

KNN

11 February 2022 20:56

Classification

K-Nearest Neighbours

K-number?

Nearest Neighbour

distance

Euclidean

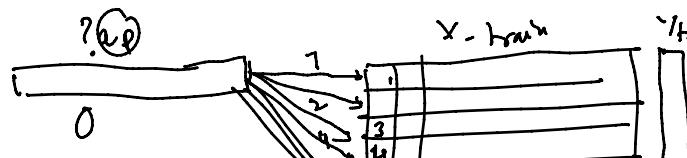
$$d(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + (p_3 - q_3)^2 + \dots} = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

$$(x, y) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

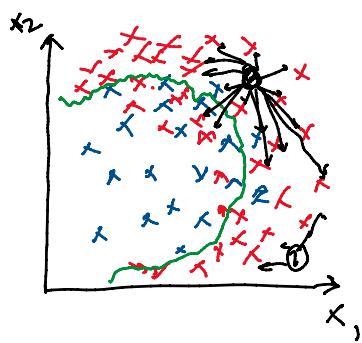
$$\text{Manhattan} = \sum_{i=1}^n |p_i - q_i|$$

$$\text{Minkowski} = \left(\sum_{i=1}^n (|p_i - q_i|)^p \right)^{\frac{1}{p}} \quad p=1$$

$$K = 3\sqrt{r}$$



Binary



(m) dimensional

Neighbours (metric)

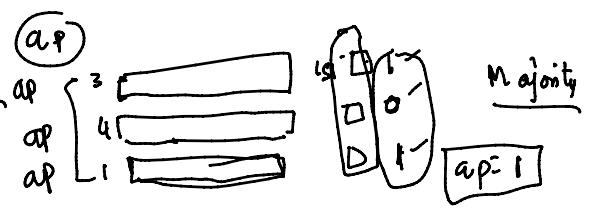
(1) distance (Euclidean, Manhattan)

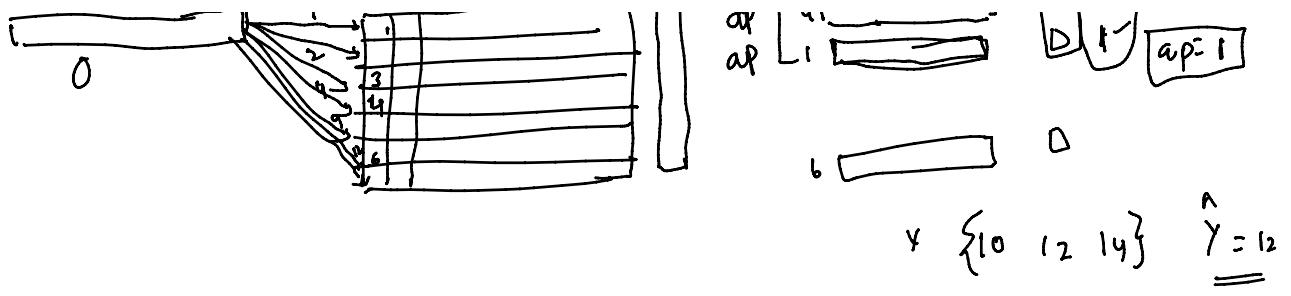
(2) Similarity (Cosine)

$$\sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

$p=1 \rightarrow \text{Manhattan} \checkmark$

$p=2 \rightarrow \text{Euclidean}$





Why?

Training data. $\text{Obs} = \underline{1000}$ $\text{Var} = \underline{10}$ Try L.R. data

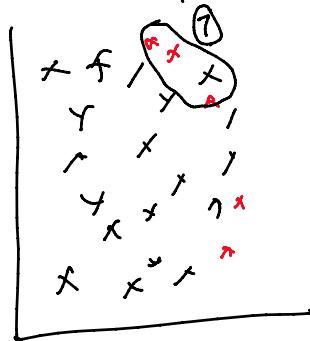
$\boxed{\text{Obs} = \underline{\underline{1000}}}$ $\text{Var} = \underline{\underline{500}}$ X Try L.R. data

"dimensionality curse"

BioTech $P = 500$ $N = 500$ $[1000]$ [Gene
1000 variables]

