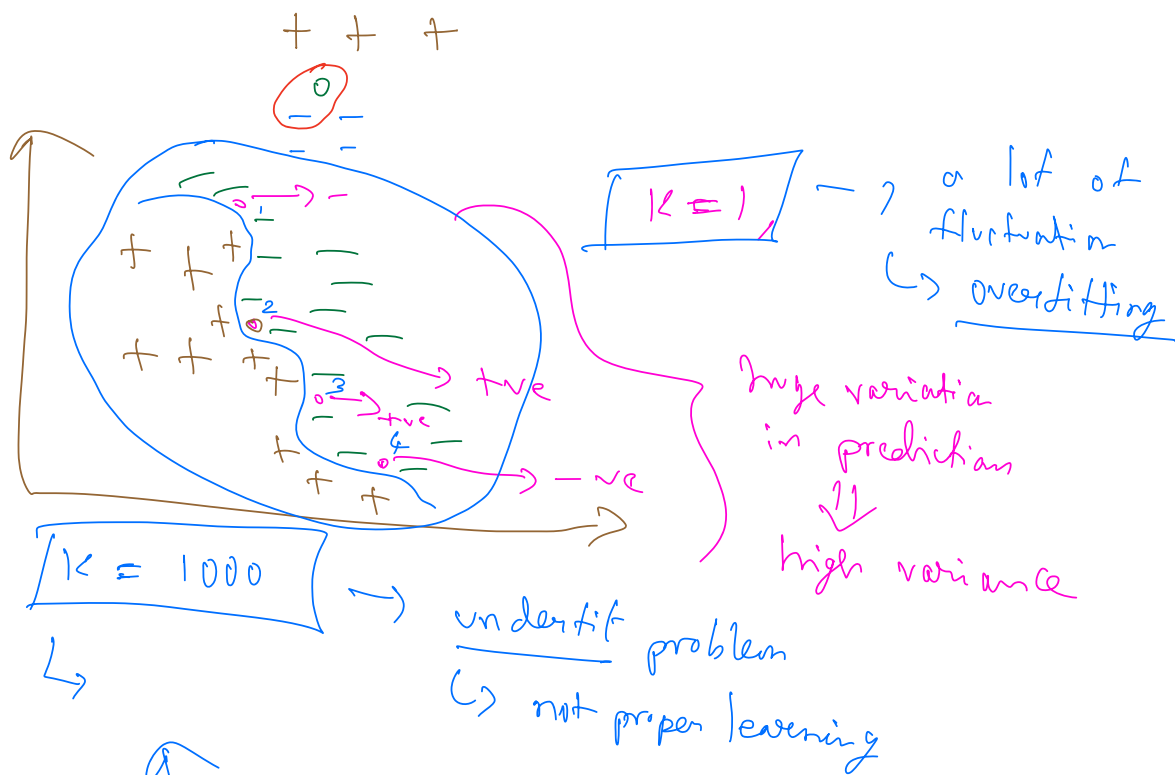
$$d_{1q} = 0.5 = d_{2q} = d_{3q}$$

$$d_{4q} = 0.1 = d_{5q}$$

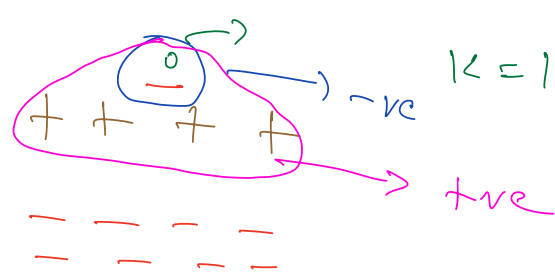
$$\left. \begin{array}{l} \underline{W_{1q}} = \frac{1}{0.5} = 2 = \underline{W_{2q}} = \underline{W_{3q}} \\ \underline{W_{4q}} = \frac{1}{0.1} = 10 = \underline{W_{5q}} \end{array} \right\} \begin{array}{l} +ve \rightarrow 2 + 2 + 2 = \underline{6} \\ -ve \rightarrow 10 + 10 = \underline{20} \end{array}$$

Extend it to more than 2 classes

$20 > 6 \Rightarrow \text{pred} \rightarrow -ve$



$k=5$
 k is small,
 outliers can impact
 predictions



small k values tries to fit
 the noise \Rightarrow one definition of
 overfitting

Break till 10:20 pm

How do we handle categorical features

1) One hot encoding \rightarrow

2) Label encoding \rightarrow

3) Target encoding \rightarrow

OHE \rightarrow 3 classes

Dog \rightarrow 01
Cat \rightarrow 10
Horse \rightarrow 00

\rightarrow classes \rightarrow N-1 dimensions

Cat Dog
0 1 \rightarrow Dog
1 0 \rightarrow Cat
[Horse] 0 0 \rightarrow neither Cat nor Dog

