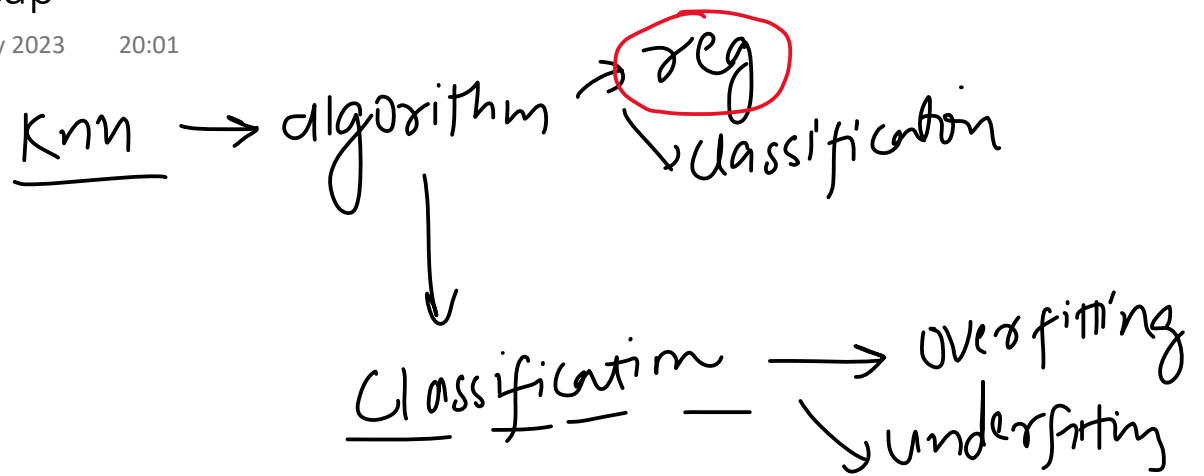


Recap

24 May 2023 20:01



failure case of knn

KNN Regressor

24 May 2023 14:44

Knn \rightarrow regression

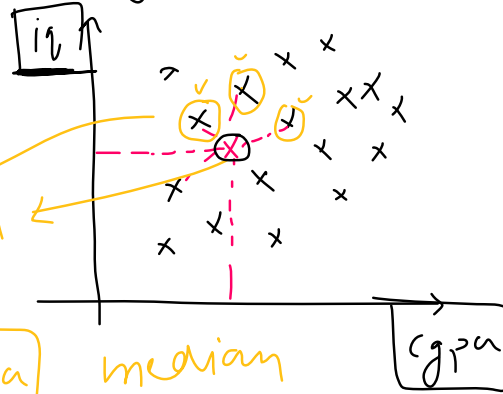
100
cgpa | iq package
K=3

K
X \rightarrow Y

3
5
4

$12/3 = 4$ lpa

median



1) find distance
 x_q and all the
 x_{train}
sort \rightarrow (3)

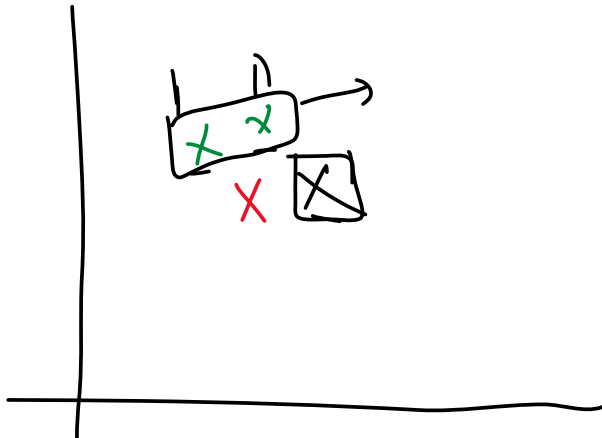
K high \rightarrow underfit

K low \rightarrow overfitting

Hyperparameters

24 May 2023 14:44

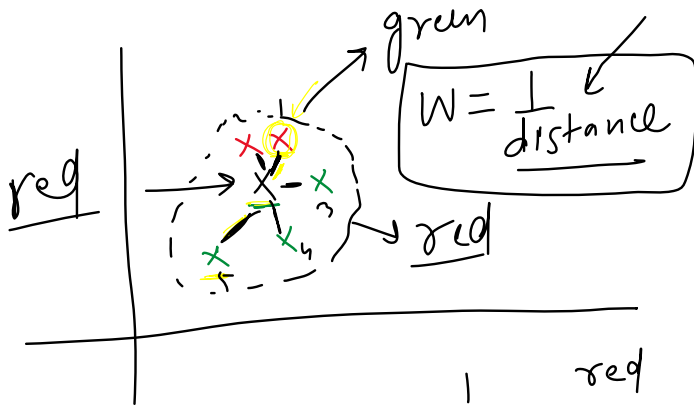
$k=3$



uniform

weighted knn → (knn)

