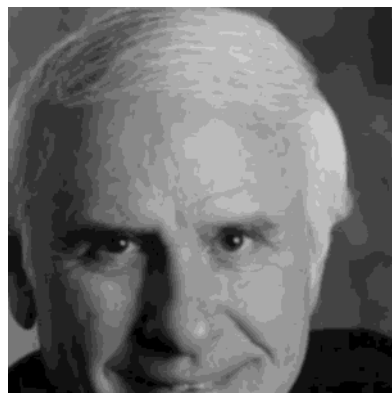


KNN Intuition

22 May 2023 19:14



**You are the
average of the five
people you spend
the most time
with**
- Jim Rohn

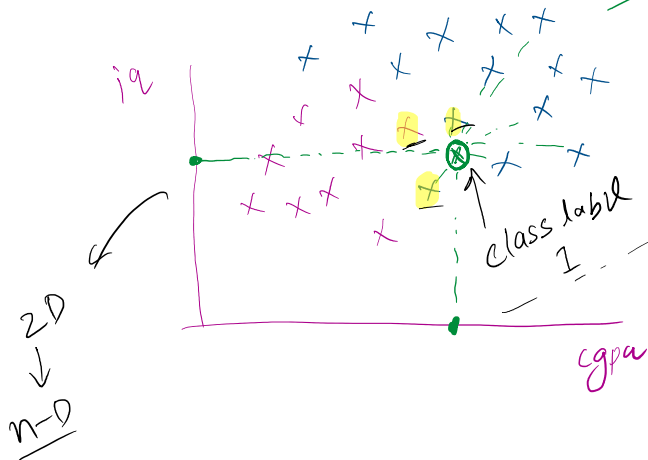
Knn → simplest
elegant

100 student

cgpa	iq	placement
------	----	-----------

(2D) ← ml model
↓
cgpa → placement
iq

euclidean

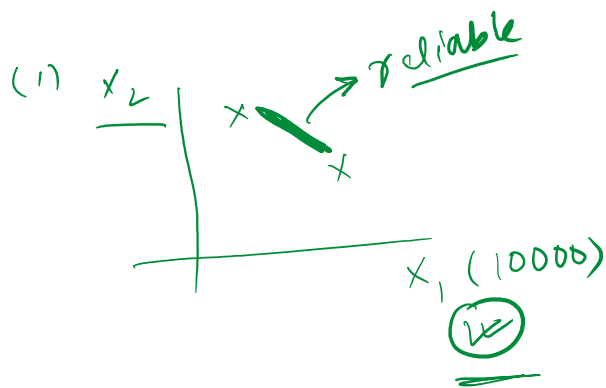
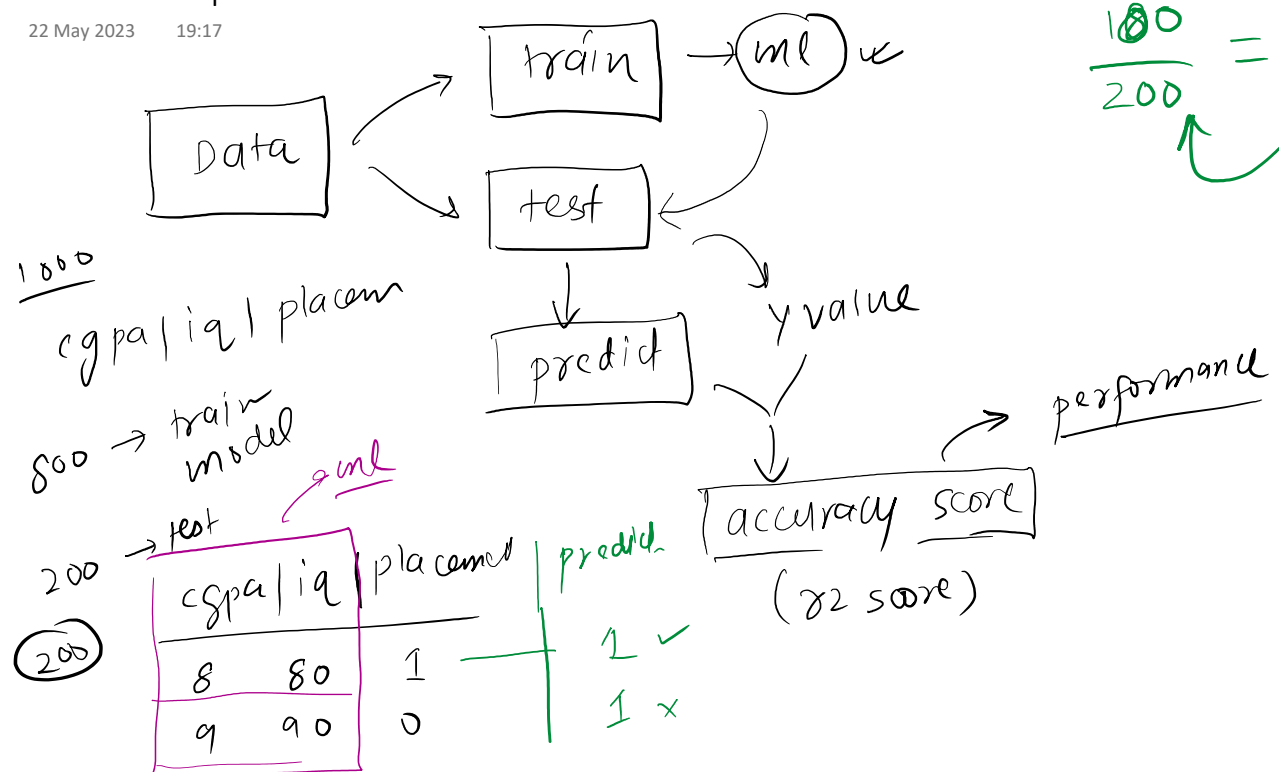


