Agenda → Project (Supervised ML Project)
*

(House Price Prediction)

(Baseline ML model)

I/P → Model → output

Price of the House

model → Push this to Production
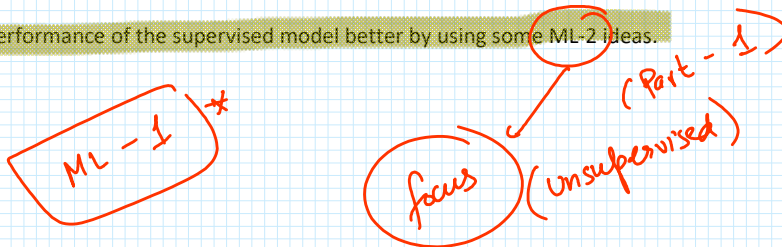
Baseline

main
(main code)
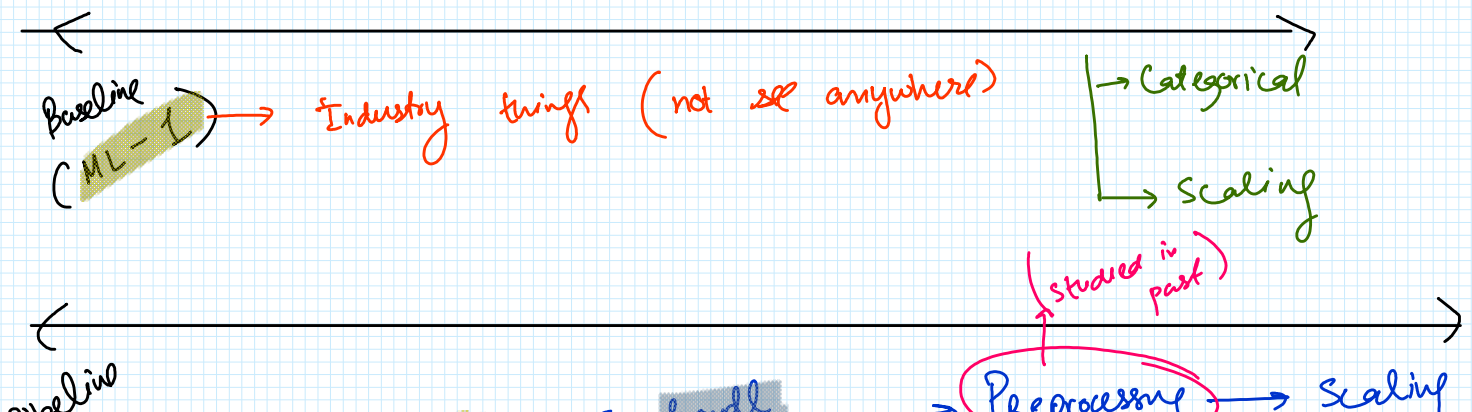(merge) (unsupervised - testing)

(Can we try something Better?)
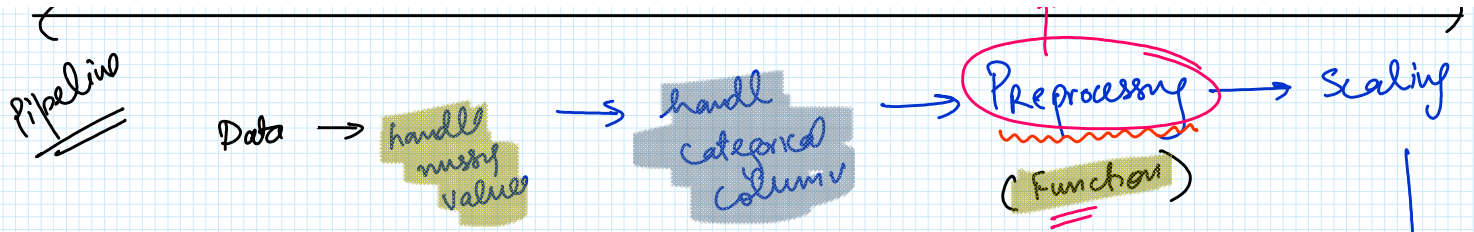→ (Unsupervised)

(Problem Statement)

Given a set of features, predict the prices of the houses.

ML → Supervised
ML-2 Unsup.

performance of the supervised model better by using some ML-2 ideas.

