# Agenda for today
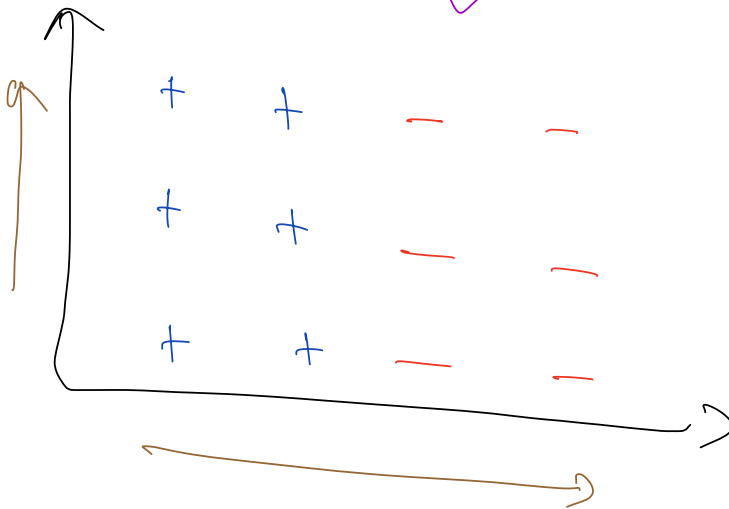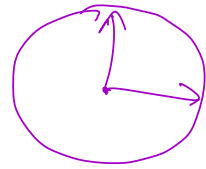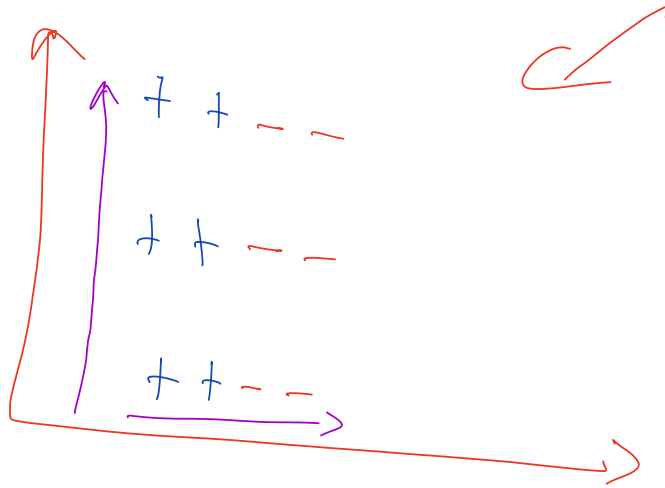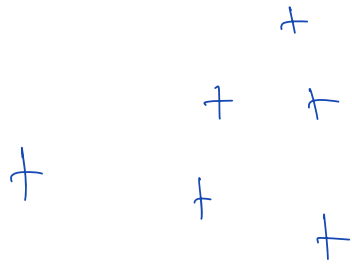
1) Blinkit Problem Statement
   ↳ issue with logistic regression

2) Geometric Intuition

3) kNN Algorithm

4) kNN scratch code

5) Assumptions of kNN

6) Sklearn's kNN implementation

7) Bias-Variance trade-off

8) Train & Test time Complexity ~

9) kNN for categorical data

10) LSH

11) kNN based imputation

$k = 9$

+ve → 3
0 → 3
-ve → 3

randomly choose
1 ot them

$$x_{q_1}, \quad x_{q_2}$$

$$x_{i_1}, \quad x_{i_2} \quad \left[ (x_{q_1} - x_{i_1})^2, \quad (x_{q_2} - x_{i_2})^2 \right]$$

$$\left( x_{q_1} - x_{i_1}, \quad x_{q_2} - x_{i_2} \right)^2$$

$$S = \left(x_{q1} - x_{i1}\right)^2 + \left(x_{q2} - x_{i2}\right)^2$$

$$ED = sqrt\_S = \sqrt{\left(x_{q1} - x_{i1}\right)^2 + \left(x_{q2} - x_{i2}\right)^2}$$

cluster counts

| idx | cluster | counts |
|-----|---------|--------|
| 0 | 2 | 3 |
| 1 | 1 | 4 |
| 2 | 3 | |

idx = 1

↓

has max count

pred = cluster [idx]

cluster [1]