$x \rightarrow 1$
 $x \rightarrow 0$
 $K=3$
↑
3 nearest neighbors
↓
sort ascending
closest
↓
majority count
 $1 \ 1 \ 0 \rightarrow 1$
↑ ↑ ↑

$x_1 \mid x_2 \mid x_3 \dots x_n$ n-dim
↓
3 vectors → majority count → class label of query point

Code Example

22 May 2023 19:17



same scale

age

72 $\rightarrow \frac{72 - \mu}{\sigma}$

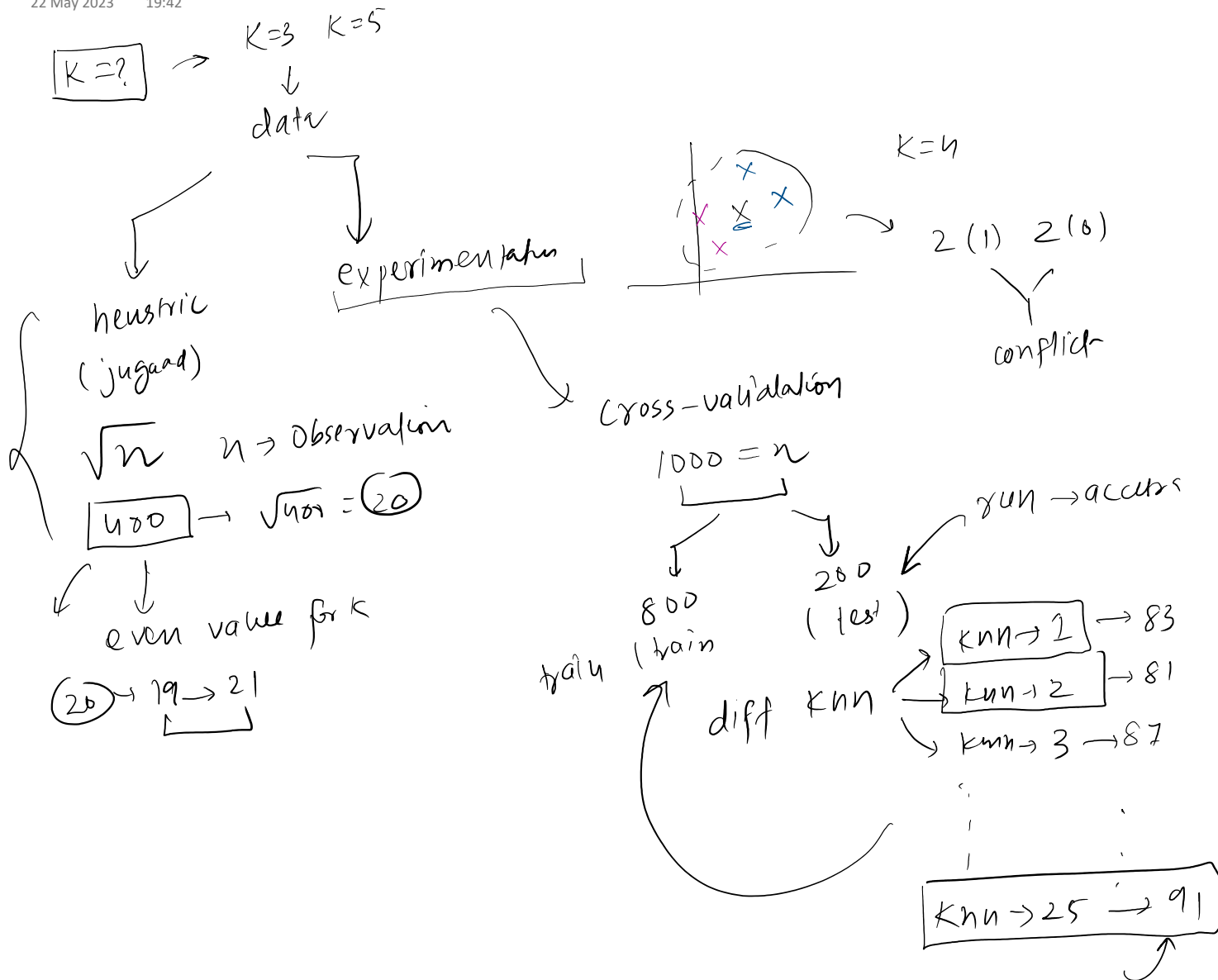
73 $\rightarrow \frac{73 - \mu}{\sigma}$

