

15-11-2025

Agenda:

previous class:

1. Decision Tree Regressor (Code Implementation) ●
 2. Random Forest Classifier & Regressor (Code Implementation) ●
 3. Gradient Boosting (concept)
 4. XGBoost (concept)
- Regression (concept)

Today's class:



- Hierarchical clustering → concept & implementation
- DBSCAN

Gradient Boost (classification)

| ID | Size | Beds | price | $y(\text{price} \geq 170)$ | P0 |
|----|------|------|-------|----------------------------|-----|
| 1 | 800 | 2 | 120 | 0 | 0.9 |
| 2 | 900 | 2 | 130 | 0 | 0.9 |
| 3 | 1000 | 3 | 150 | 0 | 0.9 |
| 4 | 1100 | 3 | 170 | 1 | 0.9 |
| 5 | 1200 | 4 | 200 | 1 | 0.9 |

$$\text{Step 1: } f_0 : \text{log-odds} = \log \frac{p}{1-p} \Rightarrow \log \frac{0.4}{1-0.4} \Rightarrow \log \frac{0.4}{0.6}$$

$$p = \frac{2}{5} = 0.4 \quad \approx -0.40546$$

$$p_0 = \sigma(f_0) = 0.40$$

Step 2:

$$\text{using logit loss } L = -y \log p + (1-y) \log (1-p)$$

where $p = \sigma(f_0)$

$$\text{gradient } g = \frac{\partial L}{\partial p} = p - y$$

$$\text{hessian, } h = p(1-p)$$

$$\text{compute } g = p - y, h = 0.4(0.6) \Rightarrow 0.24$$

| ID | Size | Beds | price | $y(\text{price} \geq 170)$ | p_0 | $g(p-y)$ | $h=0.24$ |
|----|------|------|-------|----------------------------|-------|----------|----------|
| 1 | 800 | 2 | 120 | 0 | 0.4 | 0.4 | 0.24 |
| 2 | 900 | 2 | 130 | 0 | 0.4 | 0.4 | 0.24 |
| 3 | 1000 | 3 | 150 | 0 | 0.4 | 0.4 | 0.24 |
| 4 | 1100 | 3 | 170 | 1 | 0.4 | -0.6 | 0.24 |
| 5 | 1200 | 4 | 200 | 1 | 0.4 | -0.6 | 0.24 |

Step 3: creation of stump (depth=1) or decision tree

→ mid point

size: 850, 950, 1050, 1150 price: 125, 140, 160, 185

Beds: 2.5, 3.5

few possibilities: Red ≥ 1.5 | size ≥ 250 | size ≥ 1050 | price ≥ 160

we selected $\underline{\text{size} \geq 1050}$ or $\underline{\text{size} < 1050}$:

$1, 2, 3 (\text{FD'})$ $4, 5 (\text{FD'})$

$$G_L = \sum g_i$$

$$G_R = \sum g_i$$

$$G_p = G_L + G_R$$

$$G_L = 0.4 + 0.4 + 0.9 \\ = 1.2$$

$$G_F = -0.6 + -0.6 = -1.2$$

$$H_L = 0.24 \times 3 = 0.72$$

$$H_R = 0.24 \times 2 = 0.48$$

$$G_p = G_L + G_R \\ = 1.2 + (-1.2) \\ = 0$$

$$H_p = 0.72 + 0.48 = 1.2$$

$$\text{gain: } \frac{1}{2} \left(\frac{G_L^2}{H_L} + \frac{G_F^2}{H_R} - \frac{G_p^2}{H_p} \right) \Rightarrow \frac{1}{2} \left(\frac{(1.2)^2}{0.72} + \frac{(-1.2)^2}{0.48} - 0 \right)$$

$$\Rightarrow \frac{5}{2} \Rightarrow \underline{\underline{2.5 \text{ gain}}}$$

fl

size ≥ 1050

$$\gamma_L = -\frac{G_L}{H_L} = -\frac{1.2}{0.72} = -1.66$$

$$\gamma_{size} = -\frac{G_{size}}{H_{size}}$$

$$\gamma_R = -\frac{G_R}{H_R} = -\frac{-1.2}{0.48} = \underline{\underline{2.5}}$$