$n = 1000$, $m = 1000$ \times mL

Imbalanced data



Sal 4 years after

distance based metric :-

* Normalize data ✓

(1)	25000		$\frac{3}{7}$
(2)	30000		7

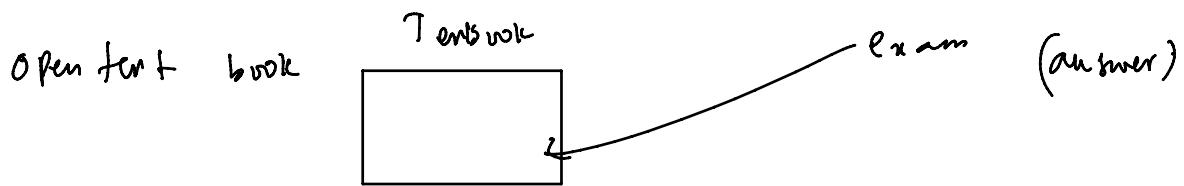
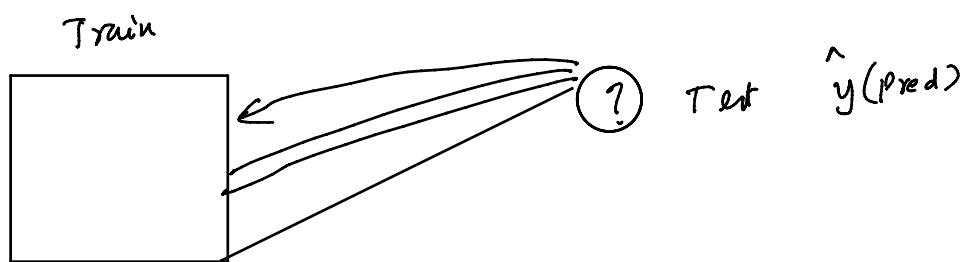
$$d = \sqrt{(30000 - 25000)^2 + (7 - 3)^2}$$

Numeric data (One hot encoding)

- ① Normalize data
- ② Distance (Euclidean, Manhattan)
- ③ K ?



disadvantages :- lazy learning

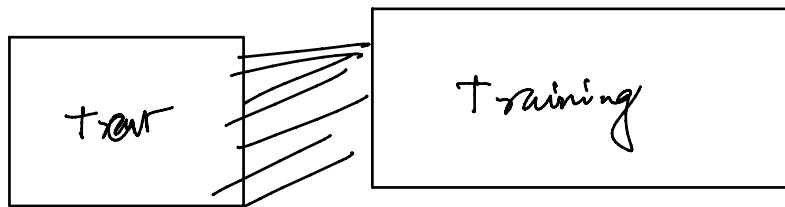


Complexity

Train $O(1)$

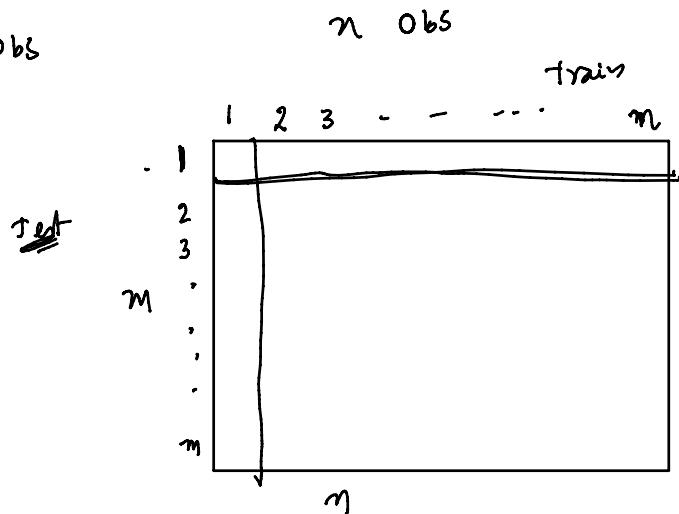
Calculation $O(n)$

Test(m) $\rightarrow O(n \times m)$

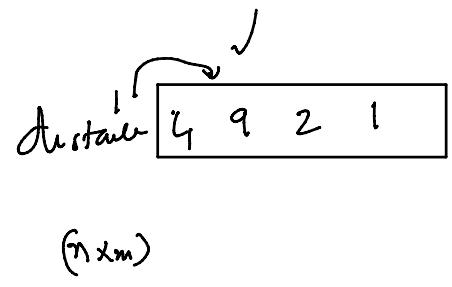


m obs

n obs



(distante)

