

DAV-3

HYPOTHESIS TESTING

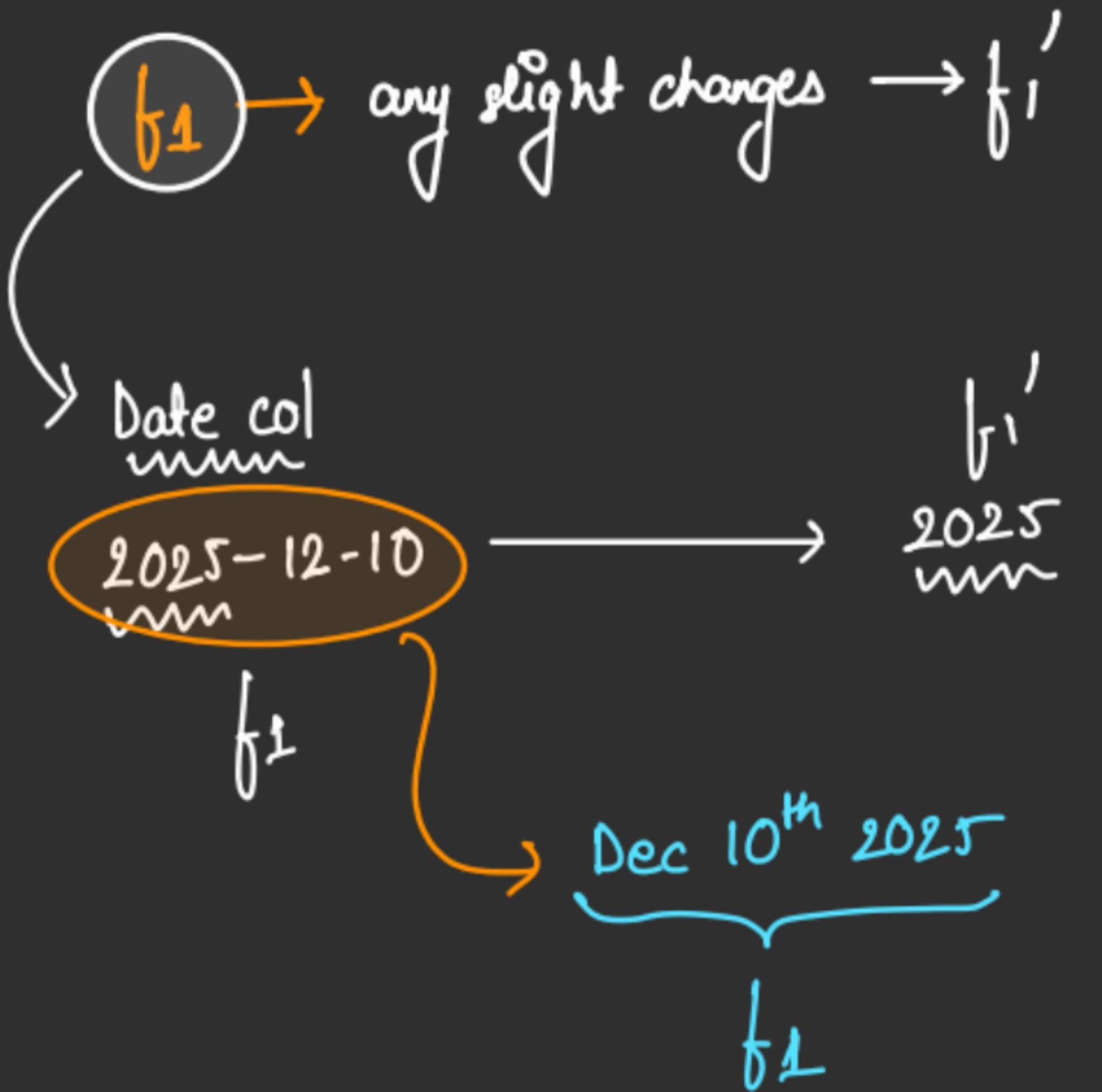
(Class starts
@ 9:10 PM)



Lecture 10: Feature Engineering 2

#Agenda

- ① Creating New Features
 - ② Missing Values
 - ③ Outlier Treatment
 - ④ Categorical Encoding & Types → TE left
 - ⑤ Normalization & Standardization
- ↓ Next Class



f_1 wrt Target \rightarrow Univariate

b_1, b_2 wrt Target \rightarrow Bivariate

Missing Values

① Identifying Null Values

↳ $\{ \text{isnull().sum()}\}$
 $\{ \text{isna().sum()}\}$



C_1	C_2
1	x
2	y
Null	z
3	a
4	b

② Why to fix NULL Values?