Input → size 800 bed 2 price 120 →

$$\eta = 0.1$$

$$f = f_0(x) + \eta \cdot f_1(x)$$

$$= -0.40546 + 0.1 \cdot (-1.66) = -0.40546 - 0.166$$

$$\approx -0.5721$$

$$p_{q_n} = \sigma(f) = \sigma(-0.5721) \approx 0.38$$

Input → size 1800 bed 2 price 120 →

$$f = f_0(x) + \eta \cdot f_1(x)$$

$$= -0.40546 + 0.1 \cdot (2.5) = -0.40546 + 0.25$$

$$\approx -0.15546$$

$$p_{q_n} = \sigma(f) = \sigma(-0.15546) \approx 0.46$$

| ID | size | bed | price | $y(\text{price} \geq 170)$ | p_0 | $g(p-y)$ | $\eta = 0.24$ | f_1 | $g(f_1 - y)$ |
|----|------|-----|-------|----------------------------|-------|----------|---------------|-------|--------------|
| 1 | 800 | 2 | 120 | 0 | 0.4 | 0.4 | 0.24 | 0.26 | |
| 2 | 900 | 2 | 130 | 0 | 0.4 | 0.4 | 0.24 | 0.26 | |
| 3 | 1000 | 3 | 150 | 0 | 0.4 | 0.4 | 0.24 | 0.26 | |
| 4 | 1100 | 3 | 170 | 1 | 0.4 | -0.6 | 0.24 | 0.46 | |
| 5 | 1200 | 4 | 200 | 1 | 0.4 | -0.6 | 0.24 | 0.46 | |

Xgboost (Classification)

| ID | Size | Beds | price | $y(\text{price} \geq 170)$ | p_0 | $g(p-y)$ | $h = 0.29$ | f_1 | $g(f_1 - y)$ |
|----|------|------|-------|----------------------------|-------|----------|------------|-------|--------------|
| 1 | 800 | 2 | 120 | 0 | 0.9 | 0.9 | 0.29 | 0.36 | |
| 2 | 900 | 2 | 130 | 0 | 0.9 | 0.9 | 0.29 | 0.36 | |
| 3 | 1000 | 3 | 150 | 0 | 0.9 | 0.9 | 0.29 | 0.36 | |
| 4 | 1100 | 3 | 170 | 1 | 0.9 | -0.6 | 0.29 | 0.46 | |
| 5 | 1200 | 4 | 200 | 1 | 0.9 | -0.6 | 0.29 | 0.46 | |

Step 1: $f(0) \rightarrow$ initial model with log of odds $\rightarrow \log \frac{P}{1-P} = \frac{2}{5} = 0.9$

Step 2:

$$\text{using logistic loss } L = -y \log p + (1-y) \log (1-p)$$

we selected $\underline{\text{size} \geq 1050}$ or $\overline{\text{size} < 1050}$:

$$\begin{array}{ccc} \swarrow & & \searrow \\ 1, 2, 3 (\text{FD'}) & & 4, 5 (\text{FD'}) \end{array}$$

$$G_L = \sum g_i \quad G_R = \sum g_i$$

$$\downarrow \quad \downarrow$$

$$G_P = G_L + G_R$$

$$G_L = 0.4 + 0.9 + 0.9 \\ = 1.2$$

$$H_L = 0.29 \times 3 = 0.87$$

$$G_R = -0.6 + -0.6 = -1.2$$

$$H_R = 0.29 \times 2 = 0.48$$

$$\begin{aligned} G_P &= G_L + G_R \\ &= 1.2 + (-1.2) \\ &= 0 \end{aligned}$$

$$H_P = 0.72 + 0.48 = 1.2$$

Xgboost gain (include λ, γ):

$$\lambda = 1, \gamma = 0$$

$$\text{gain} = \frac{1}{2} \left(\frac{\omega_L^2}{H_L + \lambda} + \frac{\omega_R^2}{H_R + \lambda} - \frac{\omega_P^2}{H_P + \lambda} \right) - \gamma$$

$$= \frac{1}{2} \left(\frac{1.44}{0.72+1} + \frac{1.44}{0.48+1} - 0 \right) - 0$$