(Part - 1)

ML - 1 *

focus (unsupervised)

Baseline
(ML - 1) → Industry things (not see anywhere)
→ Categorical
→ Scaling
(studied in past)

Baseline
...model... → Preprocessing → Scaling

**Pipeline**

Data → handle missing values → handle categorical columns → **Preprocessing** (Function) → Scaling

↓

Fit the model ←

**Handle Categorical** → one hot encoding
                       → ordinal

→ Frequency encoding

(10 – 20 columns)

Target encoding (300 column) (multiple column) → Data Leakage

| Zip code | | |
|---|---|---|
| 22 | → | 10 |
| 22 | | 10 |
| 23 | → | 15 |
| 23 | | 15 |
| 23 | | ⋮ |
| 38 | | 15 |

13

Garage
S
S
S
S
M
M
M
L
L

Price
→ mean (S)
→ mean (M)
→ mean (S)

Advance (NLP)
Embedding method

(D.L)

## Handling Categorical Features

Categorical features can be **ordinal** (have a meaningful order) or **nominal** (no intrinsic order). The approach depends on the type:
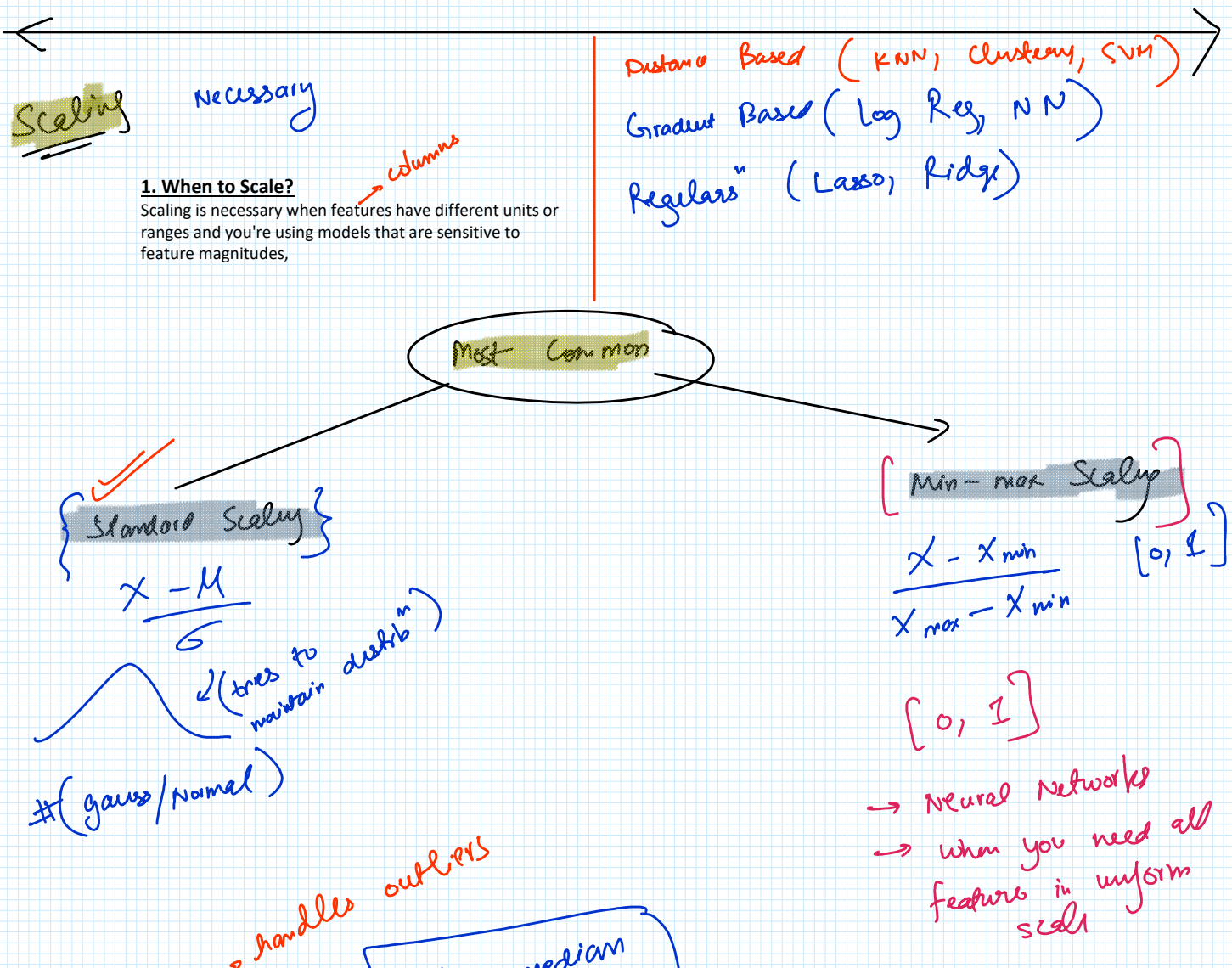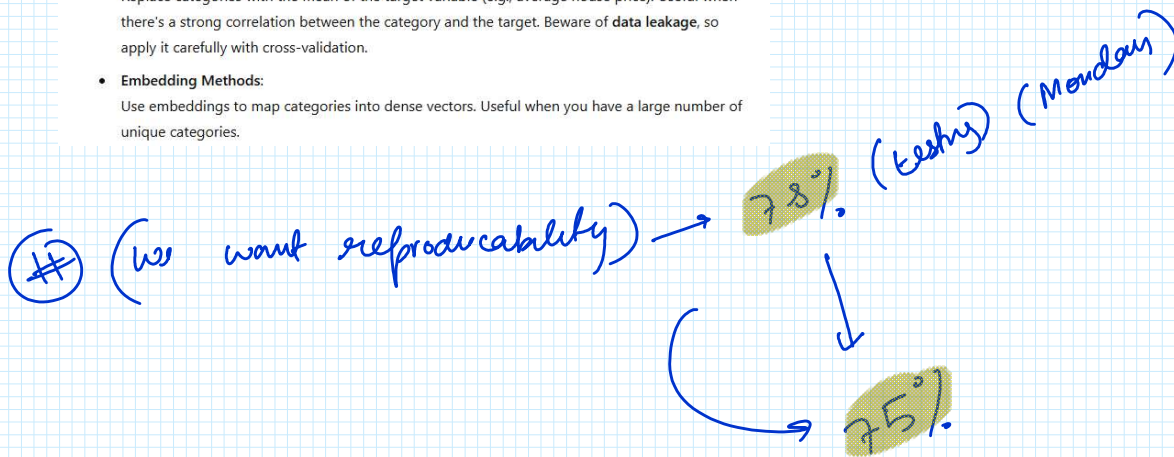
### 1. Nominal Features

