

* Data Preprocessing

Interview Question

~~completed~~ ①

Handling missing value - Measure of central

Encoding concept

Tendency

✓ ②

Handling

Char | Object

variable

Label Encoder

One-hot-encoder

dummy variable

③

Handling outlier

- Transformation approach

log, sqrt, cuberoot etc.

very

removal / Trimmer

capping -

✓ ④

Feature scaling

Normalization

Standardization → by default

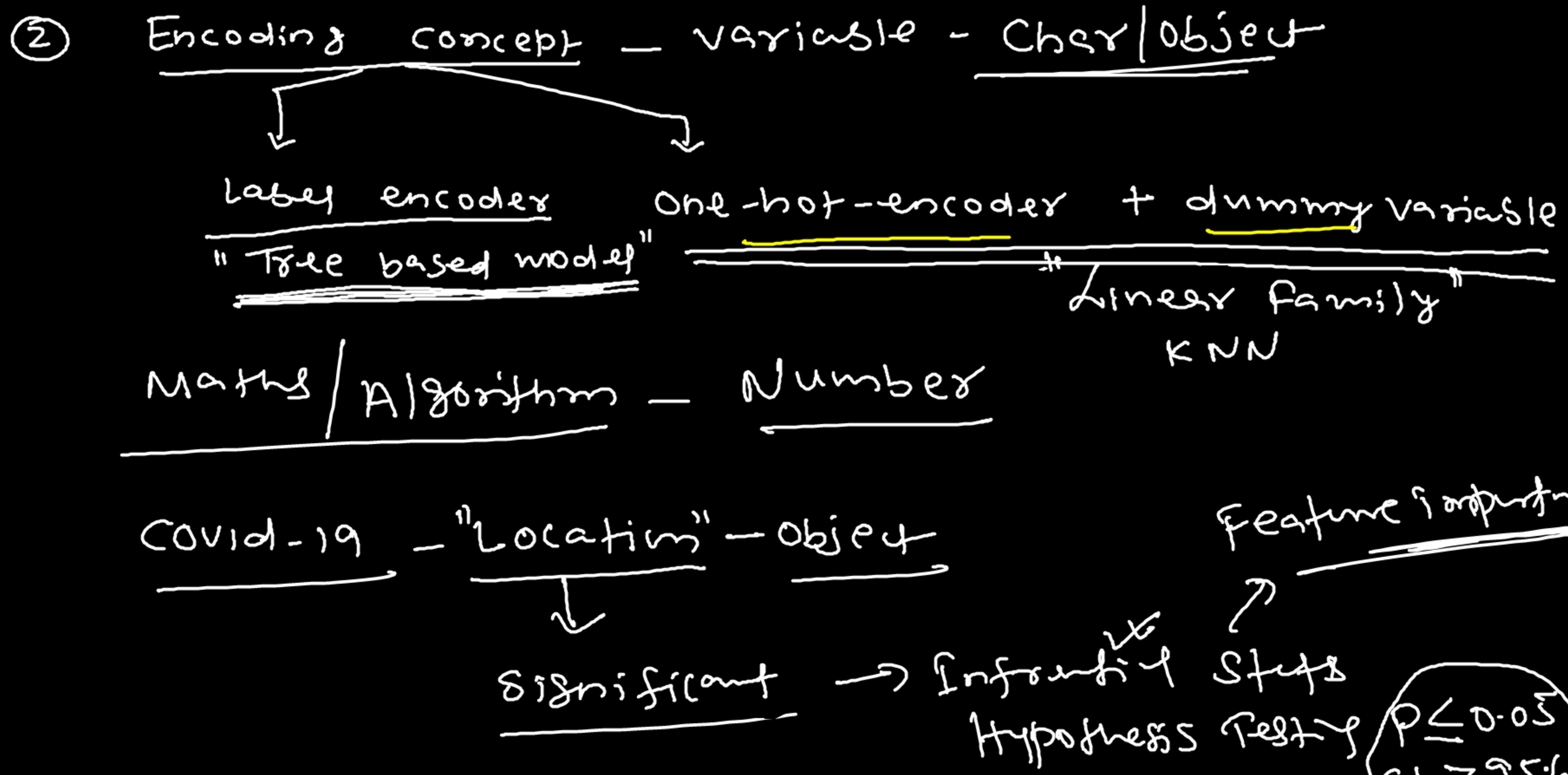
IQR

z-score

✓ ⑤

Imbalance Treatment

→ This is only applicable with Classification problem



You are screen sharing

Stop Share

Location	Label Encoder
Delhi	3
Mumbai	5
Pune	6
Bangalore	0
Chennai	2
Gujrat	4
U. P.	7
Bihar	1

Dummy Variable

$$= n - 1$$

(Total no. of column - 1)

benefit "Multicollinearity"