K=5 classification



KNN → uniform

red → 2

green → 3

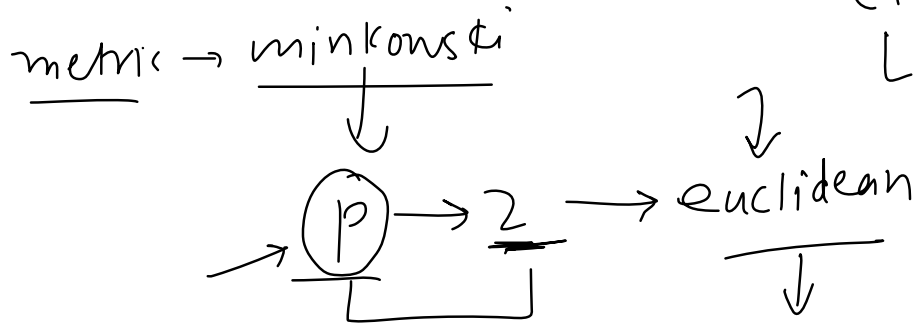
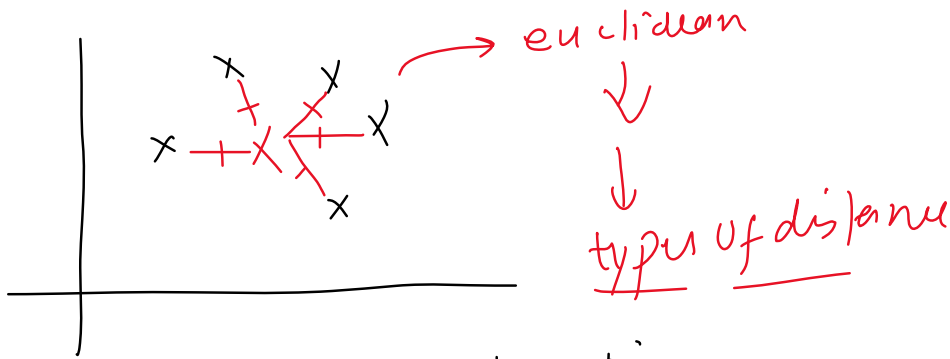
dist wt

label → green

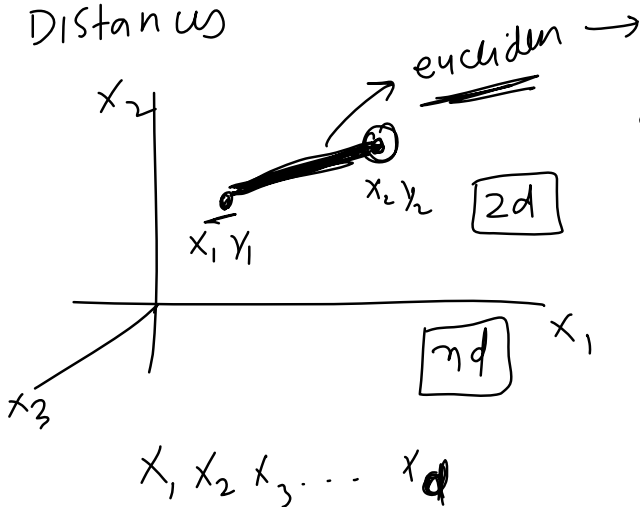
1	red	—	0.2	5
2	red	—	0.5	2
3	green	—	1	1
4	green	—	2	0.5
5	green	—	3	0.33

red → 5 + 2 = 7

green → $\frac{1 + 0.5 + 0.33}{1.83}$



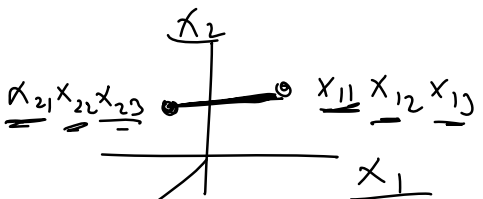
Distances



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

$$\text{dist} = \sqrt{\sum_{i=1}^d (x_{2i} - x_{1i})^2}$$

3 dim x_1, x_2, x_3



$$\sqrt{(x_{21} - x_{11})^2 + (x_{22} - x_{12})^2 + (x_{23} - x_{13})^2}$$

$$\text{dist} = \left(\sum_{i=1}^d (x_{2i} - x_{1i})^2 \right)^{\frac{1}{2}}$$

n dim points

Vectors

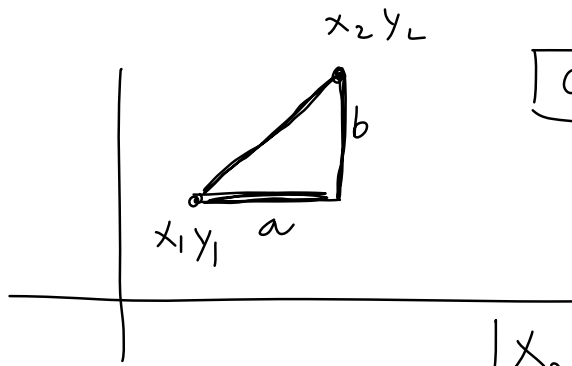
norm

L2 norm

$$m = \sum_{i=1}^d |x_{2i} - x_{1i}|$$

Manhattan distance

Manhattan distance

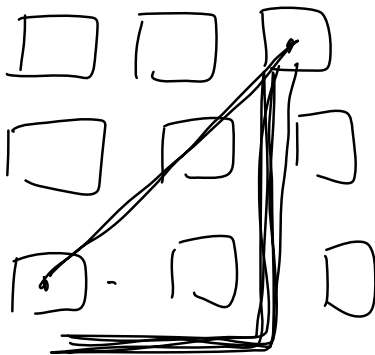


$$a+b$$

→ manhattan distance

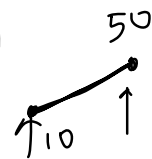
↓
taxi cab distance

$$|x_2 - x_1| + |y_2 - y_1| \rightarrow \text{man}$$



(0-cr)