→ Data Loss
→ ML model don't understand
Null values

③ How to fix them?

→ Statistical way → Mean, Median
Mode, Max, Min

→ Library Based

Column
~~~~~

Numerical  
~~~~~

Mean

Median

Mode

Constant Value

Categorical
~~~~~

Mode

Constant Value

Num & has outliers → Median

Num & has almost No outliers → Mean

Cat → Mode

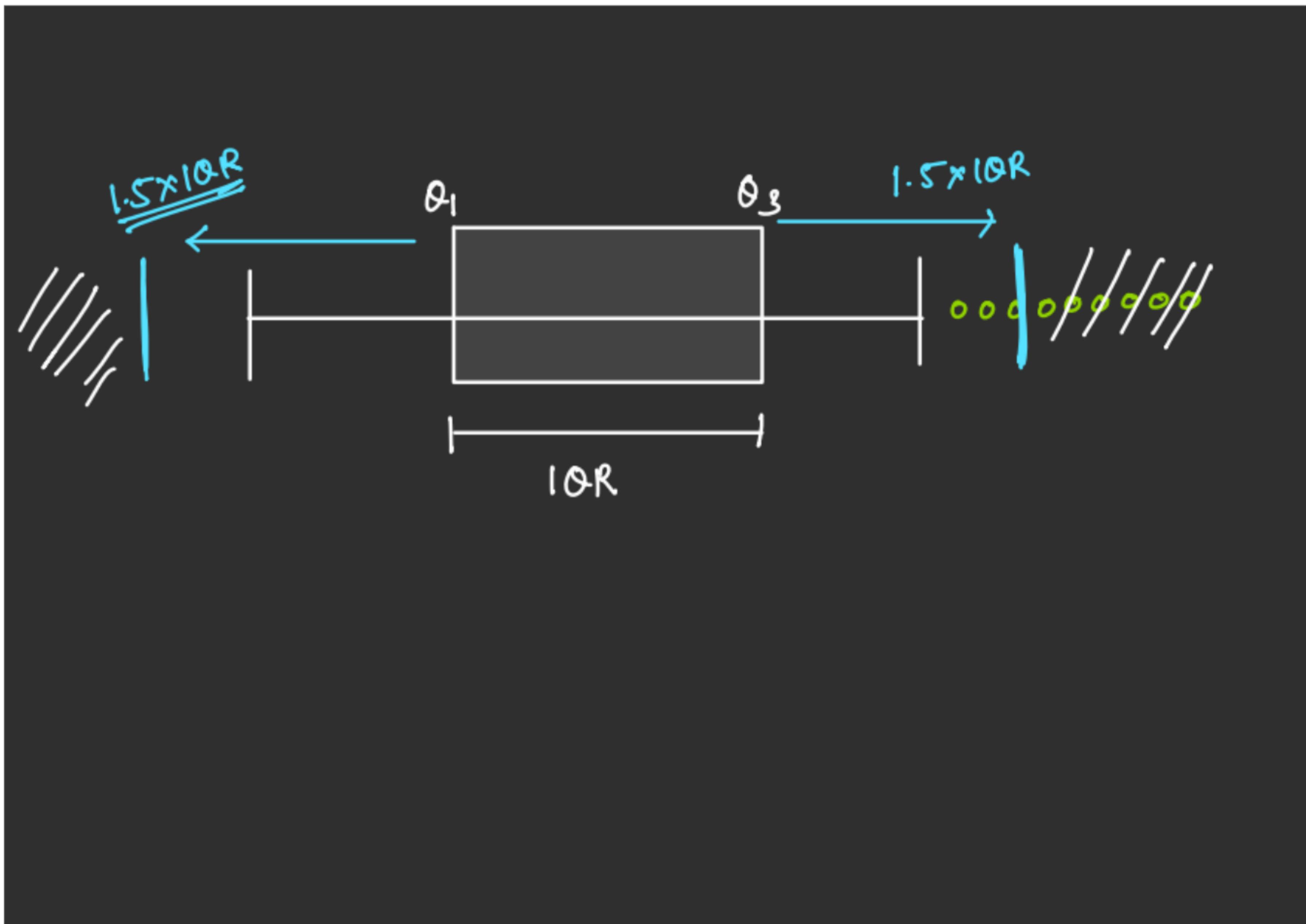
→ Const Value

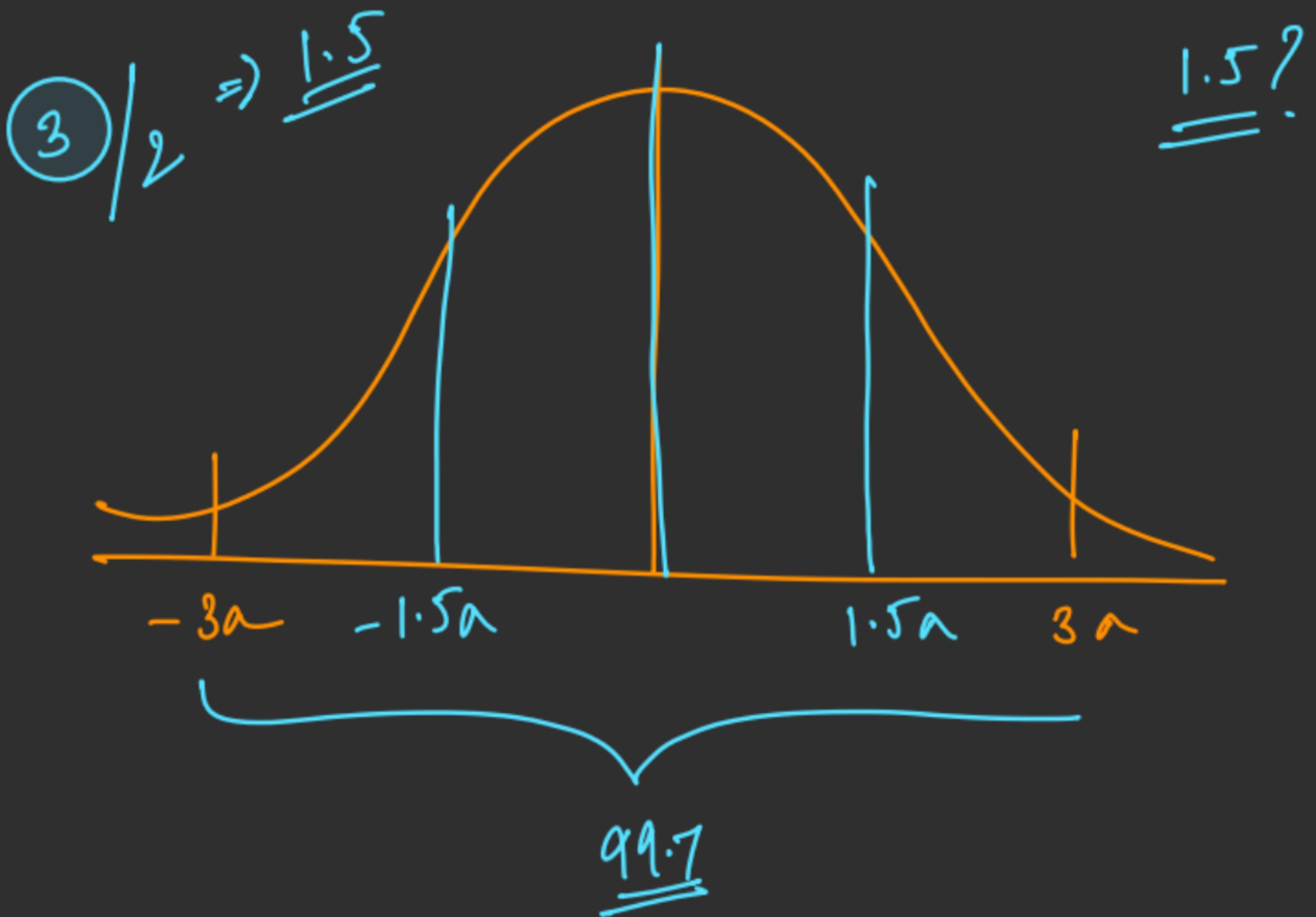
Num & you want to fill with your own value → Const Value

## #Outlier Treatment

- o) How to identify outlier? (Check the distribution → Boxplot  
→ Histogram)
- o) What are some ways to remove outliers? (IQR, Z Score)







## #Categorical Encoding

- ① What is categorical Encoding? [Convert Cat to Num]
- ② Why do we need this? [ML model don't understand Cat Variable]
- ③ What are the Types? [There are 3 Types]

## # Types of Encoding

- ① One hot Encoding (OHE)
- ② Label Encoding
- ③ Target Encoding.

# ① One Hot Encoding (OHE)

| Gender | M | F |