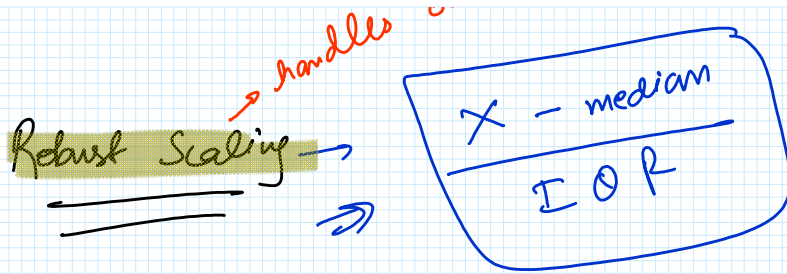Examples: City, House Style, or Neighborhood.

- **One-Hot Encoding:**
  Converts each category into a binary column. Suitable when the number of categories is small. Use libraries like `pandas.get_dummies()` or `OneHotEncoder` from `sklearn`.

- **Frequency Encoding:**
  Replace each category with the frequency of its occurrence. This helps reduce dimensionality compared to one-hot encoding.

- **Target Encoding:**
  Replace categories with the mean of the target variable (e.g., average house price). Useful when there's a strong correlation between the category and the target. Beware of **data leakage**, so apply it carefully with cross-validation.

- **Embedding Methods:**
  Use embeddings to map categories into dense vectors. Useful when you have a large number of unique categories.

