# ML: K Nearest Neighbours KNN-2

Sumit.shukla_1@Scaler.com

# Summary

(Supervised ML-Algo)

## KNN

| | Age ($x_1$) | Income ($x_2$) | Default? ($y$) | Target variable (Dependent var) |
|---|---|---|---|---|
| A | | | Y ● | |
| B | | | N ● | |
| C | | | Y ● | |
| D | | | N ● | |
| E | | | Y ● | |

old customer / existing customer

→ Historical Data

New customer → $(x_q)$ Age = 19, Income = 10K

Income ($x_2$) vs Age ($x_1$) scatter plot with points f, g, a, c, b, e, d

$x_q$ ------

## Steps of KNN

① Compute distance from the query data point to all other points

② Sort the data points according to the corrosponding dinstances

$x_q \longrightarrow [a, b, d, e, c, f, g \ldots\ldots]$

③ choose k-neast neighbour (K=3)

④ voting

$a \longrightarrow Default = Y$
$b \longrightarrow Default = N$
$d \longrightarrow Default = Y$
$\left. \right\} < \begin{array}{l} Y = 2 \\ N = 1 \end{array} > [x_q \longrightarrow default = Yes]$

$d(x_q, a) \longrightarrow 1$
$d(x_q, b) \longrightarrow 1.5$

$\boxed{k = 4} < \begin{array}{l} Y = 2 \\ N = 2 \end{array} >$ split Brain situation

$k = default = 5$
$3 \leq k \leq 11$
$k = odd$ value

$\underline{KNN}$ — Non-parametric algo

| m. Budget | Sales |
|-----------|-------|
|           |       |

① Linear Regression $\longrightarrow$ Train Data $(x, y)$ $\longrightarrow$ [Linear Reg. Algo] $\longrightarrow$ (parameters) Best weights $\overline{W}, w_0$ $\longrightarrow$ prediction

② Logistic Regression

(Learning the best fit line)
line $\rightarrow \overline{W}x + w_0$

Sales $= w$ (m. Budget) $+ w_0$

estimate sales $\underset{1000}{\searrow}$

$\downarrow$ Parametric

Train Data $(x, y)$ $\longrightarrow$ [Logistic Reg. Algo] $\longrightarrow$ (Parameter) Best weights $\overline{W}, w_0$ $\longrightarrow$ predictions

(Learning the best Separating hyperplane)
$\overline{W}x + w_0$

non-parametric

KNN $\longrightarrow$ No parameter is being learnt

Hyperparameter

$\hookrightarrow$ In the case of KNN to predict the outcome of any new data point, we need to iterate through the entire training data and here the algo. doesn't learn any parameter

① parameter $\longrightarrow$ this is what the algo learns from the data ($\overline{W}, w_0$)

② Hyperparameter $\longrightarrow$ Has to be provided by the user and the algo will behave differently for different values of Hyper-parameter

# K impact on Bias & variance

K=1 (we are dependent on few neighbours)

overfitting



$2B = modi$
$3 = R.G$

In case k is very small, we are being dependent on few neighbours and if any of the neighbour change, the outcome may get impacted ⟶ overfitting [High variance, Low Bias]

K = 100  In this case when k is very high, we are being dependent on so many neighbours and the change in one or two neighbours will not impact my find outcome but the results are over generalised and hence very high error is expected
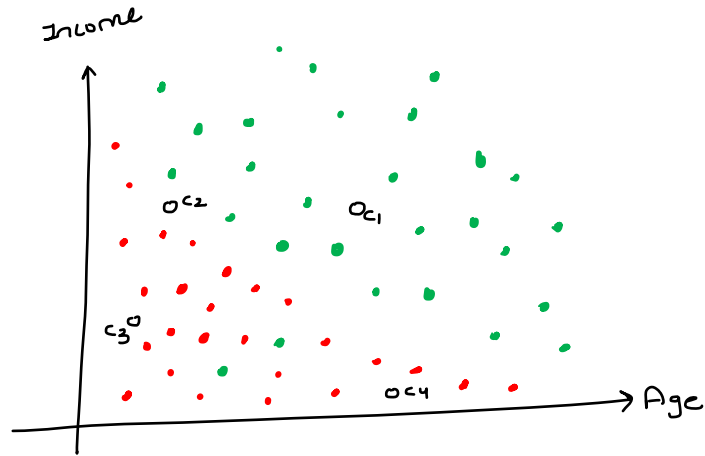
underfitting [Low variance, High Bias]

Hyperparameter Tuning

|  | Age | Income | y |
|----|-----|--------|---|
| C1 | | | |
| C2 | | | |
| C3 | | | |
| ⋮ | | | |
| C100 | | | |

Training Data



Income

$O^{C_2}$    $O_{C_1}$

$C_3 O$

$O C_4$    Age

|  | Age | Income | y |
|-----|-----|--------|---|
| C1 | | | |
| C2 | | | |
| C3 | | | |
| ⋮ | | | |
| C10 | | | |

Testing

for each data point in test data
we will run the entire KNN
steps

$C_1$ ⟶ KNN (Training data)

$C_2$ ⟶ KNN ( " " )

$C_3$ ⟶

⋮

# Time complexity of KNN

| $x_1$ | $x_2$ | $x_3$ | . . . . | $x_d$ |
|-------|-------|-------|---------|-------|
|       |       |       |         |       |
|       |       |       |         |       |