$\rightarrow 61$

$\rightarrow 37$

How to select K?

22 May 2023 19:42



[Decision Surface]

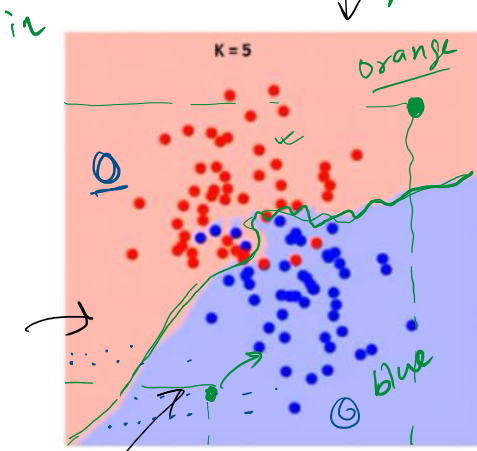
22 May 2023 19:15

tool → classification → knn → svm
 lr
 dt
 nn

1D 2D 3D

4Dmings → mixtend

plot decision
surface



coordinate
 disc video

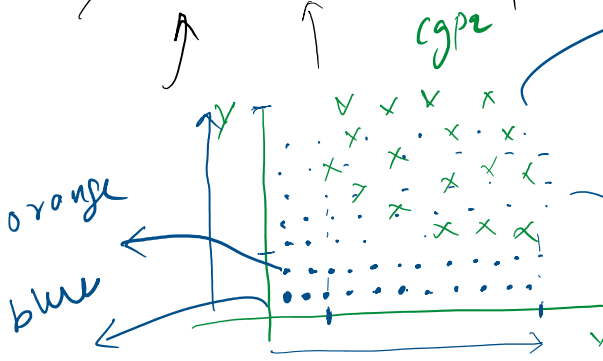
cgpa 1/2 | place binary
 [0,1]
 ↑ ↑

training points

pixel on
 anim

generate
 a numpy
 meshgrid

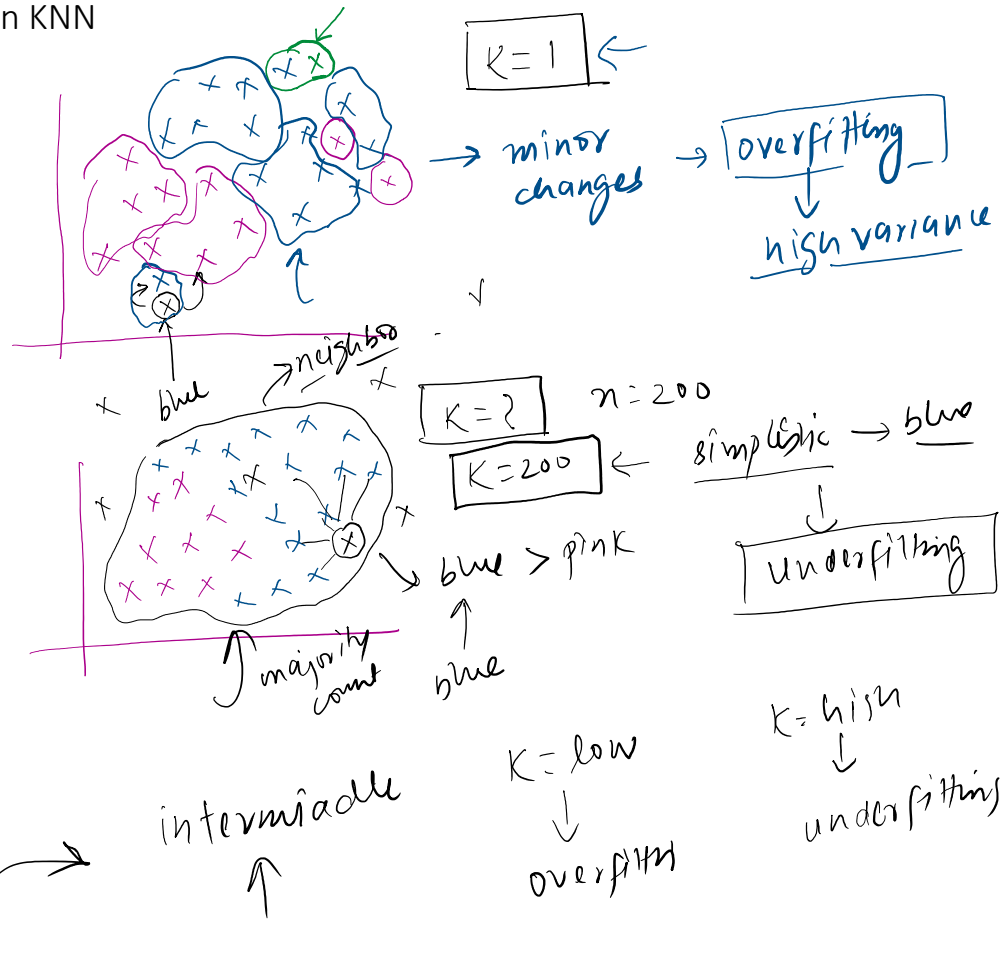
knn train → 1
 → 0



Overfitting and Underfitting in KNN

22 May 2023 19:15

cgpa | iq | placed
 2.5 student



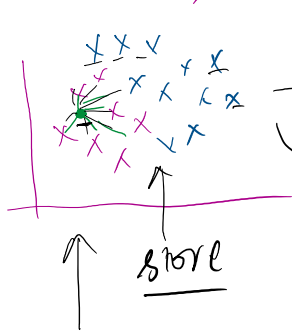
Limitations of KNN

22 May 2023 19:32

failure case

1) large datasets → $n = 5L$, $f = 100$

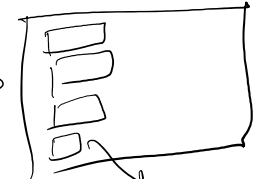
↳ Knn is a lazy learning techn