### 2. Ordinal Features
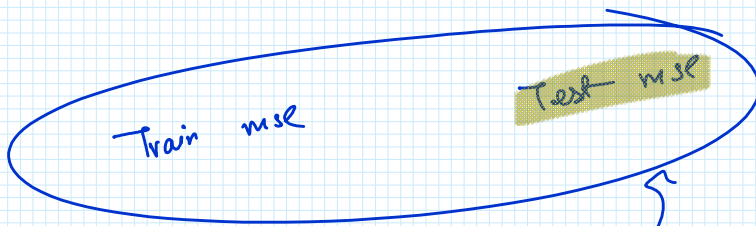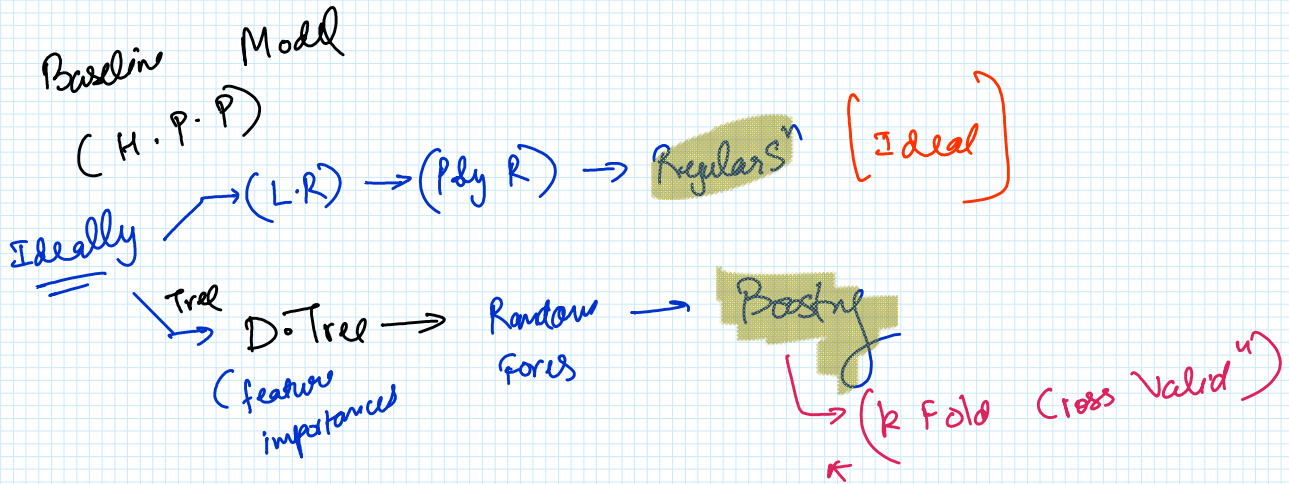
Examples: Quality Ratings (Low, Medium, High).

- **Label Encoding:**
  Assign integer values to categories based on their rank or logical order (e.g., Low=1, Medium=2, High=3). Use `LabelEncoder` from `sklearn`.
  Ensure the order is meaningful; otherwise, consider other encodings.