$n$

$x_q \rightarrow [x_{1q}, x_{2q} \cdots x_{dq}]$

① Distance calculation $\longrightarrow$ $O(nd)$

② Sort distance ascending $\longrightarrow$ $O(n\log n)$

③ pick k-values $\rightarrow O(k)$ $\Big\}$ very low

④ voting $\longrightarrow$ very less

Time complexity = $O(nd + n\log n)$

$D = |x_1 - x_{1q}| + |x_2 - x_{2q}|$
$\qquad + \cdots + |x_d - x_{dq}|$

$\vdots$

$\widehat{n}$

## problem with KNN

① As n increase KNN will go slow

② As d increase KNN will go slow

## Advantages

① KNN con generate non-linear boundary

② Simple to understand

③ very powerful on smaller data

④ very easily applied in multi-class classification.
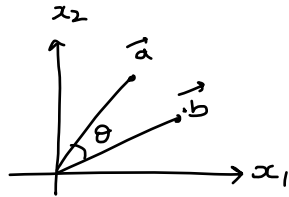
# Various Distance consideration
## KNN

① minkowski Distance

$$d = \left[ \sum_d |x_{ad} - x_{bd}|^p \right]^{1/p} \longrightarrow \text{for custom distance metric}$$
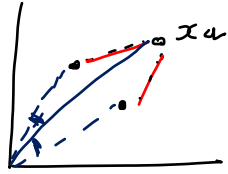
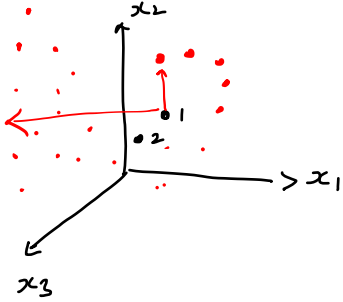② $P = 1 \, [L1]$

$$d = \sum_d |x_{ad} - x_{bd}|$$

③ $P = 2 \, [L2]$

$$d = \left( \sum_d |x_{ad} - x_{bd}|^2 \right)^{1/2} \longrightarrow \text{Data is less}$$

④ Cosine Similarity $\quad \cos(\theta) = \dfrac{\vec{a}^T \cdot \vec{b}}{\|\bar{a}\| \, \|\bar{b}\|} \longrightarrow \text{In the case of highly dimensional data}$

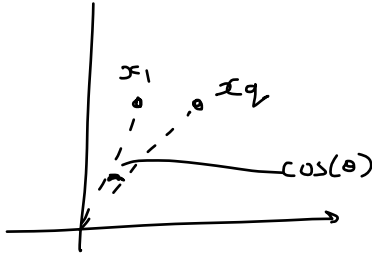$$\mathcal{E} \cdot d = \left[ (x_q - x_1)^2 + (x_q - x_2)^2 + \cdots\cdots + (x_q - x_{100})^2 \right]^{1/2}$$

$$\mathcal{E}d(1) \stackrel{\backsim}{=} \mathcal{E}d(2) \quad \text{Because of higher dimensions}$$



cosine similarity in higher dimension
the angles $\cos(\theta)$ appears very
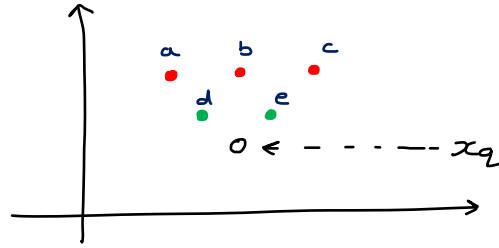different and easy for algo to differentiate.
=

# Weighted KNN

$d(x_q, a) = 0.5$

$d(x_q, b) = 0.5$

$d(x_q, c) = 0.5$

$d(x_q, d) = 0.2$

$d(x_q, e) = 0.2$



$K = 5$

$x_q \longrightarrow$ Red

Idea $\longrightarrow$ Data point close to the query data point should be given more weight in the voting

weight $\propto \dfrac{1}{\text{distance}}$

$$\omega(i) = \dfrac{1}{d(i)}$$

$\omega_1 = \dfrac{1}{0.5} = 2$

$\omega_2 = \dfrac{1}{0.5} = 2$   Red

$\omega_3 = \dfrac{1}{0.5} = 2$

$\omega_4 = \dfrac{1}{0.2} = 5$   green

$\omega_5 = \dfrac{1}{0.2} = 5$

Red $= [2+2+2] = 6$

$x_q$

green $[5+5] = 10$

$\downarrow$

$\boxed{\text{green}}$

KNN- Imputation (for missing value imputation)

| | Age | Income | Tenure | Distance |
|---|---|---|---|---|
| A | 21 | 30 | 2 | 4.24 |
| B | 22 | 34 | 5 | 2.24 |
| C | 24 | 38 | 1 | 5.00 |
| D | 25 | 32 | 2 | 1.41 |
| E | 24 | 31 | 4 | 2.00 |
| q | 24 | 33 | NA | |

The missing data point can be predicted
by taking avg of closest neigbors

$$K = 3$$

$$Tenure(q) = \frac{Tenure(B) + Tenure(E) + Tenure(D)}{3}$$

$$= \frac{5 + 2 + 4}{3} = 3.0$$