(Black box model)   K- Nearest Neighbour (KNN) (classification)

iq



new query

gpa

Blue → placement no gaya (1)

Black → No placement (0)

Suppose a new query comes then we will calculate its distance (euclidean) from all points. Now we sort it in ascending order. Now do majority count

Lets for pt 1 dist of query → 0                    $k=3$
$\quad\quad\quad\quad\quad\quad$ 2 $\quad$ — 1
$\quad\quad\quad\quad\quad\quad$ 3 $\quad$ — 1

Since among the 3 nearest pt majority has placed so new query also get placed

• How to select $k$ ?

heuristic

∴ $k = \sqrt{n}$     $n = $ no. of observation

$\sqrt{400} = \dfrac{20}{1}$

don't take even
either take 19 or 21

experimentation

cross validation

$n = 1000$

$\underbrace{\quad\quad\quad}$
800               200

apply diff KNN
KNN = 1
$\quad\quad\quad$ 2
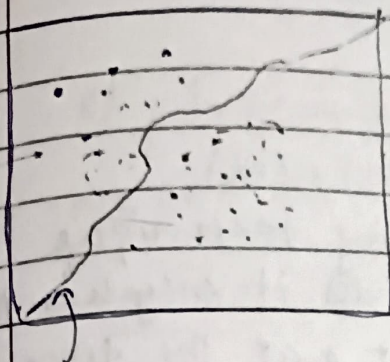$\quad\quad\quad$ 3
$\quad\quad\quad$ |

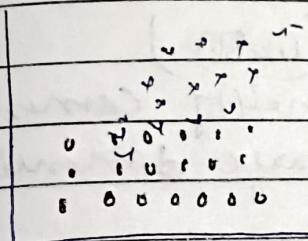The one which gives best result take that value of k

• Decision Surface
   ↳ It is a tool used in classification algorithm


Training Data

Suppose we considered a 2D data.
So There are two surface one for
0 & other for 1.
Now if any new pt. comes, we
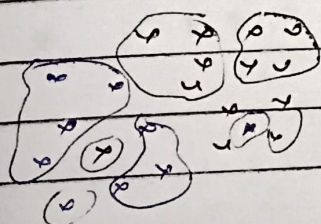can tell weather it will
fall in region 1 or reg. 0

Decision Boundry

→ how its made -



→ we plot the training data
→ then generate a numpy meshgrid
→ then apply knn which tells
   weather 0 or 1
→ Then we make pixels &
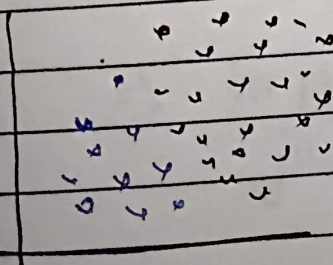   Show it as an image.

• Overfitting and Underfitting in KNN

n = 200

K : 1



→ This is overfitting as
   ↳ high variance

K = 200



↳ If a new pt comes then
now we will see its distance
from all pts.
Since black is in majority
So it will always win
→ underfitting

low value of k → overfitting
high  "  "  → under fitting

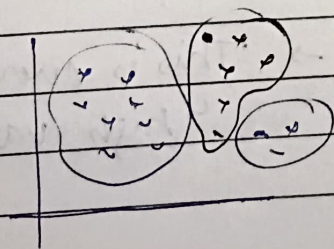- Limitations of KNN -

1) Large dataset .. (50000, 100)

   KNN is lazy learning technique means training is fast as it only has to store the points. but @ the time of prediction all calculation (Calculating dist from all pt, sorting etc) happens so it became slow

2) High Dimensional data (like 500)
   $\mathrel{\raise.3ex\hbox{$\llcorner$}}$ Curse of dimensionality comes into play which says at large dimension dist are not reliable.
   Since KNN is based on dist So it is also not reliable

3) Doesn't works good with outliers.



4) Non homogeneous Scales-

   | exp | Sal |
   |-----|-----|
   | 0-25 | 10k - 1lack |

   small scale    Large scale

5) Imbalanced data set

  ↳ 98% Yes

  ↳ 2% No

6) ↓ Inference    So (Black Box model)

not good for

## • Decision Boundry :-

  ↳ we can draw decision boundry for all classificam algo including Neural Network

  ↳ It can be both Linear or non-linear

  ↳ for high dimension problems, it act as an hyperplane.

→ Voronoi Diagram - ✓

→ Meshgrids

$y$
```
7
6                    17   27   37   47
5      =)            16   26   36   46
                     15   25   35   45
x      1 2 3 4
           x              ⇓
```

```
1   2   3  4        7   7   7   7
1   2   3  4        6   6   6   6
1   2   3  4        5   5   5   5
      x y                 y y
```

=) for plotting Decision Boundry just watch lecture