query point → prediction
training → nothing

5L distance → sort → majority

low latency



3 sec → slows

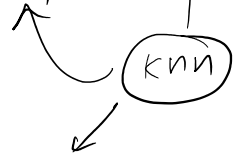
prediction
↓
slow → dataset

2) High dim data

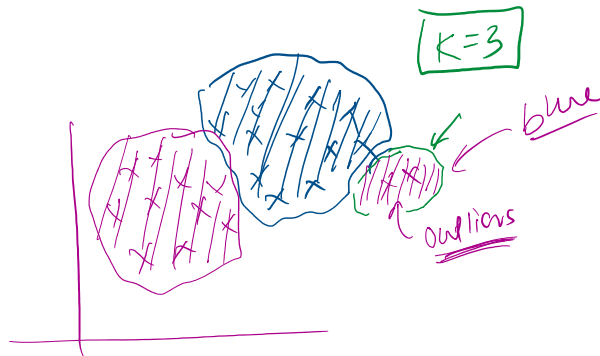
→ $f = 500$

↓
curse of dimension

distance → reliable
concept

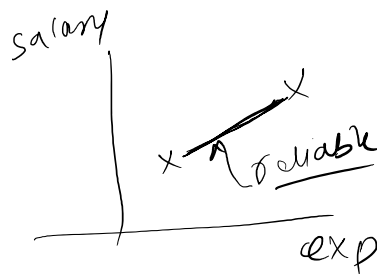


3) Outliers



4) Non-homogeneous scales

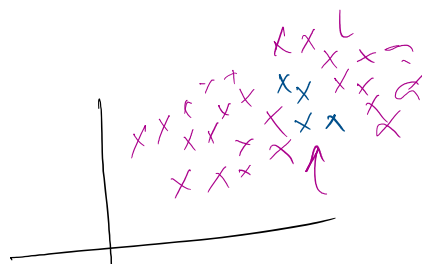
exp | salary | fire
0-25 | 20K-1cr



→ scaling

5) Imbalanced dataset

↳ Yes → 98%

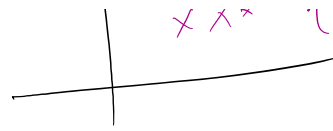


5) Imbalanced

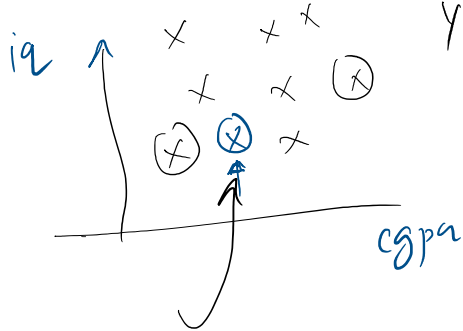
→ Yes → 98%

No → 2%

→ biased → print

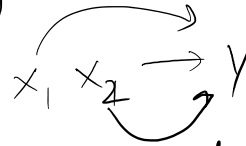


6) Inference and not for prediction



$$y = f(x)$$

↑



black box model →