

→ Handle missing values

→ Outliers (Distribution of Data)

Symmetric &

Relationship b/w

mean, median & mode

skewed

Categorical data

Outliers

mode

↓

→ Descriptive Statistics

Median

(More robust to
the outliers)

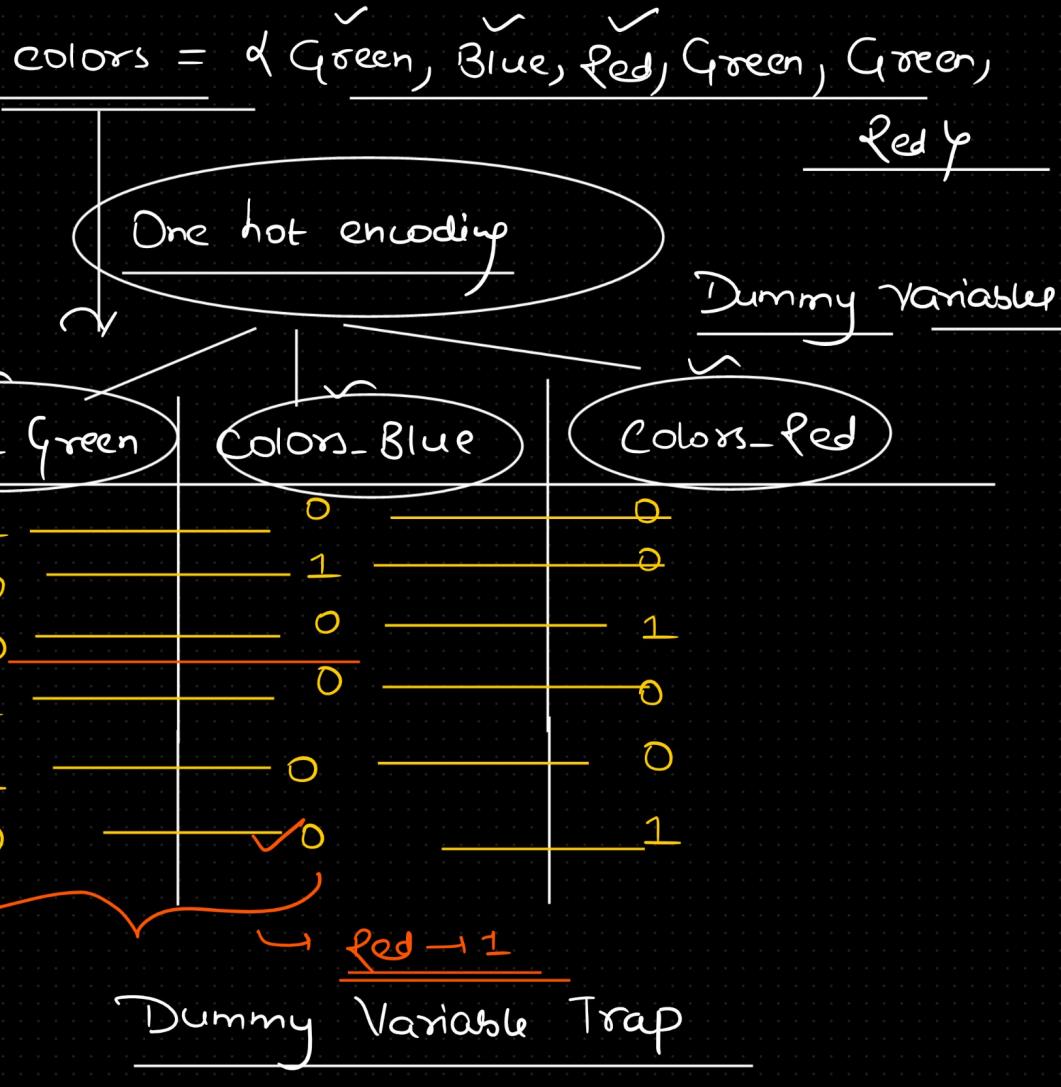
Handle Categorical data

→ numeric form (One hot encoding)

8

Label encoding)