$$= \frac{1}{2} \left(0.83720 + 0.97297 \right)$$

$$\approx 0.90509 \leftarrow \underline{\text{gain}}$$

Gain (gradient boost)
↓

2.5

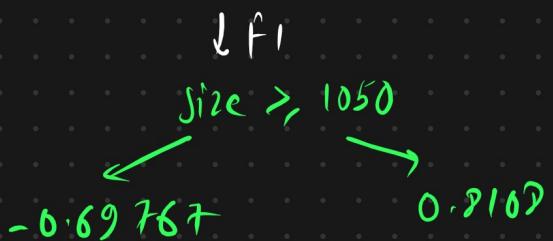
Gain (xgboost)
↓ λ

0.9051

leaf weight: $w = -\frac{\gamma}{H + \lambda}$

$$\text{left } w_L = -1.2 / 0.72 + 1 = -\underline{\underline{0.69787}}$$

$$\text{right } w_R = -(-1.2) / 0.48 + 1 = \underline{\underline{0.8108}}$$



$$\begin{array}{c} \downarrow f_1 \\ \text{size} > 1050 \\ -1.66 \quad 2.5 \\ (\text{GB}) \end{array}$$

$$\begin{array}{c} \downarrow f_1 \\ \text{size} > 1050 \\ -0.69767 \quad 0.8102 \\ (\chi_{\text{GB}}) \end{array}$$

$$\rightarrow F(x) = f_0(x) + \eta f_1(x)$$

$$x=800 \rightarrow -0.40546 + 0.1 \cdot (-0.69767) = -0.4752$$

$$x=1000 \rightarrow -0.40546 + 0.1 \cdot (0.8102) = -0.3243$$

$\underbrace{\hspace{10em}}$ sigmoid \downarrow

$$\approx 0.41 \qquad \approx 0.38$$

| ID | Size | Beds | price | $y(\text{price} \geq 170)$ | P_0 | $y(P_0)$ | $h=0.29$ | $f_1(\text{GB})$ | $f_1(\chi_{\text{GB}})$ |
|----|------|------|-------|----------------------------|-------|----------|----------|------------------|-------------------------|
| 1 | 800 | 2 | 120 | 0 | 0.4 | 0.4 | 0.29 | 0.36 | 0.38 |
| 2 | 900 | 2 | 130 | 0 | 0.4 | 0.4 | 0.29 | 0.36 | 0.38 |
| 3 | 1000 | 3 | 150 | 0 | 0.4 | 0.4 | 0.29 | 0.36 | 0.38 |
| 4 | 1100 | 3 | 170 | 1 | 0.4 | -0.6 | 0.29 | 0.46 | 0.41 |
| 5 | 1200 | 4 | 200 | 1 | 0.4 | -0.6 | 0.29 | 0.46 | 0.41 |

HIERARCHICAL CLUSTERING

↳ Unsupervised

$$\begin{array}{l} \text{Blue } 0-0.5 \\ \text{Red } > 0.5 \rightarrow 1 \end{array}$$

| x_1 | x_2 | x_3 |
|-------|-------|-------|
| 1.0 | 0.5 | 0.1 |
| 0.8 | 0.4 | 0.1 |
| 0.5 | 0.8 | 0.7 |
| 0.2 | 0.9 | 0.8 |



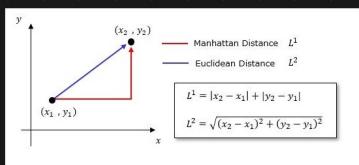
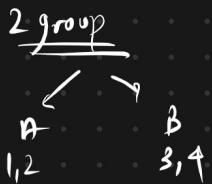
| x_1 | x_2 | x_3 |
|-------|-------|-------|
| 1.0 | 0.5 | 0.1 |
| 0.8 | 0.4 | 0.1 |
| 0.5 | 0.8 | 0.7 |
| 0.2 | 0.9 | 0.8 |

Step 1: iterate over data

Step 2: using distance formula check data closest to the current data.

Step 3: Merge the closest data with current data & save them as cluster.

Step 4: repeat step 1-3.



