→ → Regression

⇒ Classification

Clustering

MLA₁ → M₁ → P₁

MLA₂ → M₂ → P₂

MLA₂ → M₃ → P₃

Champion - Challenger Model

Raw Data → Pre-processing → Prepared Data

Scaling

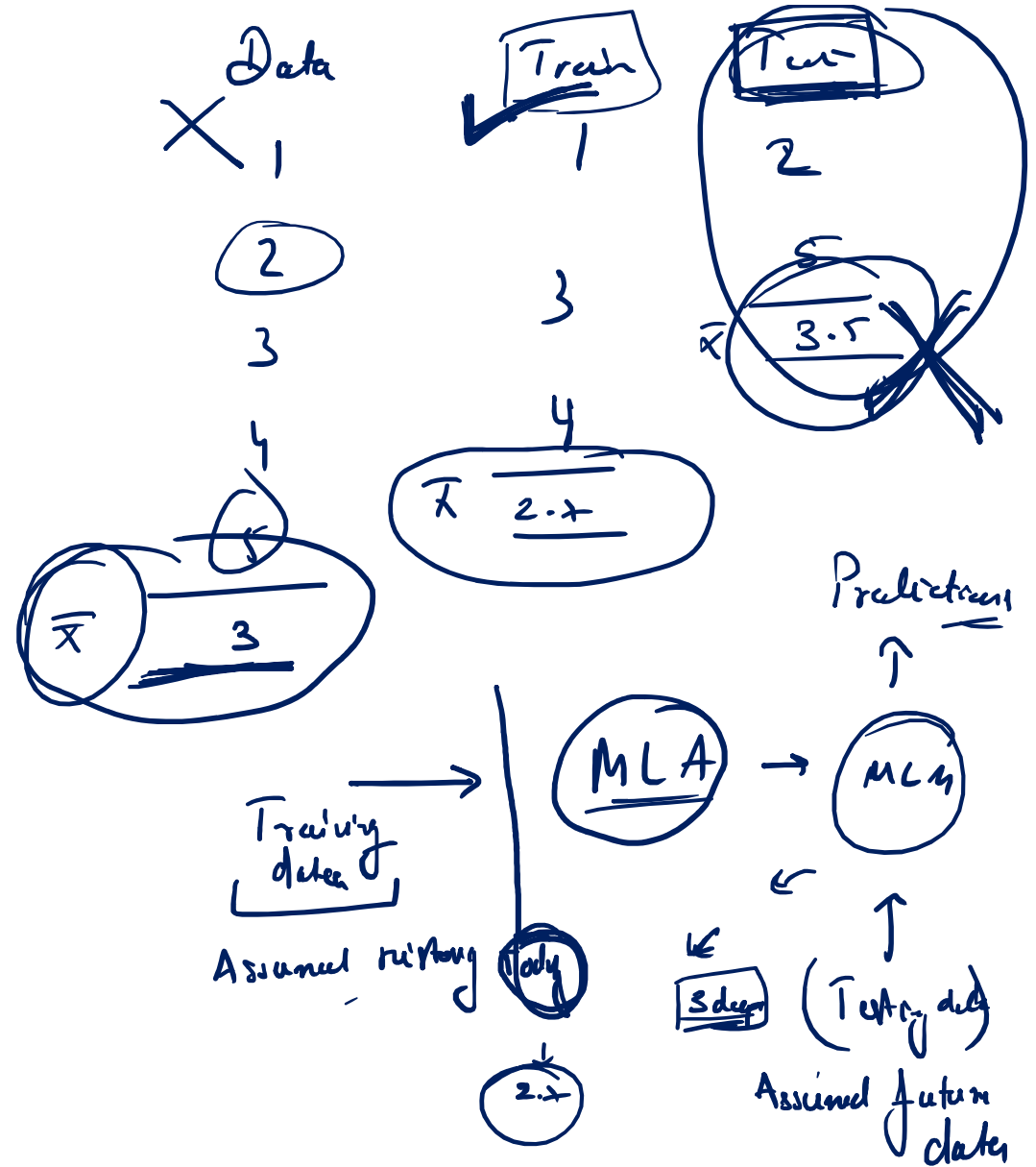Missing value    Outlier    Transform

joins

↳ MLA → MLA

08:29PM

Data preprocessing

(1) Import libraries

(2) load the data into python Env.