- **Map to Numeric Scales:**
  If the categories represent intervals (e.g., Bad=1, Good=3, Excellent=5), map them directly to these numeric values.
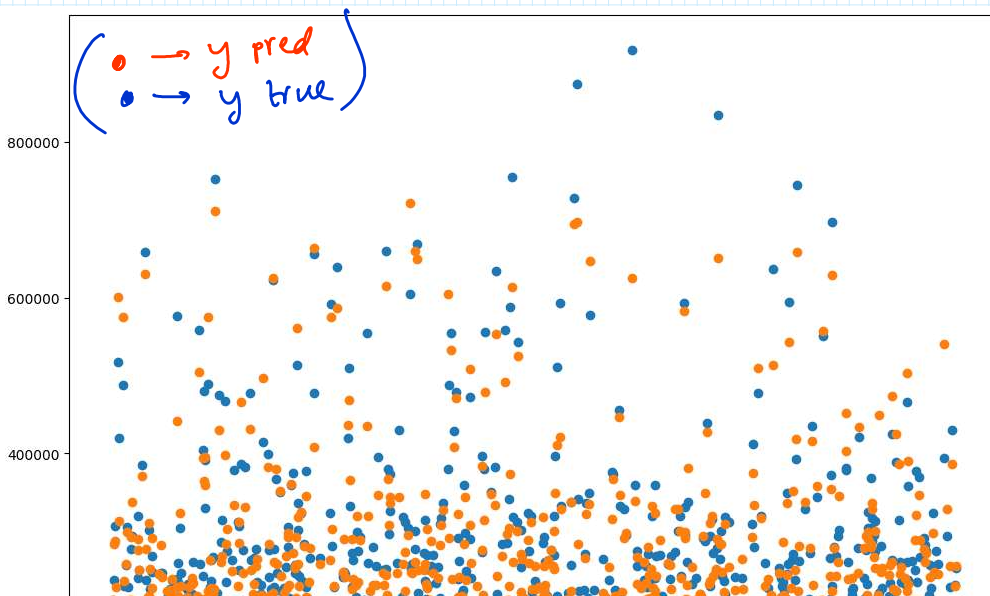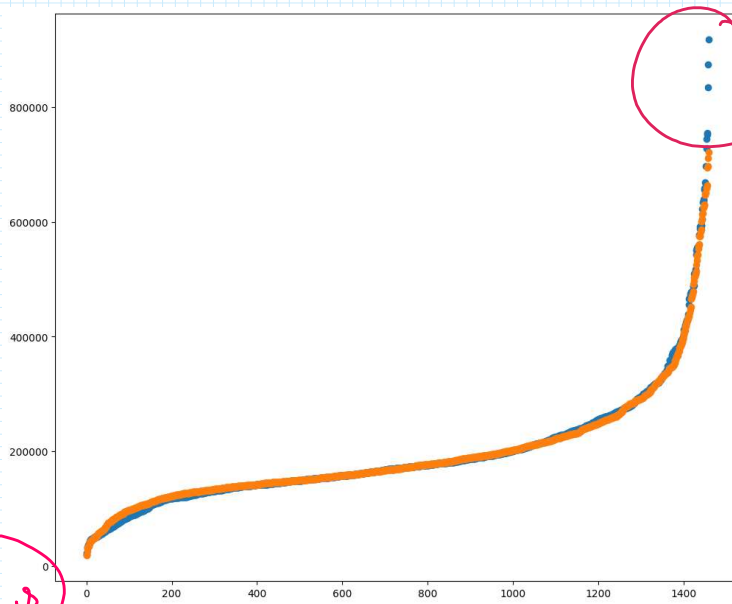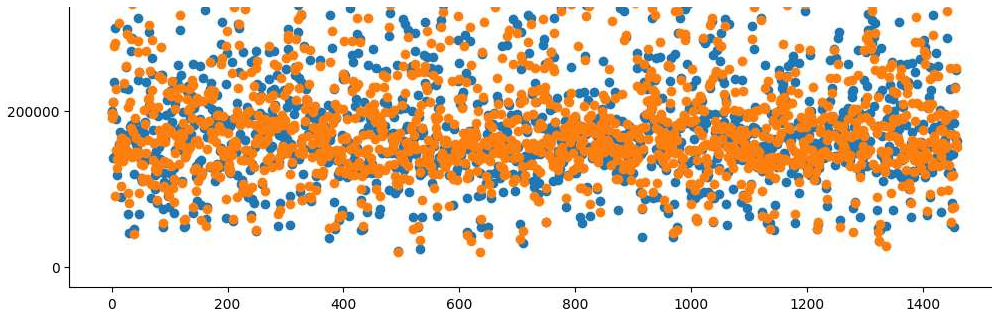
# (we want reproducability) → 75% (testing) (Monday)

↓

75%

---

Scaling    Necessary

Distance Based ( KNN, Clustery, SVM)

Gradient Based ( Log Reg, NN)

Regularis" ( Lasso, Ridge)

### 1. When to Scale?

← columns

Scaling is necessary when features have different units or ranges and you're using models that are sensitive to feature magnitudes,

Most Common

**Standard Scaling**

$$\frac{X - \mu}{\sigma}$$

( tries to maintain distrib")

# (gauss/Normal)

→ handles outliers

median

**Min - max Scaling**

$$\frac{X - X_{min}}{X_{max} - X_{min}}$$    [0, 1]

[0, 1]

→ Neural Networks
→ when you need all features in uniform scale

→ handles

features in scale ✓

Robust Scaling →

$$\frac{X - median}{IQR}$$

log transform$^n$

---

Baseline Model
(H.P.P)

Ideally → (L.R) → (Poly R) → Regular$^n$ [Ideal]

Tree → D.Tree → Random → Boosting
(feature importance) forest

→ (R Fold Cross Valid$^n$)

Train mse        Test mse

[Apply ML - 2 techniques)



( • → y pred
  • → y true )

800000

600000

400000

→ y true is high
but y pred is low

(current status)