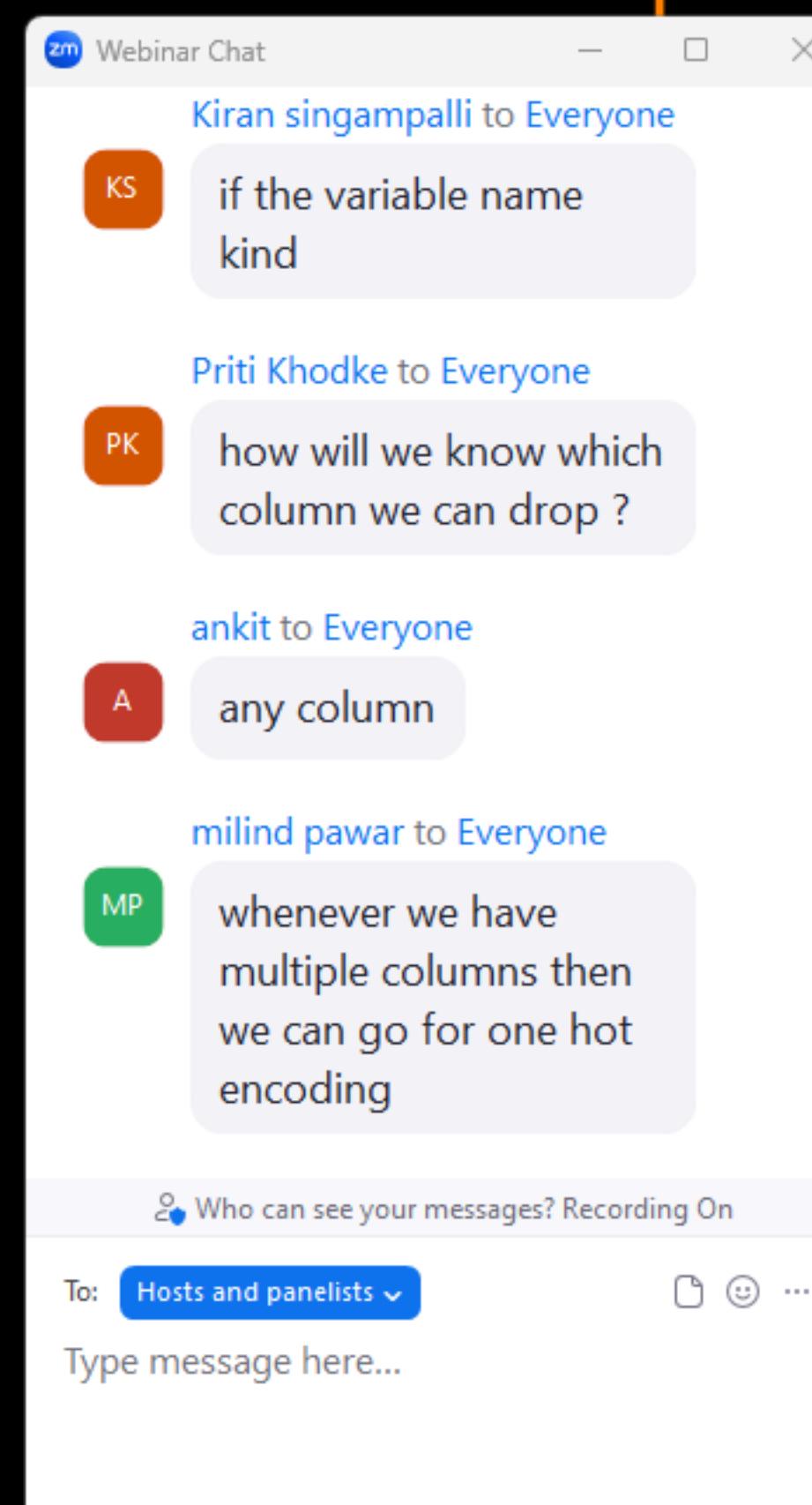
Linear Regression

One-hot-encoder

Gender	Label Encoder	One-hot-encoder			Dummy variable	
		G_F - Female	G_M - Male	G_T - Trans	G_M	G_T
Male	1	0	1	0	1	0
Female	0	1	0	0	0	0
Female	0	1	0	0	0	0
Male	1	0	1	0	1	0
Transgender	2	0	0	1	0	1
M	1	0	1	0	1	0
F	0	1	0	0	0	0
T	2	0	0	1	0	1

You are screen sharing

Stop Share



Feature Scaling

- Algorithm

Age	Salary
20	50K - 50000
25	60K
30	80K
35	110K
45	150K
55	200K
60	250K - 250,000