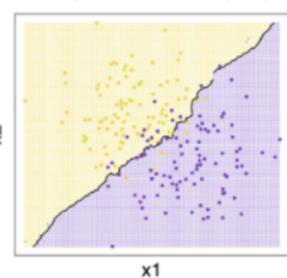
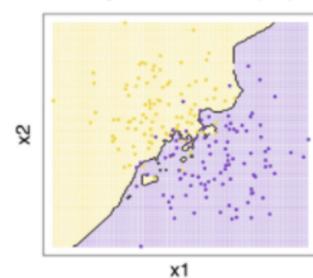
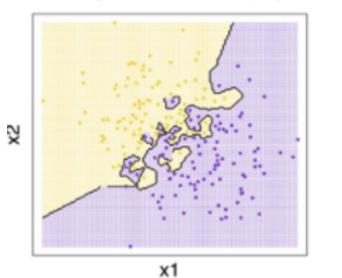


(n x m)

Binary kNN Classification (k=1)

Binary kNN Classification (k=5)

Binary kNN Classification (k=25)

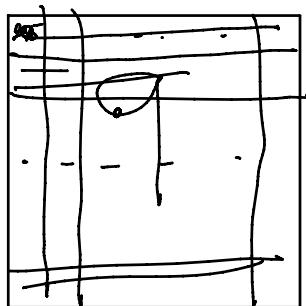


28

9

$$1 \times 1 \cdot 1 \times 2 \cdot 1 \times 3 \cdots \cdots 1 \times 28 \cdot 2 \times 1 \cdot 2 \times 2 \cdots \cdots 2 \times 28 \cdots 28 \times 1 \cdot 28 \times 2 \cdots 28 \times 28$$

28



784

Multi class Problem

$$V \sim S_n - a?$$

Multi class problem

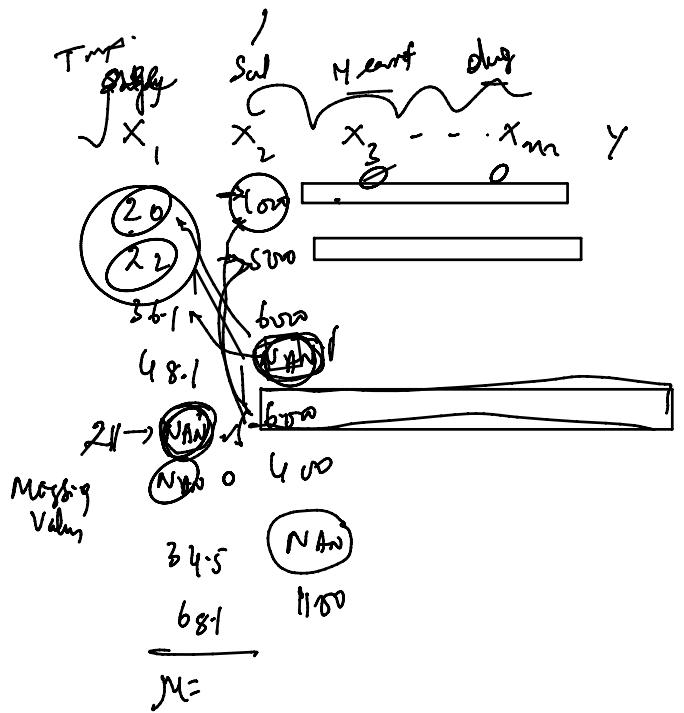
$$y \sim 0, 1$$

$$Y \sim \{0 - 9\}$$

Imputation? $\{KNN \text{ Imputation}\}$

Handling missing values

$$L \cdot R = \theta_0 + \theta_1 \overset{\text{miss}}{\underset{\text{miss}}{\circlearrowleft}} x_1$$



Replacing missing -(Imputation)

Imputation \rightarrow Mean Imputation

Median Imputation

Zero Imputation

Placeholder $\leftarrow (99)$