Salary



$$\sqrt{100 + (50)^2}$$

② absolute

exp (1-20)

problems → euclidean

1) same scale ←

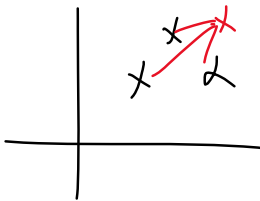
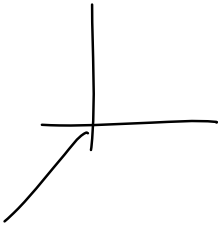
2) curse of dimension

↑ manhattan

euclidean

reliable

$$d = 100$$



p=2 → euclidean

p=1 manhattan

minkowski

$$\left(\sum_{i=1}^d (x_{2i} - x_{1i})^2 \right)^{1/2}$$

p=2 L2

$$\left(\sum_{i=1}^d |x_{2i} - x_{1i}| \right)^{1/1}$$

p=1 L1

eucl →

man

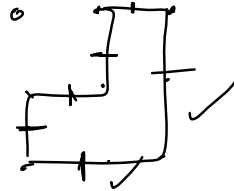
man
↓
general
0 → higher

$$\left(\sum_{i=1}^d |x_{2i} - x_{1i}| \right)$$

$$\left(\sum_{i=1}^d (|x_{2i} - x_{1i}|)^p \right)^{1/p}$$

$$p > 0$$

$\left(\sum_{i=1}^d (|x_{2i} - x_{1i}|)^p \right)^{1/p} \leftarrow 2 \text{ (each)}$
 $\leftarrow 1 \text{ (man)}$
 $\underline{3} \quad 13$
 $\underline{4} \quad 44$
 $2.5 \quad \underline{12.5}$

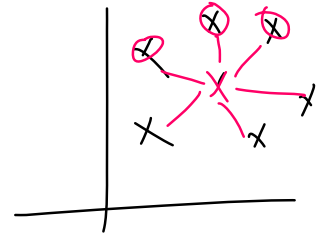


Knn \rightarrow slow algorithm

↓
predict

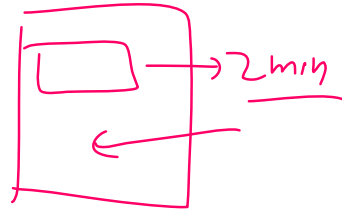
→ train

↑
production



Sort
↓
k nearest
↓
major

← result



time complexity

outline $f(n)$ double
↑
no of rows

Time complexity

$O(nd)$
↑
 $O(nd)$

$n \rightarrow$ # of rows in training data
 $d \rightarrow$ # features

$x_1, x_2, x_3, \dots, x_d | y$
↓
(10L, 1000)

⊗

euclid

$$\sqrt{\sum_{i=1}^d (x_{2i} - x_{1i})^2}$$

↑
d times (2) \rightarrow 1 cr

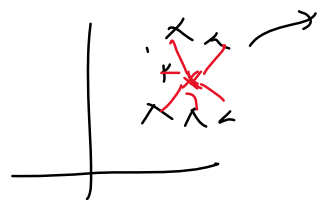
$2 \times 10^8 \rightarrow$ distance

space complexity

(1GB) \rightarrow 20m

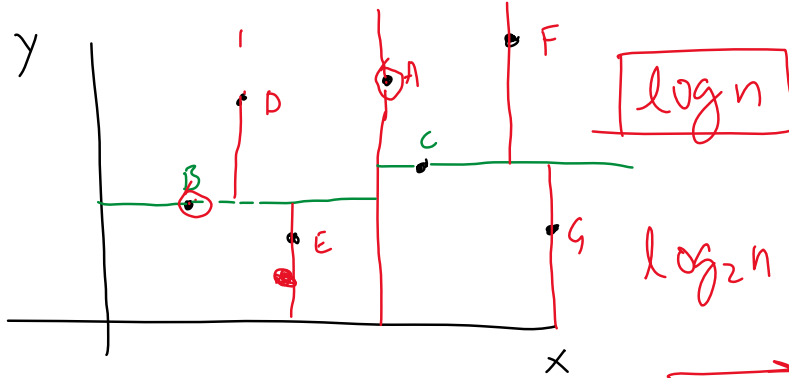
$O(nd)$
↑

train \rightarrow 1 lac



Binary search tree

4-d rce


$$x_9 \rightarrow \begin{pmatrix} 5, 10, \\ \uparrow \quad \uparrow \end{pmatrix}$$

K-d tree

2d-tree