↓

value 1

$q_1 = +$

$q_2 = +$

$q_3 = +$

K = 11

+ve → 5

−ve → 3

0 → 3

―――

11



underfit

$k = 1 \longrightarrow$ overfitting

$q_1 = 0$

$q_2 = -ve$

$q_3 = +ve$

$q_1 \quad to \quad q_3$

$\downarrow \qquad\qquad \downarrow$

$0 \qquad\qquad +ve$

slight change in query point $\rightsquigarrow$ change in pred

$k = large \implies underfit$

$k = small \implies overfit$



Acc ↑

Train acc

Val acc

overfit

underfit

overfit         $k \longrightarrow$         underfit

$$\underline{ED} = \left[ \sum_{i=1}^{d} \left( x_{q,i} - x_{j,i} \right)^{\boxed{2}} \right]^{1/2}$$

$$\left( x_q, \ x_j \right)$$

$$\left[ \ x_{q1}, \ \ x_{q2}, \ x_{q3}, \ \cdots \ x_{qd} \right]$$

$$\left[ \ x_{j1}, \ \ x_{j2}, \ \cdots \ x_{jd} \right]$$

Minkowski Distance

$$= \left[ \sum_{i=1}^{d} \left( x_{q,i} - x_{j,i} \right)^{p} \right]^{1/p}$$

*l2 norm*

$\underline{P = 2}, \ \Rightarrow$  Euclidean Distance

$\underline{P = 1}, \ \Rightarrow$  Manhattan Distance

*l1*

$$\left[ \sum_{i=1}^{d} \left| x_{q,i} - x_{j,i} \right| \right]^{1}$$

$\swarrow$ $L1$ metric

$M.D = d_1 + d_2$

$E.D = \sqrt{d_1^2 + d_2^2}$

$\curvearrowright$ $L2$ metric

$L1$ norm $\quad \sum |w_j|$

$L2$ norm $\quad \sum w_j^2$

N sample-points with d features

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ \vdots & & & \\ x_{N1} & x_{N2} & \cdots & x_{Nd} \end{bmatrix}$$

N rows, d columns

$N \times d$

## Test - time

Dist Calculation $\quad O(Nd) \longrightarrow$ for $N$ data-points in the training set

Sort : $\qquad O(N \log N)$

choosing Top k : $\qquad O(k)$

histogram: counts $\qquad O(k)$

majority vote : $\qquad O(k)$

$$O(Nd + N\log N + k + k + k)$$

$d << N$ ?? $\quad k << N$

$- ? X$

$$O\left(Nd + N\log N\right)$$
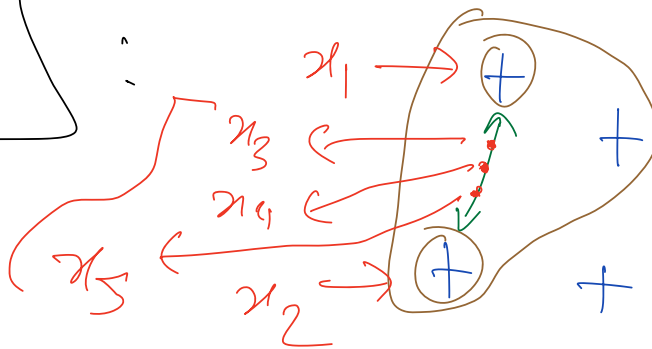
d is close to n
↳ is significant

K is very small compared to N

## Train time Complexity

$$O(1)$$

SMOTE :



$N = 5$
$K = 3$

$K = 3$

$x_1$    $x_2$            0.3

random value between 0 — 1

$$x_{new} = x_1 + 0.3 \times (x_2 - x_1)$$