# Machine Learning

شوال ۱٤٤۳ – أيار/مايو ۲۰۲۲

| M | T | W | T | F | S | S |
|---|---|---|---|---|---|---|
| | | | 1 | 2 | 3 | 4 | 5 |
| 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 | 17 | 18 | 19 |
| 20 | 21 | 22 | 23 | 24 | 25 | 26 |
| 27 | 28 | 29 | 30 | | | |

## K Nearest Neighbours Algorithm.

Let

| CGPA | iq | Placement |
|------|-----|-----------|
| 8 | 80 | 1 |
| 7 | 70 | 0 |
| . | . | . |
| . | . | . |
| . | . | . |

100

Our goal is to build a machine learning model

$$\Downarrow$$

cgpa | iq $\rightarrow$ prediction

iq



$x \rightarrow 1$

$x \rightarrow 0$

cgpa

$x \rightarrow$ query point

we decide a value of K let say we assume $K = 3$

3 nearest nei[ghbours]

generally we find the euclidian distance.
we have calculated 100 distances.
then sort in ascending order
then majority count.    1 1 0 $\longrightarrow$ ①

M T W T F S S
May 30 31 1
2 3 4 5 6 7 8
9 10 11 12 13 14 15
16 17 18 19 20 21 22
23 24 25 26 27 28 29

اثنان/ مايو ٢٠٢٢ - شوال ١٤٤٣

MAY 2022

# How to select K?

$$\boxed{K = ?}$$

It depends on data

heuristic approach (jugar)

$\sqrt{n}$

$n \Rightarrow$ no. of observations

↓

generally avoid even value for K.

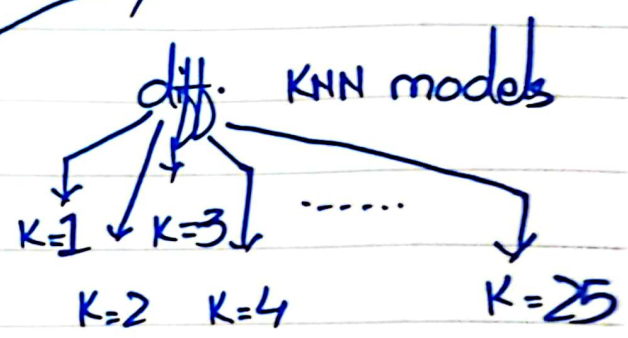experimentation

↓

cross validation
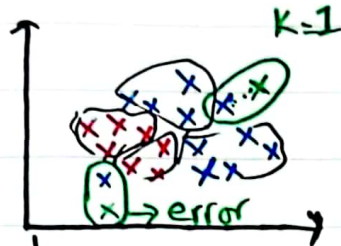
$n = 1000$

800 (train)      200 (test)

diff. KNN models

$K=1$  $K=3$  ......  $K=25$

$K=2$  $K=4$

train all these models on above dataset
and select that model
having highest accuracy score.

شوال ۱٤٤۳ هـ – ايّار/مايو ۲۰۲۲

MAY 2022

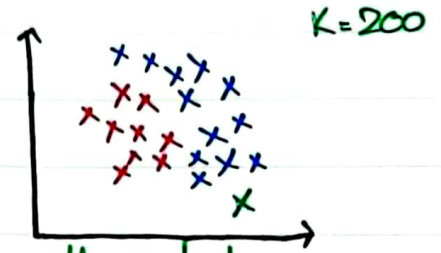| M | T | W | T | F | S | S |
|---|---|---|---|---|---|---|
| | | | 1 | 2 | 3 | 4 | 5 |
| 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 | 17 | 18 | 19 |
| 20 | 21 | 22 | 23 | 24 | 25 | 26 |
| 27 | 28 | 29 | 30 | | | |

June

# Overfitting and Under fitting in KNN:

lel we have data for 200 students.

Consider two scenario.



K=1

x → error

decision surface
would be divided
in multiple regions
overfitting → minor
        changes



K=200

all points lie in neighbor
        blue > red
i-e result always blue

## Limitations:

1) large datasets → $n = 5L$, $p = 100$      prediction
        it is lazy learning technique.      slow
2) High dimensional data
        curse of dimension
3) Outliers
4) Non-homogenous scales
5) Imbalance dataset      Yes – 98%,   No – 2%
6) Inference and not for prediction.