Linear Regression \rightarrow PdV methods

$$Y = m_1 x_1 + m_2 x_2 + C$$

↓ ↓ ↓ ↓
Dep Slope Slope Intercept

$$Y = m_1 * \text{Age} + m_2 * \text{Salary}$$

$$m_1 \text{ & } m_2 = 1$$

$$Y = \text{Age} + \text{Salary} * 2$$

You are screen sharing

EN

Stop Share

Range

where variable's range
if known

Feature Scaling

through the value

Normalization

0 1

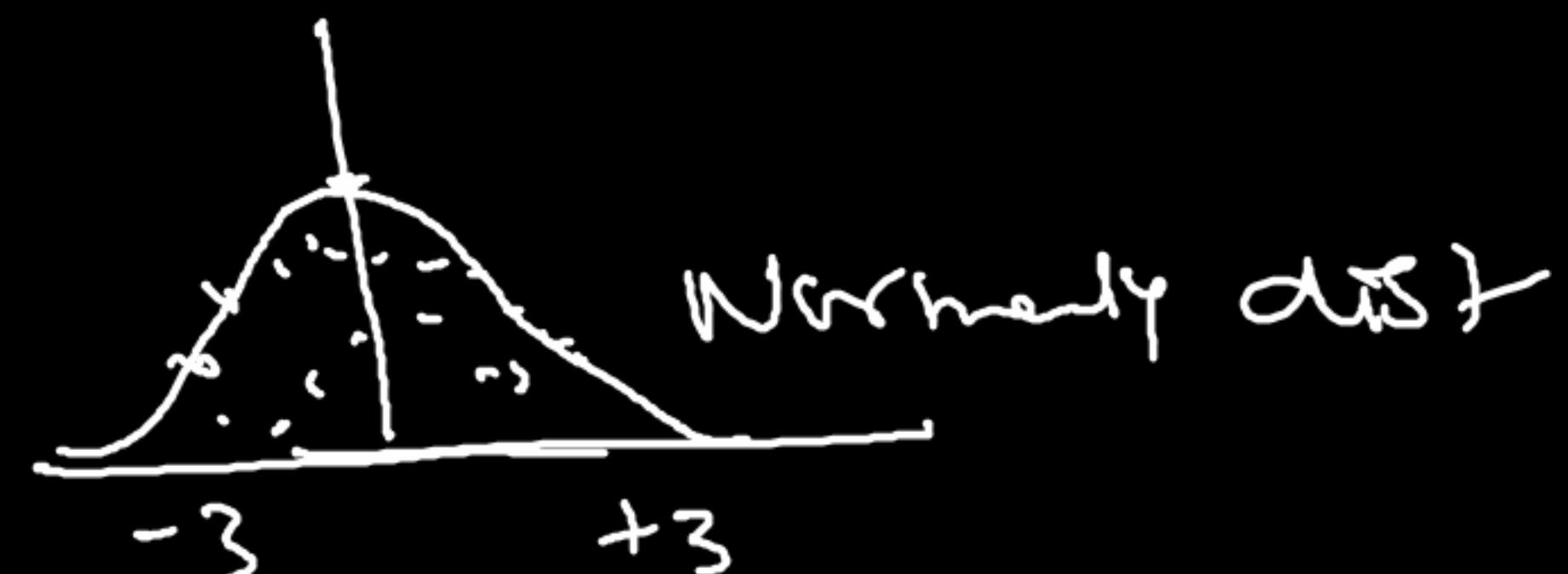
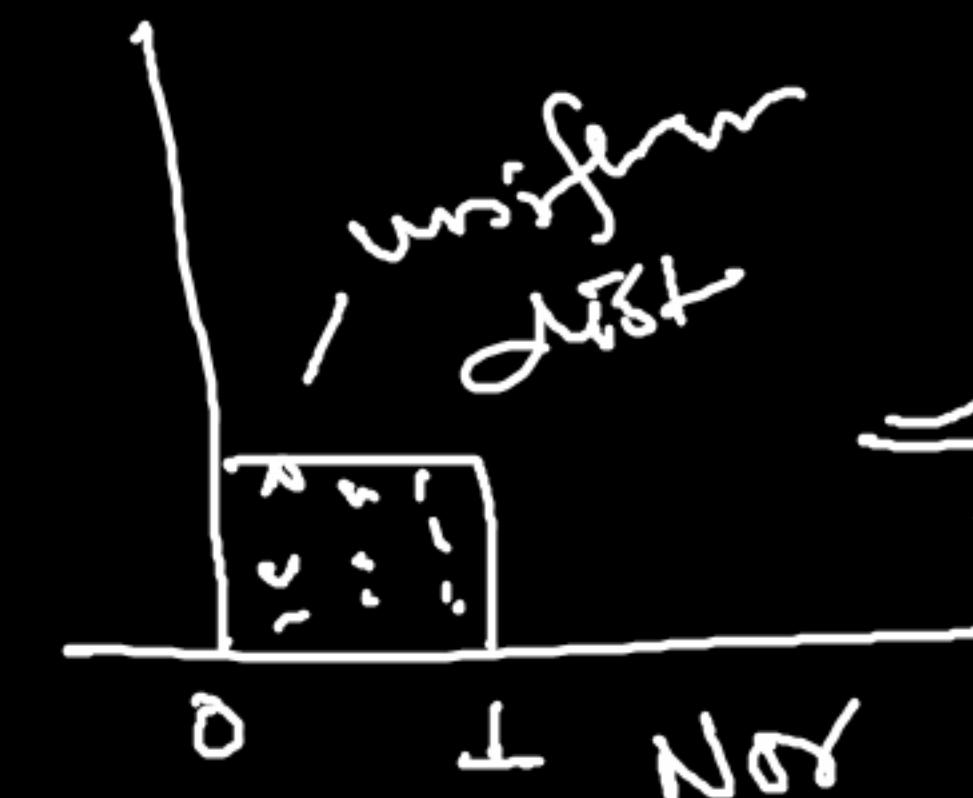
$$Nor = \frac{x_i - \text{min_value}}{\text{max_value} - \text{min_value}}$$

NO outliers

Standardization

there is an outlier

$$Std = \frac{x_i - \text{mean}(x)}{\text{std}(\sigma)}$$



Std