Baseline Model ——→ good preds for cheap houses
but error for expensive houses is high
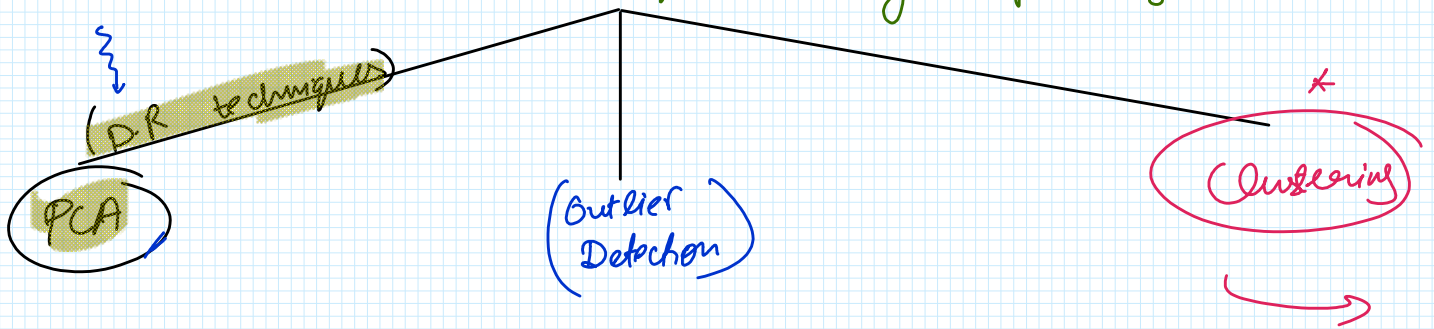
solution? → No more
Tr data for
expensive houses

→ Outlier
Removed

10:09 pm

ML-2 (use ML-2 techniques & try to improve the RMSE)
* (There is no guarantee)

What techniques do you know?

{ (D.R techniques)

PCA

(Outlier Detection)

* (Clustering)

---

PCA → * (Reduces dimension)

Reduce dim

error (Bug ↑)

(79) → 31,000 $

↓

35 → 33,000 $

---



explained variance

components

79

---

(Outlier Detection)

Z score
IQR (Box Plot)

(Outlier

→ IQR (box

PCA (n-comp = 2)

# (circled)

Visualisation → t-SNE (checked visually)
→ UMAP

# (circled)

→ (Read 2 research papers) → LOF
→ Isol$^n$ forest

global str of data



Industry
79 ──→ 2 → 2 dim
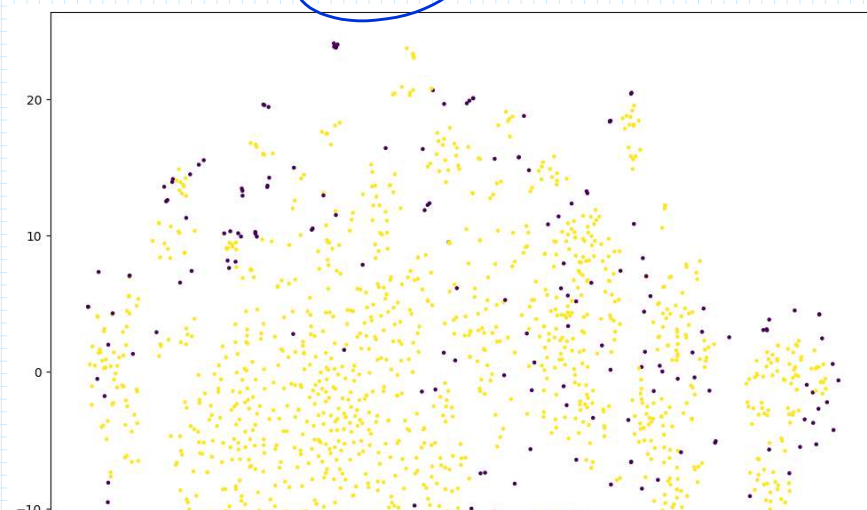
→ How much variance

Meat [30%]

t-SNE (circled)

LOF / I Forest *

I forest



LoF

Clustering ~~~→ Next Session

↳ (EDA on clustering → some more steps)

* 

Come prepared *

*(Not good)

(Dunn Index) & [Silhoutte Score]