Label encoding

Dog \rightarrow 1
Cat \rightarrow 2
Horse \rightarrow 3

Target encoding

$p(y=1 / c_1 = \text{Dog})$
 $\hookrightarrow p_1$

$p(y=1 / c_1 = \text{Cat})$
 $= p_2$

$p(y=1 / c_1 = \text{Horse})$
 $= p_3$

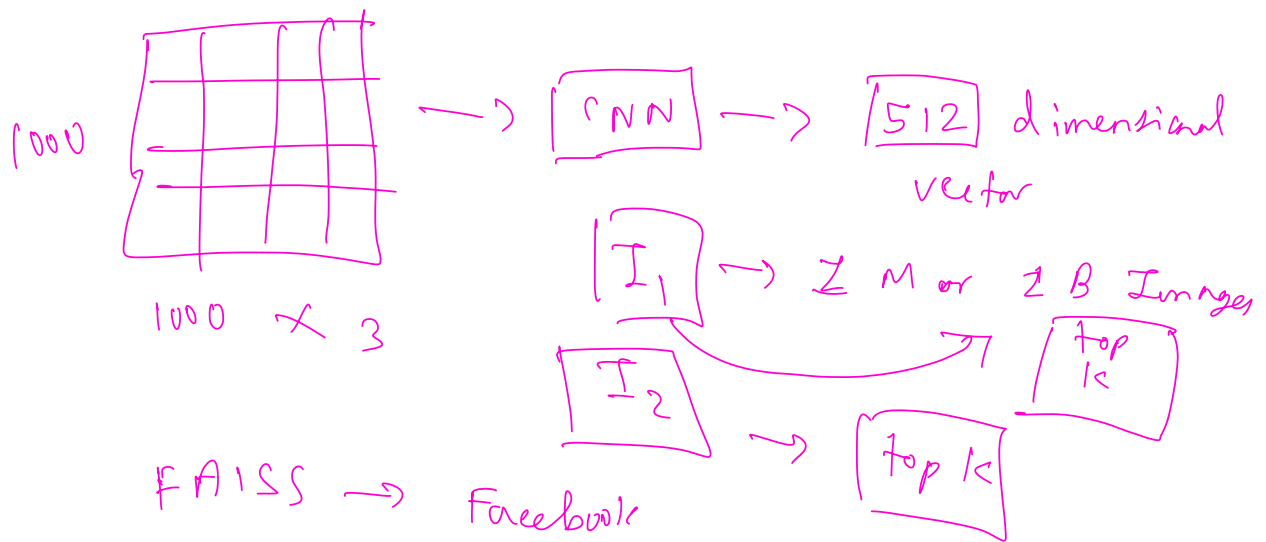
c_1
Dog $\rightarrow p_1$

Cat $\rightarrow p_2$

Horse $\rightarrow p_3$

y-label \rightarrow 0/1

Categorical Embeddings



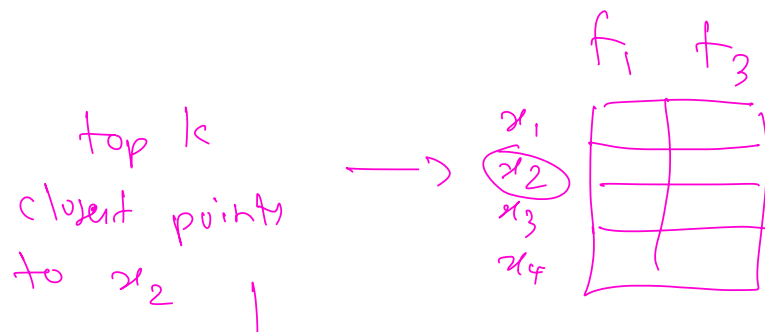
Youtube? Similar Content Search

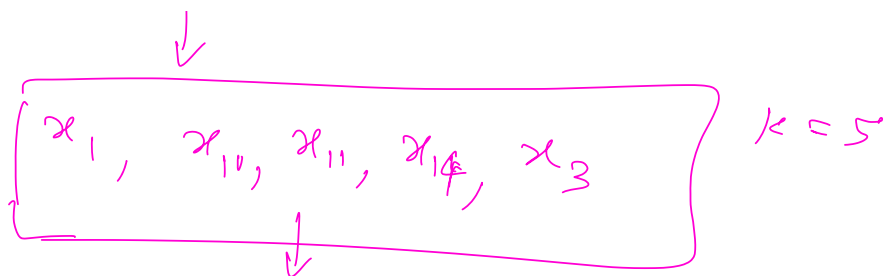
Imputation

	f_1	f_2	f_3	y
x_1	10	17	18	0
x_2	11	11	4	0
x_3	15	1	2	1
x_4	9	23	3	1

\downarrow