(3) Making sure the dtypes are correct.

(4) Extract target & features

(5) Split the data into Train & Test

(6) → Data Cleaning

Data

X 1

(2)

3

4

Train
1

Test
2

3

4

$\bar{X}$ 2.7

3

5

$\bar{x}$ 3.5

$\bar{X}$ 3

Predictions

Training data

MLA → MLM

Assumed missing body

5 days

Testing data

Assumed future data

2.7

→ Data Cleaning

100° ←

① Missing value Treatment

↳ executed col-wise

→ If missing value ≤ 5% → Impute / replacing the values

Mean → Continuous data | Median ⇒ discrete data | mode = Categorical data

→ If missing values ≥ 10% → Drop the col | random Imputation (IQR)
↳ Business logic to replace

5%.

④%.   If missing value b/w 5% – 10%.
then it is upto the data scientist
to choose from any above case

If cot A has
var = 1
↳ LalB

② Outlier Treatment $\longrightarrow$ Outlier is good $\longrightarrow$ Nothing to be done, fraud detection $\qquad$ w/o outlier

$\qquad$ "-"

$\qquad$ abc

$\longrightarrow$ Outlier is bad $\longrightarrow$ general pattern $\qquad$ with #

If outliers $\geq 10\%$ $\longrightarrow$ Split the dataset row-wise & Create 2 Models

with out $\qquad$ w/o outlier

$\longhookrightarrow$ If outliers $\leq 5\%$ $\longrightarrow$ Impute/replace them with either LL or UL

**Median method**

$$LL = Q_1 - 1.5(IQR)$$

$$UL = Q_3 + 1.5(IQR)$$

**Mean method**

$$LL = \mu - 3\sigma$$

$$UL = \mu + 3\sigma$$

$\longhookrightarrow$ If outlier is b/w $5\% - 10\%$, then upto the data scientist to decide

(7) → Analyze the Statistical Summary

(8) → Feature Selection ( Corr Analysis )

(9) → OPTIONAL STEP → Data Scaling

↳ To bring data points in Close Proximity

$(\mu, \sigma) \to$ fit()

a) Standard Scalar → Z-Score ⇒ $\dfrac{X - \mu}{\sigma}$ → Convert your data in a Standardized → $N.D. (0,1)$
$\mu \quad \sigma$

b) Min-Max Scalar →

$(X_{min}, X_{max})$ → fit()

Normalization ↳ into a range ⇒ $\dfrac{X - X_{min}}{X_{max} - X_{min}}$ ⇒ $[0, 1]$
of $[0, 1]$

↳ Apply is done three transform()

→ Data Encoding → Convert Categorical cols to Numerical values

such that the learning is not impacted

**Label Encoding** ~~only for Target~~

↳ Label Encoder() → Alphabetic order

→ Map() → Custom order

**Categorical Variables**

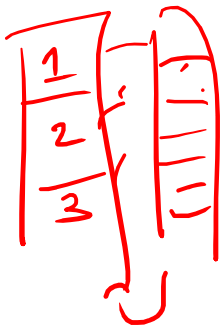Ordinal        Nominal

↓           ↓

has an in-built Rank    No in-built Rank

⇅

Numerical

↳ just use the Rank

1 { high Rating    3

2 { medium Rating   2

3 { low Rating     1

**Label Encoding**

| 1 | ' ' | ' ' |
|---|-----|-----|
| 2 | ' ' | ' ' |
| 3 | ' ' | ' ' |

**Dummy -Variables**

↳ One-hot Encoding

| ' A ' |
|-------|
| ' E ' |
| ' H ' |

⇒

| A | E | H |
|---|---|---|
| 1 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 0 | 1 |

2

3

No Revision

A
B
C



gender: M, F, F, F, M, M

t male: 1, 0, 0, 0, 1, 1

t fmale: 0, 1, 1, 1, 0, 0

$-1$   $0.7$   $0$   $0.7$   $+1$

$r = -1$   High correlate

F vs M

$r \rightarrow [-1, +1]$

No corr | $r \rightarrow$ ≈ 0
$(-0.1, +0.1)$

weak corr | $r \rightarrow (-0.7, -0.1)$
$(0.1, 0.7)$

High corr ⇒
$(-0.7, -1)$
$(0.7, +1)$