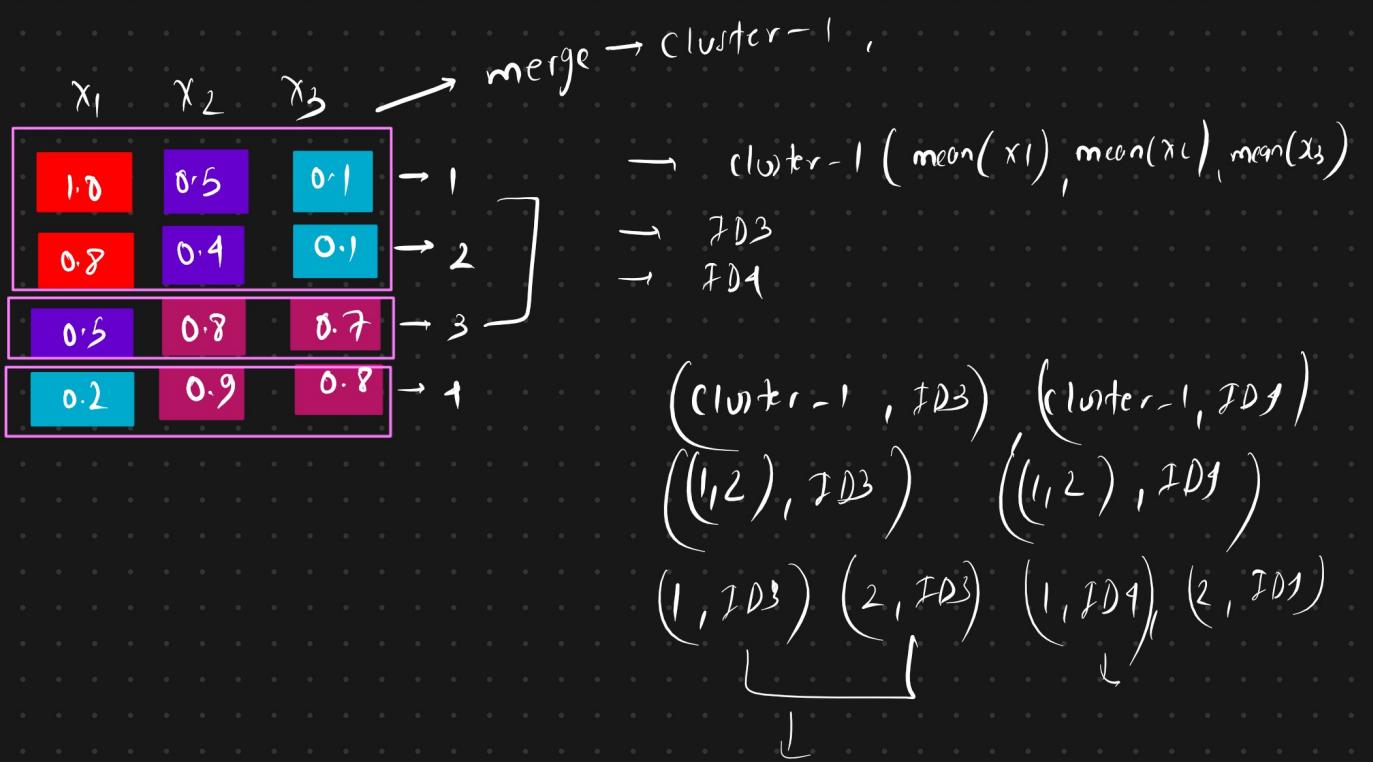
$$\begin{aligned} \boxed{\text{Id}(1)} &\rightarrow 1.0, 0.5, 0.1 \\ \text{Id}(2) &\rightarrow 0.8, 0.4, 0.1 \end{aligned}$$

| x_1 | x_2 | x_3 |
|-------|-------|-------|
| 1.0 | 0.5 | 0.1 |
| 0.8 | 0.4 | 0.1 |
| 0.5 | 0.8 | 0.7 |
| 0.2 | 0.9 | 0.8 |

$$\left(\begin{array}{l} \text{euclidean}(\text{Id}(1), \text{Id}(2)) \rightarrow \text{out } 1 \\ \text{euclidean}(\text{Id}(1), \text{Id}(3)) \rightarrow \text{out } 2 \\ \text{euclidean}(\text{Id}(1), \text{Id}(4)) \rightarrow \text{out } 3 \end{array} \right)$$

$$\sqrt{(1.0 - 0.8)^2 + (0.5 - 0.4)^2 + (0.1 - 0.1)^2} \rightarrow \underline{\text{out } 1}$$

$$\min(\text{out } 1, \text{out } 2, \text{out } 3) \rightarrow \underline{\underline{102}}$$



Distanz
 \rightarrow Euclidean
 \rightarrow Manhattan
 \rightarrow cosine

Linkage:

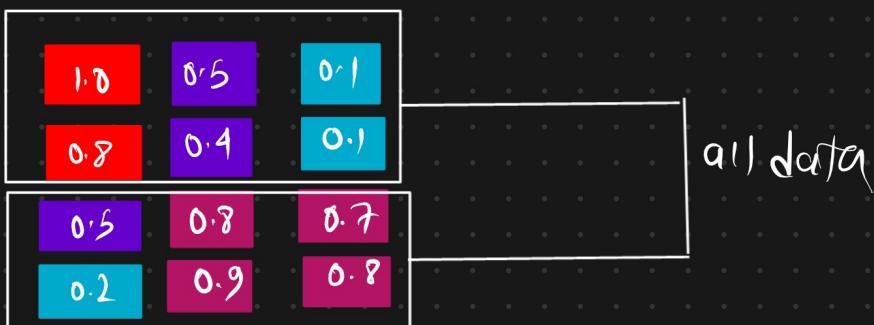
complete \rightarrow farthest

max distance

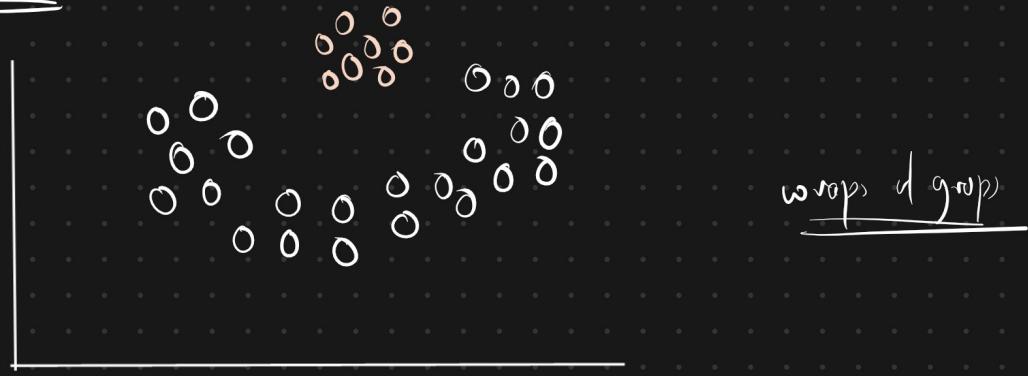
ward'