|--------|---|---|
| Male   | 1 | 0 |
| Female | 0 | 1 |
| Male   | 1 | 0 |
| Female | 0 | 1 |
| Female | 0 | 1 |

When  $\text{cat} = 2$

OHE is good

When  $\text{cat} = 5$

Sparse data

Majority of your data has zeros

| Travel | Low | Medium | High |
|--------|-----|--------|------|
| Low    | 1   | 0      | 0    |
| Medium | 0   | 1      | 0    |
| High   | 0   | 0      | 1    |

## ② Label Encoding

| Gender | New col |
|--------|---------|
| Male   | 1       |
| Female | 0       |
| Male   | 1       |
| Female | 0       |

Categorical → Animal

| Animal   | Value |
|----------|-------|
| Cat      | 1     |
| dog      | 2     |
| elephant | 3     |
| Rat      | 5     |
| Fox      | 4     |

Nominal data

We use Label Encoding when we have ordinal data



→ Data is ordered

Nominal Data

No Order

Example

High → 3

Medium → 2

Low → 1

### ③ Target Encoding

- When  $\text{cat} = 2$  then use OHE
- When  $\text{cat} > 2 \rightarrow$  Label Encoding (Ordinal data)
- When  $\text{cat} > 2 \rightarrow$  Target Encoding (Nominal data)

→ Loan Status (The column you are predicting)

| Name   | Target |
|--------|--------|
| Akash  | 1      |
| Max    | 0      |
| ✓ John | 1      |
| Max    | 1      |
| ✓ John | 1      |

⇒

$$\text{Akash} \rightarrow 1 \rightarrow \frac{1}{1} \Rightarrow 1 \checkmark$$

$$\text{Max} \rightarrow 2 \rightarrow \frac{0+1}{2} \Rightarrow 0.5 \checkmark$$

$$\text{John} \rightarrow 2 \rightarrow \frac{1+1}{2} \Rightarrow \frac{2}{2} = 1$$

## # Column Normalization & Standardization

- ① Why do we need this? (To bring every feature into same scale)
- ② What are Types?

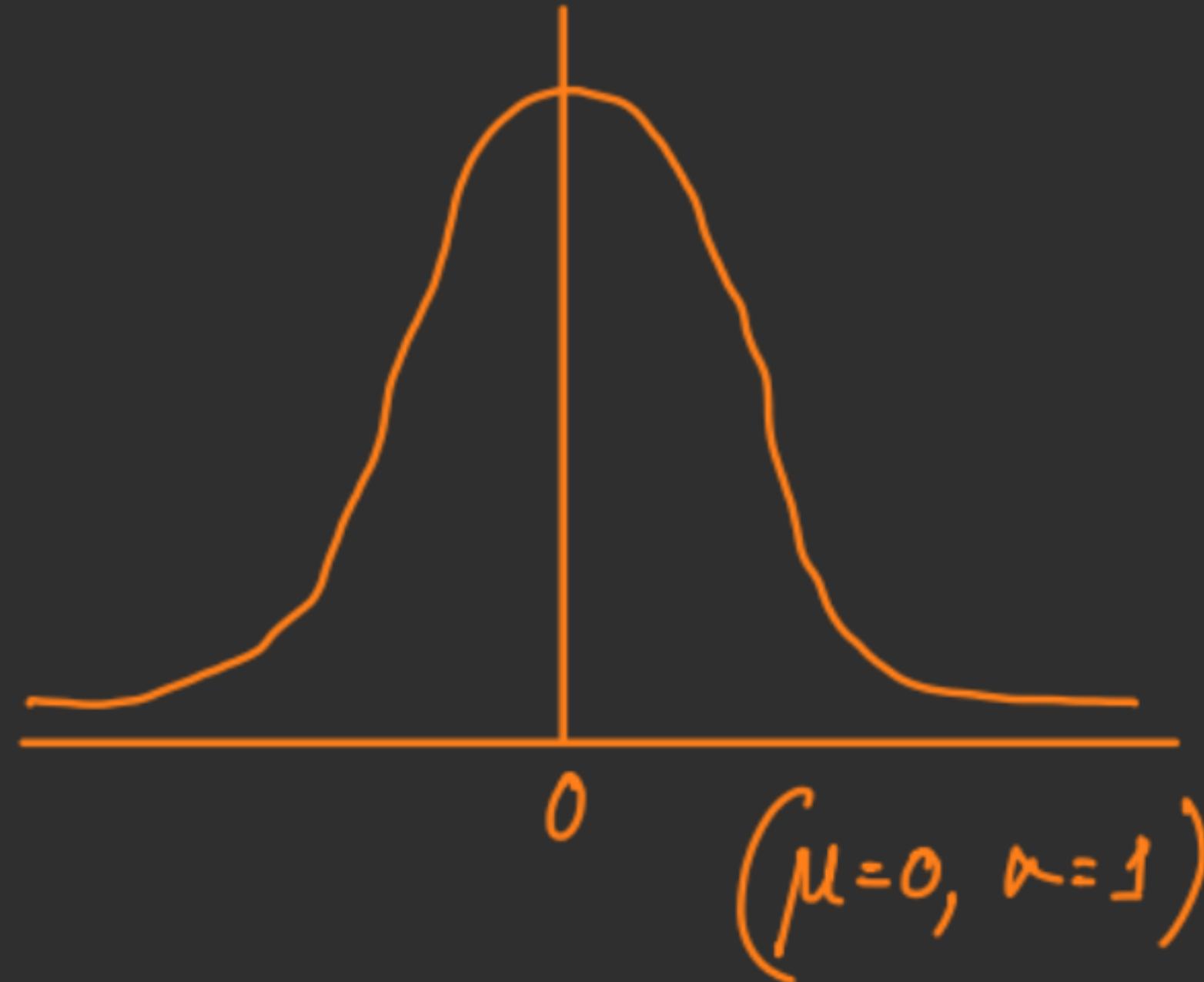


Normalization & Standardization

# Standardization  
~~~~~

"Standard Scaler" =
$$\left(\frac{x_i - x.\text{mean}())}{\text{std}(x)} \right)$$

→ Standardization
Library



Normalization

“Min Max Scaler” =
$$\left(\frac{x - x_{\min}}{x_{\max} - x_{\min}} \right)$$

Normalization
library

Range $\in [0,1]$

Normalization and standardization are both used to rescale data, but they serve different purposes and are best suited for different situations:

1. **Normalization** (scaling to a range, typically 0 to 1) is ideal when:

- You know the data has no extreme outliers, and values are already bounded or have a specific range.-
- You're working with algorithms that assume data within a fixed scale (e.g., neural networks or k-nearest neighbors).

2. **Standardization** (scaling to a mean of 0 and standard deviation of 1) is ideal when:

- Your data has a normal (bell-shaped) distribution or contains outliers.
- You're using algorithms that assume data has a Gaussian distribution (e.g., linear regression or principal component analysis).

In short, normalize when you want all values on the same scale without outliers affecting the range, and standardize when you want to center and scale data, especially when the data might have extreme values.