# Linear Regression

1. Business Problem understanding (Domain Expertise)

## Machine Learning Project Pipeline/ Workflow

2. Data ingestion/Collection → SQL
(Raw data)

3. Data understanding → Rows (Records), Column (features)

85-90% time →

4. EDA (Exploratory Data Analysis) (Statistics)

Data Normalization/ Standardization

→ Outliers ↳ deal

Quality data

Missing Values ↳ Imputation

① Remove
② focus on what kind of Cost function
③ To get the best optimal value

Visualization ↳ Correlation

feature Selection →

heatmap

feature engineering

Data | 1000 Records — Split | Plots $\longrightarrow$ seaborn, matplotlib

Train Set (80%) — 800 Records

Test Set (20%) — 200 Records

⑤ Model Training $\longrightarrow$ which specific algorithm suits the most??

⑥ Model Testing $\longleftarrow$ unseen data for the model

Model Predictions $(\hat{y})$

Actual values $(y)$

⑦ Model Evaluation

Regression
- MSE
- MAE
- RMSE
- R-Squared
- Adjusted R-squared

classification tasks
- Confusion matrix
- Accuracy
- Precision
- Recall
- f1-score

⑧ Model Deployment

## Data Standardisation

Biased Model

2

2500

'MedInc', 'HouseAge', 'AveRooms', 'AveBedrms', 'Population', 'AveOccup', 'Latitude', 'Longitude'

data will be within one specific range

↳ model to be unbiased

① Standard Scaler ——→ Standard normal form

② Min-max Scaler

$\mu = 0$
$\sigma = 1$

fit_transform ——→ train data
$(\mu, \sigma)$

transform ——→ test
$(\mu, \sigma)$      data

SNF

$z\text{-score} = \dfrac{x_i - \mu}{\sigma}$

| Age |  |
|-----|---|
| 23  | → |
| 45  | → |
| 67  | → |
| 49  | → |
| 62  | → |
| 22  | → |

fit()
↳ $\mu, \sigma$

transform()
↳
$z\text{-score} =$
$\dfrac{x_i - \mu}{\sigma}$

$$\hat{y} = m_1 x_1 + m_2 x_2 + \text{-----} + m_8 x_8 + c$$

Coefficients = 8

intercepts = 1

MAE   vs   MSE

amplify error

$$\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

$$\frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|$$

→ Residuals

↳ More robust to outliers

$$RMSE = \sqrt{MSE}$$

↳ ideal ⟶ house price

↓ 5 lakh

↳ high error

temp prediction in celsius

$RMSE = 2°c \longrightarrow$ may be acceptable

0.5 lakh ↳ decent

$10°C \longrightarrow$ poor prediction

Lower $RMSE \longrightarrow$ model
prediction
will be
better