average

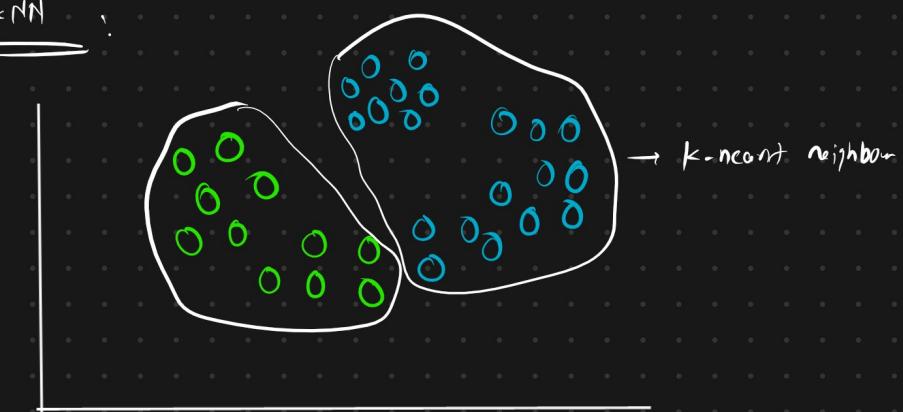
single \rightarrow nearest neighbour



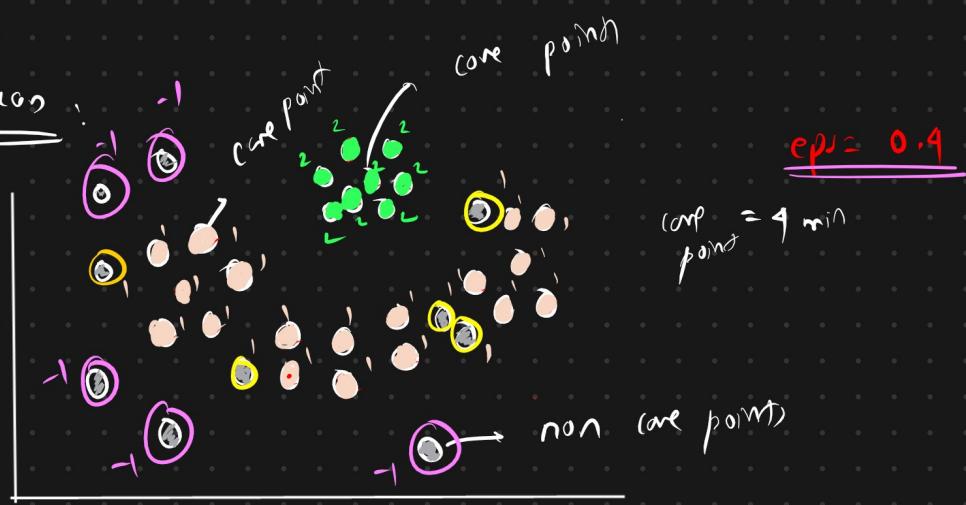
DBSCAN



kNN



DBSCAN

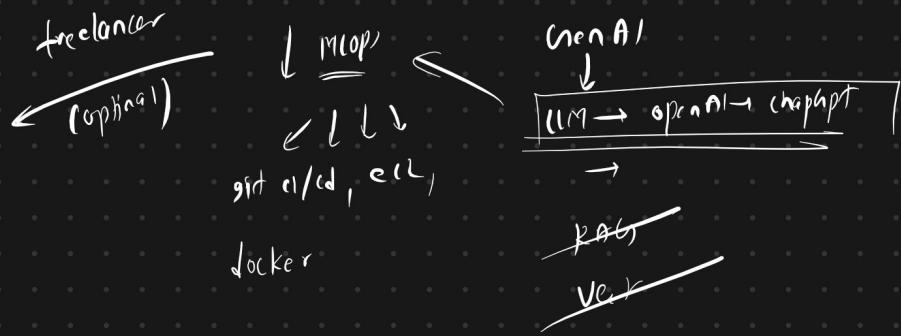
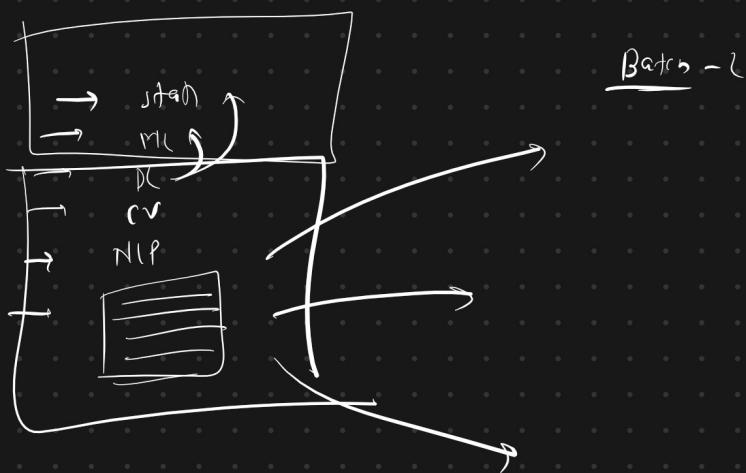


- we initialize core point randomly
- we check if it is surrounded by min-point or not!
- if Yes →
 - make all the point it is touching as core point
- if No →
 - proceed to another random point

- we will start with random core point with epk at area.
- assign data point as similar group which are core point
- for data point which are similar but non-core will be merge in group but not allowed to further extend the group.

(1) step 1

- create core & non-core point
- from core point → create group using data near to point
- note non-core point which are near to the above process
- group them, but they are still non-core
- groups (output) + outliers



start

- Python → recursion, OOP, decorator, list, dict comprehension, comp prog question
- SQL → SQL query (last thing you should do)
- Stats → p-value, mean, median, mode, distribution
 - pdf, cdf, hypothesis testing
- Machine learning →
 - Logistic regression : classification
 - loss function → \rightarrow regression
 - optimization → gradient descent
 - Activation function → sigmoid
 - optimal value of k in k-means, k-nearest
 - SVM → margin
 - kernels
 - cross-validation
 - Bootstrapping
 - Regularization ($L_1 \& L_2$)
 - Gradient Boost vs Xgboost
 - Decision tree → impurity \rightarrow entropy
 - \rightarrow gini
 - Boosting & Bagging
 - ↓
 - ↓
 - 1 algo 1 algo
 - project
 - resume → only keep project in which you are confident to talk about .
 - Biggest challenge you faced in a project /
 - project you are proud of ?
- swap 2 numbers
- fibonacci series
- factorial
- prime number or not
- leap year
- 1 simple sort method