(global)                    (data point level)
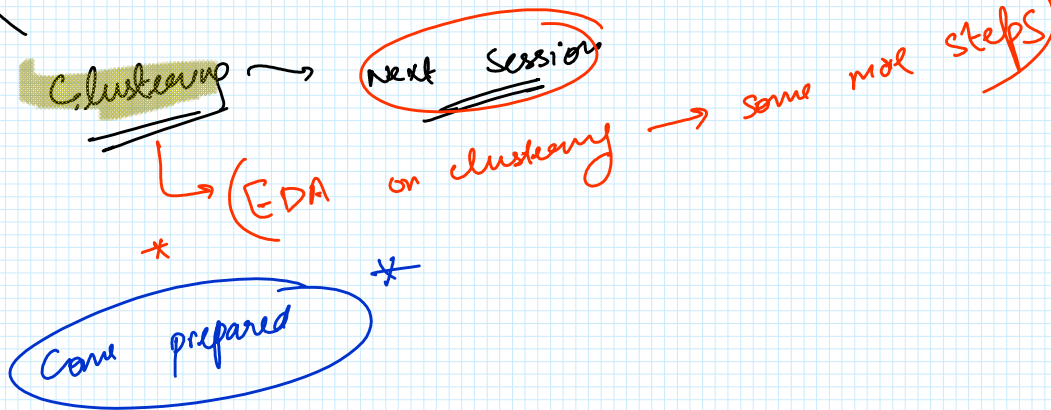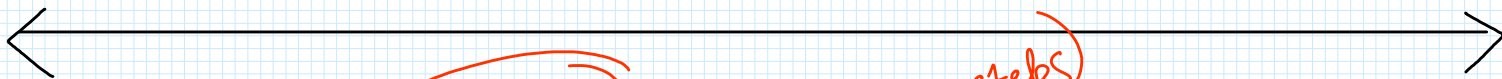
↳ on Index                 20 data pts
                           L.S. scores

**Dunn Index**
The Dunn Index evaluates the compactness and separation of clusters

$$D = \frac{\min (Inter-cluster)}{\max (Intra-cluster\ distance)}$$

→ Compact Cluster
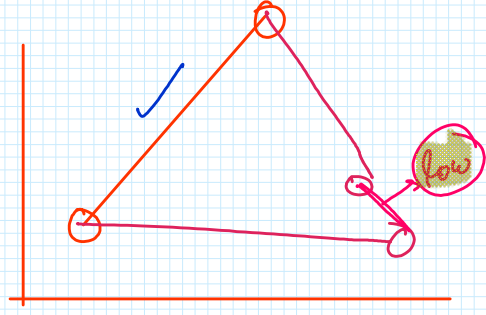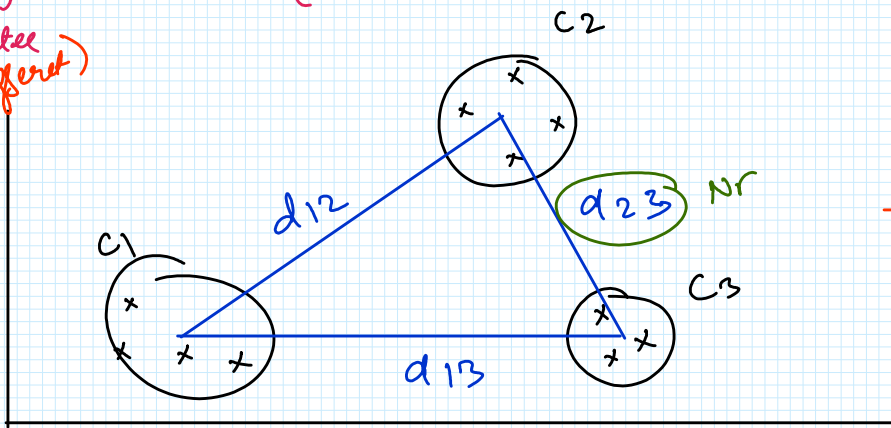
→ well separated clusters

$(> \frac{1}{a})$

$\left( \begin{matrix} >1 \\ >2 \end{matrix} \right)$

(good)

max (Inter cluster distance) → well clusters

Intra (same)

Inter (Different)