Encoding
Technique



No multicollinearity

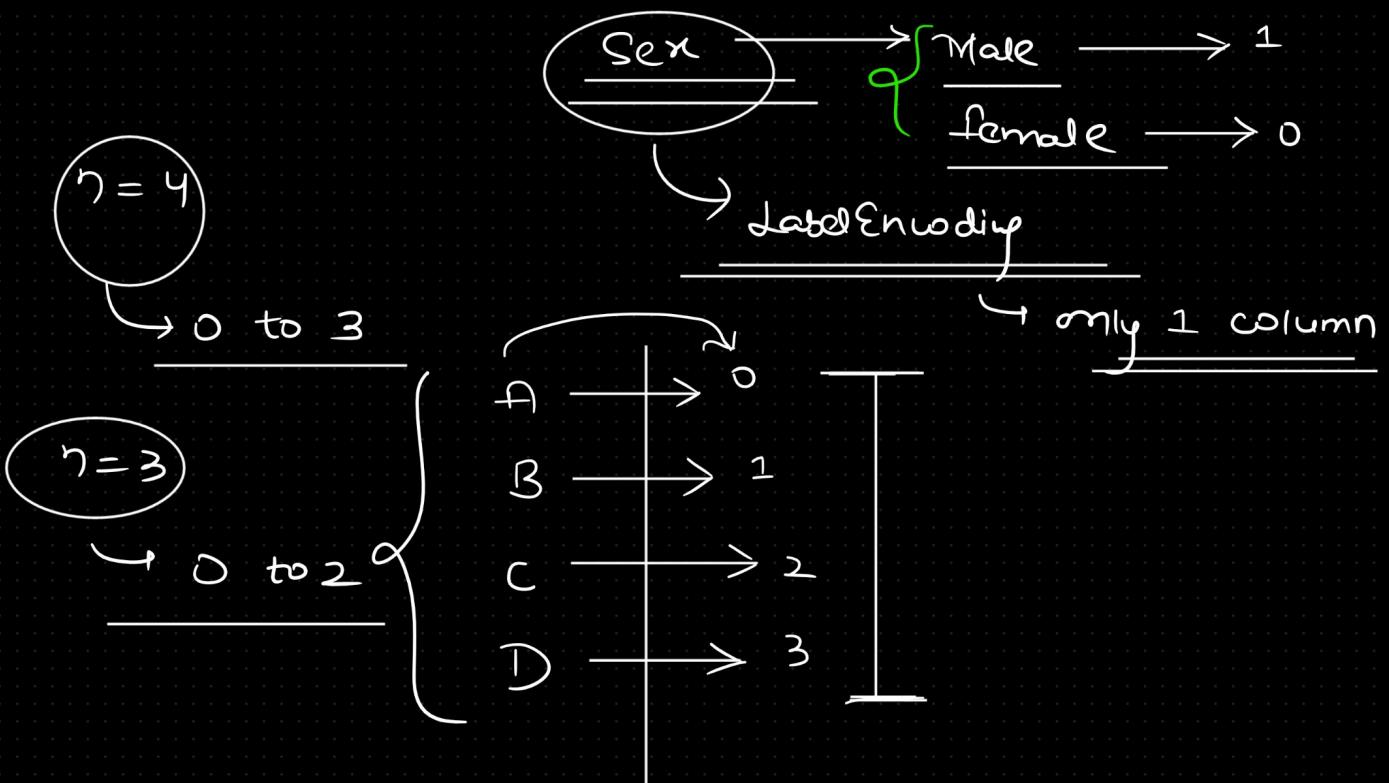
drop-first = True

→ avoid the issue

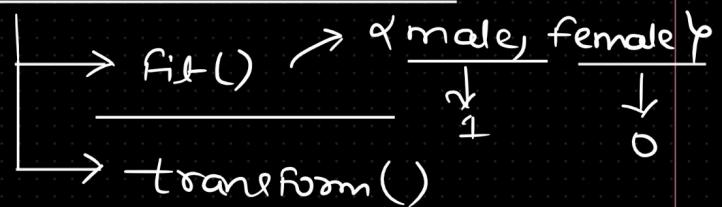
✓ dummy

variable trap

~~Sex-Male~~ | Sex-female
1 - female
0 - male



fit_transform() mapping

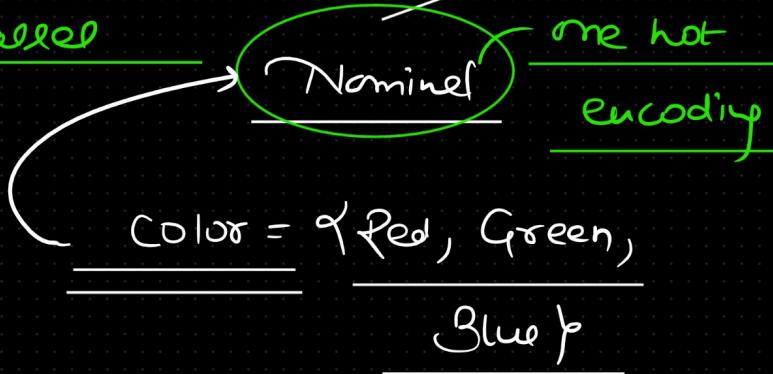


complete
feature
'Sex' will
be transformed
as per the
above mapping

Types of Categorical data

more than 2

classes



Label Encoding

Some sort of order

(Priority)