Mean, Median, Same Statistic value

Drop $\rightarrow f_2$

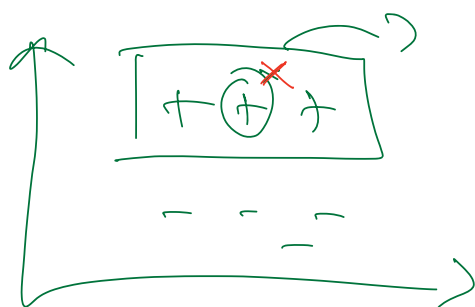
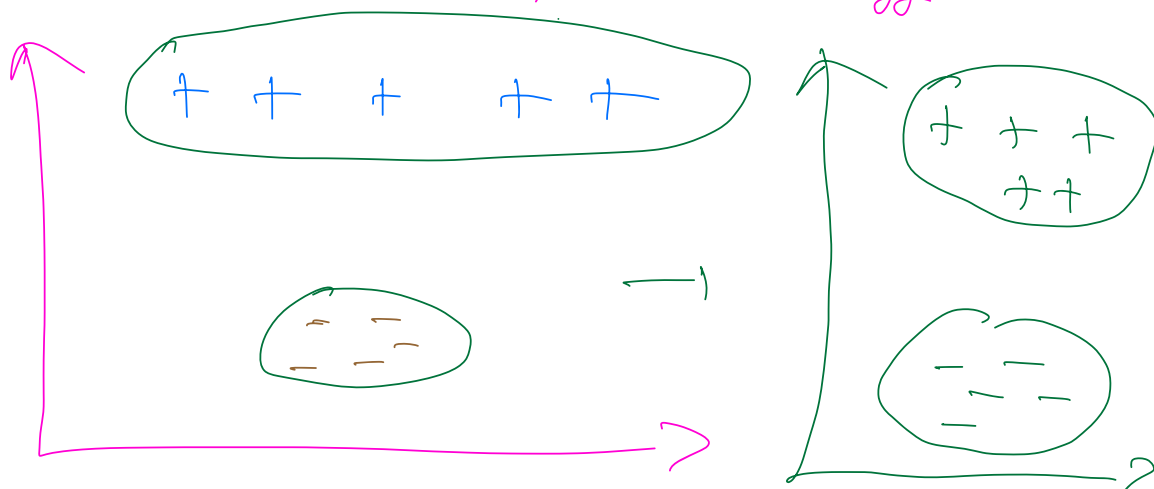




$$x_{2,2} \equiv \text{avg of } (x_{1,2}, x_{10,2}, x_{11,2}, x_{14,2}, x_{3,2})$$

↓
continuous data

For other data, other strategy



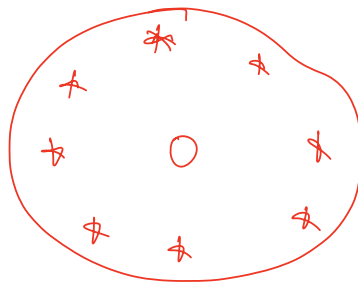
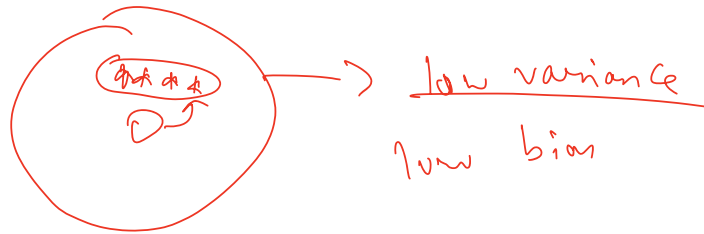
$k=1$ ~~does it train~~
 ↳ input on
 train accuracy

$$- \left[y_i \log(p_i) + (1 - y_i) \log(1 - p_i) \right]$$

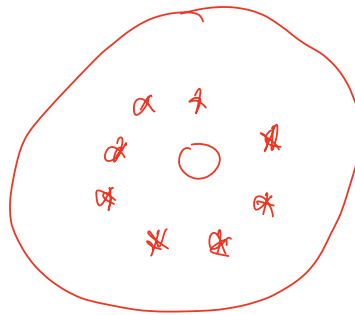
← undefined

$y_i = 0/1$

$p_i = \text{prob of falling into class 1}$



High Bias



Low Bias