C2

$d_{12}$

$d_{23}$  Nr

C1

C3

$d_{13}$

low

---

Inter (Different)

b

C2

C1

C3

$\left. \begin{matrix} d_1 \\ d_2 \\ \vdots \\ d_{10} \end{matrix} \right\}$
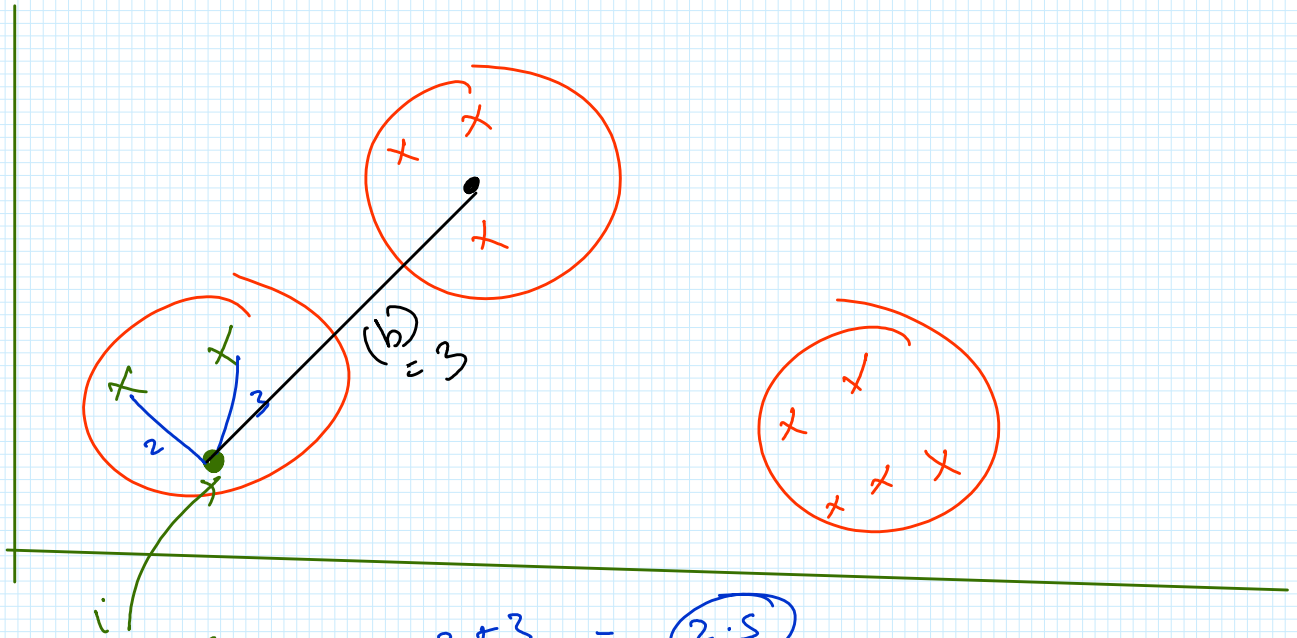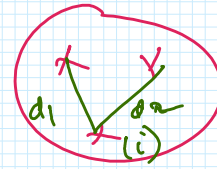
---

Silhouette Score    $(-1, 1)$

$$S(i) = \frac{b(i) - a(i)}{max(b, a)}$$

b = inter cluster distance of that pt i to nearest cluster

$a$ = Intra cluster distance

$$a = \frac{a_1 + a_2}{2}$$



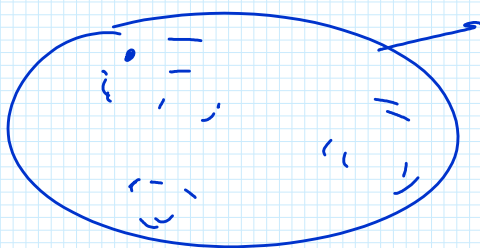$(b) = 3$

$(a) \rightarrow \dfrac{2+3}{2} = \boxed{2.5}$

$(b) \rightarrow 3$

$$S.S(i) = \frac{b-a}{\max(a,b)} = \frac{3-2.5}{3} = \frac{0.5}{3}$$

$[-1, 1]$

$S(i)$: Silhouette score for point $i$.

- $S(i) \approx 1$: Point is well-matched to its own cluster and far from other clusters.
- $S(i) \approx 0$: Point is on the boundary between two clusters.
- $S(i) < 0$: Point is closer to another cluster than its own.

$\rightarrow$ final S.S = avg ( all S.S for $i$ data pts )

$k = 3$

D.I

(high)

$k = 5$

D.I