Embaraked = {S, C, Q}

Education = {B-Tech,
M-Tech, PhD}

Ranking = {1, 2, 3, ... 25}

1
I
II
III

9

Correlation Coefficient (Heatmap)

feature

$\rightarrow -1 \text{ to } 1$

Selection

$\curvearrowleft \rightarrow$ Curse of Dimensionality

$f_1 \ f_L \ f_3 \ f_4 \ f_5 \ \dots \ f_{300}$

300 Dimensions

Subset of

Complex

the

above

features

\rightarrow feature Selection

Capture

the important

info from all

the above features

& recreate new

features all together

}

\rightarrow PCA, tSNE

Negative
Correlation

$$\uparrow f_1 \quad f_2 \downarrow$$

near to -1

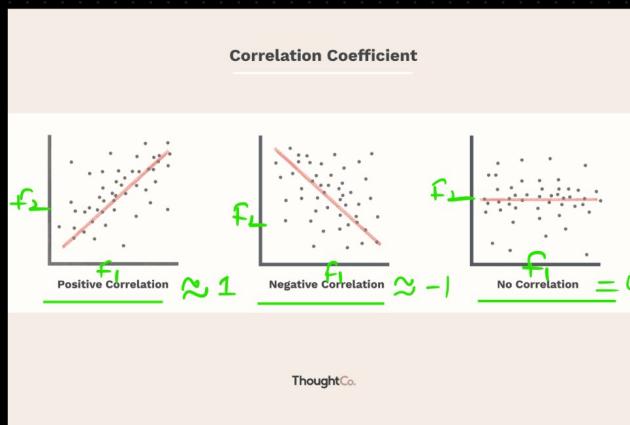
Redundant

information

Positive
Correlation

$$f_1 \uparrow \quad f_2 \uparrow \quad r = 0.98$$

near to +1



> 85% similar/redundant info

$$r = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{[n\Sigma x^2 - (\Sigma x)^2][n\Sigma y^2 - (\Sigma y)^2]}}$$

numeric values

$x \rightarrow f_1$

$y \rightarrow f_2$

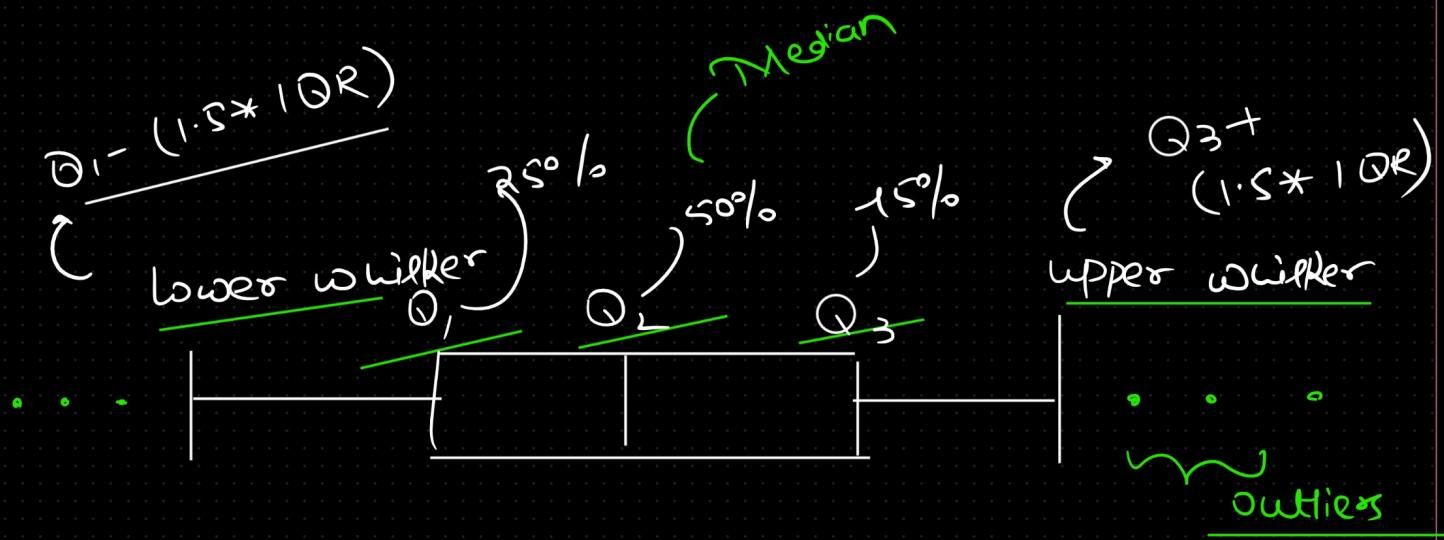


Outliers Detection

Box Plot

$$\text{Range} = \text{Max} - \text{Min}$$

$$\text{IQR} = Q_3 - Q_1$$



0 1 2 3 4 5 6 7 8 9 10 11
 (5, 40, 42, 46, 48, 49, 50, 52, 53, 55, 56, 56)
 ↑
 75% 100%
 12 13

Median → Sort the data

(Ascending order)

$$50\% \quad Q_2 = \frac{50 + 52}{2} = 51$$

$$25\% \quad Q_1 = 46$$

$$75\% \quad Q_3 = 56$$

$$IQR = Q_3 - Q_1 = 56 - 46 = 10$$

$$\text{Lower whisker} = 46 - (1.5 \times 10) \quad \text{Research}$$

$$\Rightarrow 46 - 15 \Rightarrow 31$$

$$\text{Upper whisker} = 56 + (1.5 \times 10)$$

$$\Rightarrow 56 + 15 \Rightarrow 71$$

Lower whisker

✓
.



Upper whisker ✓

.

Grouping & Aggregation

↳ SQL
