# Encoding - Theory

- 1) Label encoding
- 2) One - Hot Encoding ( pd. dummy) [0 or 1]
- 3) Ordinal Encoding

⤷ when Cat feat is not ordinal, The no. of categorical feat is less

↳ when cat feat is ordinal, The no. of categorical feat is quite large.

① Label Encoding (each label/class is assigned a
   ↓                unique integer based on
→ | alp order NOT |   alphabetical ordering)
   | order of Data |

→ from sklearn.preprocessing Import Label Encoder

→ le = Label Encoder ( )        ⎡ Drawback →     ⎤
→ Y = le.fit_transform (y) .    ⎢ False order of ⎥
                                ⎣ Data           ⎦

② One - Hot Encoding [ creats diff col for classes of
                                    a feat, and 1 → if it
Drawback → ↑ in Dimensionality              appears and
                                            0 → if not ]
————————— . × ————————— . × ——————— . × —————

* Tree-Based algos → Can work with ⎡ Cat Variables    ⎤
                                   ⎣ and Label Encoding⎦

• LR, Distance metric ( K-m, KNN) or ANN
                        ↳ works with OH
                                Spiral Encoding

= one-df = pd.get_dummies (y, drop-first = True)
→ df = pd.concat ([df, one-df], axis=1)
→ df = df.drop ([' '], axis=1).

3)• <u>ordinal Encoding</u> → as per order of Data [Label]

→ <u>Cat feat</u> → ordinal num value [ordered set].

| → eg → | | | |
|---|---|---|---|
| poor | 1 | Not alphabetic order |
| Good | 2 | But order of Data. |
| Very Good | 3 | |
| Excellent | 4 | |

# creating a dict for mapping (with values)

dt-dict = { 'poor': 1, 'Good': 2, 'Very Good': 3,
                    'Excellent' : 4 }

df ['Colomn'] = df ['Customer Rat'].map (df-dict).

———— X ———— X ———— X ————

Feature Scalling [After split on
[value range 0-1]          training Data]

→ Normalization      and   Standardazation.
    [min-max]                    [z-score]
[value end up B/w 0 and 1]   [centered around
                                         mean with value of
                                         std ]

$$\frac{X - X_{min}}{X_{max} - X_{min}}$$

$$\frac{X - \mu}{\sigma}$$   L Kitna std
                                    deviation
                                    from mean
                                    Spiral

# fit → transform

Normalize → when we know that Distribution
of our data does not follow a Gaussian
distribution    [ KNN and NN]

Standardization → where Distribution of our
data follows Gaussian
distribution

## # filter the numeric

df_num = df.select_dtypes (include = np.number)

### # normalization
from sklearn.preprocessing import minMaxScaler
norm = minMaxScaler (). fit (data _num)
data_num_norm = norm. transform [ ]  └ here this or
feature you
want to
normalize.

### # Standardization
from sklearn.preprocessing import StandardScaler
scale = StandardScaler (). fit (data _num)
data _num_ scale = scale. transform (data-num)

——————×——————×——————×——————

* **Missing values**

① [Check assigned Data types]
    └ if wrong → [change using astype ()]
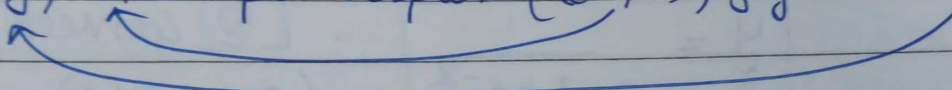
② Standard miss value → detect By Python (NaN)

## Outliers → detect using Boxplot

```
df_num = df.select_dtypes(include = np.number)
fig, ax = plt.subplots([2,2], figsize=(15,8)]
```

```
for var, subplot in zip(df.columns, ax.flatten()):
    z = sns.boxplot(x= df[var], orient= 'h', ax=subplot)
    z.set_xlabel(var, fontsize=20).
```

| ① To drop them | ② To treat Them |
|---|---|
| IQR method | Log transformation |
| Z-Slore (>,< -3) | Quantile based flooring and capping |

Measuring SLR, MLR, PR

1) MSE values
2) $R^2$ value

| | SLR → Residual plot
| | MLR → dist plot
| | PR → poly func plot

↓ MSE val    and ↑ $R^2$ val ($R^2$ → coeff of Determinant)

↳ sum of squares due to Regression

∴ SST = SSR + SSE

↓                    ↳ sum of squares due to
Total sum of squares                    error.

• $r^2 = \dfrac{SSR}{SST}$   or →   $\dfrac{SST}{SST} = \dfrac{SSR}{SST} + \dfrac{SSE}{SST}$

$1 = r^2 + \dfrac{SSE}{SST} \rightarrow r^2 = 1 - \dfrac{SS\,reg}{SS\,tot}$

• Sample Correlation Coeff

$r_{xy} = (sign\ of\ b_1)\sqrt{coeff\ of\ Determinant}$

$r_{xy} = \boxed{(sign\ of\ b_1)\sqrt{r^2}}$

——— x ——— x ——— x ———

$2^{nd}$ way of calc $\boxed{R^2}$ →

$R^2 = 1 - \left[\dfrac{MSE\ of\ reg\ line}{MSE\ of\ avg\ of\ data}\right]$

$\boxed{\dfrac{MSE}{\dfrac{\Sigma(y - \hat{y})^2}{n}}}$

$\boxed{1 - \dfrac{\Sigma(y_i - \hat{y}_i)^2}{\Sigma(y_i - \bar{y}_\bullet)^2}}$

Spiral

# Reg performance

→

• For SLR → <u>MSE value</u> + <u>$R^2$ score</u>

For MLR → <u>MSE value</u> + <u>Adj $R^2$ score</u>

Adj $R^2$ score ⇒ $1 - \left[\dfrac{(1-R^2)(u-1)}{N-p-1}\right]$

$N$ = Sample size
$p$ = No. of <u>predictors</u>.

---

* Every time we
add a new feat, $R^2 \uparrow$

(even though the feat
has low corr)

• But adj $R^2 \uparrow$ only when
feat has high corr

$\boxed{\text{Adj } R^2 \leq R^2}$ .

for var in df.unique
$\sum_{i=1}$
df-dict = { var: $c_i$ }
$i = i+1$;