**\* Data cleaning:**

**① Data collection:**

Multiple Sources Like

① web Scraping
② Databases
③ API
④ Surveys

} Data Storage

**② Raw Data:**

**• Issues in Raw data are below:**

① Missing values
② Duplicate values
③ Spelling errors
④ Outliers
⑤ Incorrect datatypes
⑥ Inconsistent formatting

③ Data cleaning Techniques:
① Summary Statistics
② Distribution plots
③ info ( data types / Null count
④ frequency → Categorical data
⑤ Visual inspect

④ Handling Missing data:

① Delete rows / columns:
→ More than 50 % of data is
missing.

② Imputation: Filling Missing values
① For Numerical columns:
→ Mean / Median / Mode

② For categorical values:
→ Mode / unknown

③ Datetime columns ( Time Series )
→ forward fill / backward fill

④ Predictive Imputation using
ML Models.
( for Most Accurate information
using ML Model to fill the data)

⑤ Inconsistent formatting :
→ Replace / lowercase / uppercase
→ Extra spaces Removal.

⑥ Handling duplicate Rows:
① Remove duplicate Records.

② check for primary key.
③ Remove near duplicate Records.

⑦ Correcting Datatypes:

① Numeric → int / float
② Categorical → object
③ Datetime → datetime

⑧ Handling outliers:

① Genuine outliers
② Errors
   ↳ Remove outliers using
      Z-score, IQR
   → Domain Knowledge
   → Transformming Data