




Data Preprocessing

By: Matt Collins, Surya Maddali



Attendance: Meeting Name “PreProcess”





Table of contents

01

**Why Preprocess
Data?**

02

Datasets & Libraries

03

Data Preparation

04

Data Cleaning

05

Exploratory Data Analysis

06

Data Transformation



Why Preprocess Data?

Data Quality

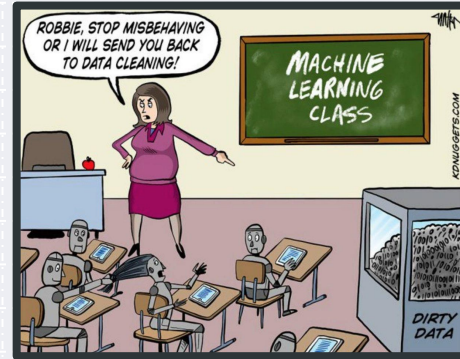
- Incomplete Data
- Noisy Data
- Data Discrepancy

Incompatibility

- Data Transformation
- Normalization

Interpretability

- Pattern recognition
- Bias removal



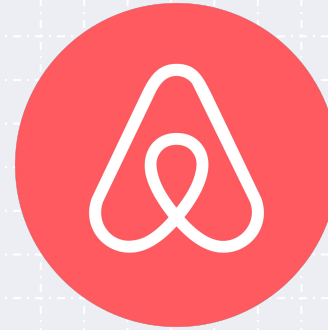
Colab



Datasets and Libraries

Datasets

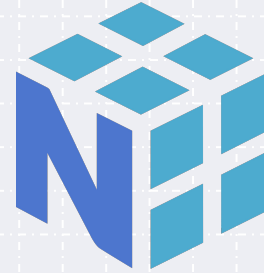
- *Ideally imported in csv form*
 - Comma-separated values
- Python - *Pandas*
- AirBnB open dataset



Datasets and Libraries

Libraries

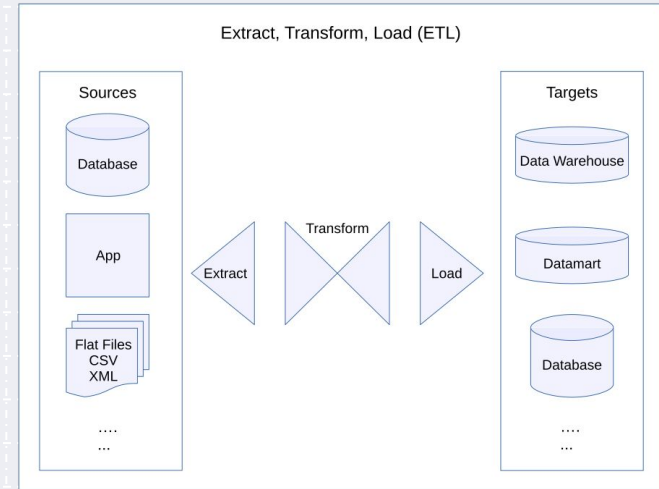
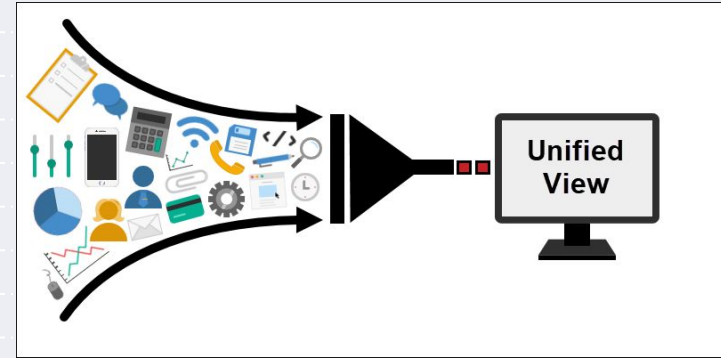
- *Pandas* library
- *Numpy* library
- *Sci-kit learn*
- *matplotlib*



Data Preparation

Data Integration

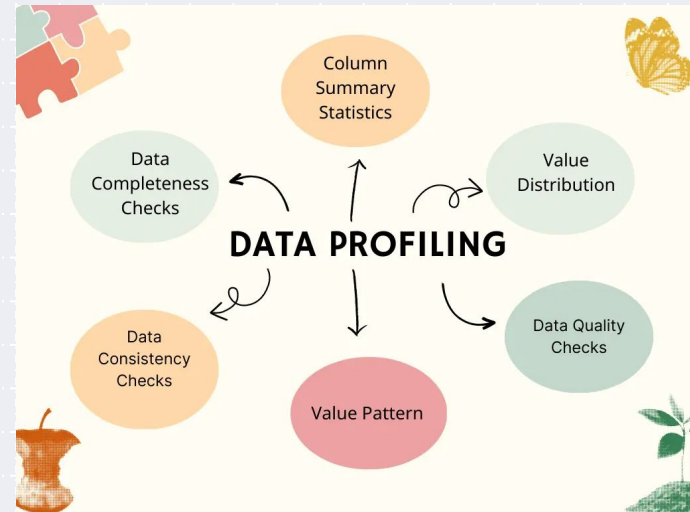
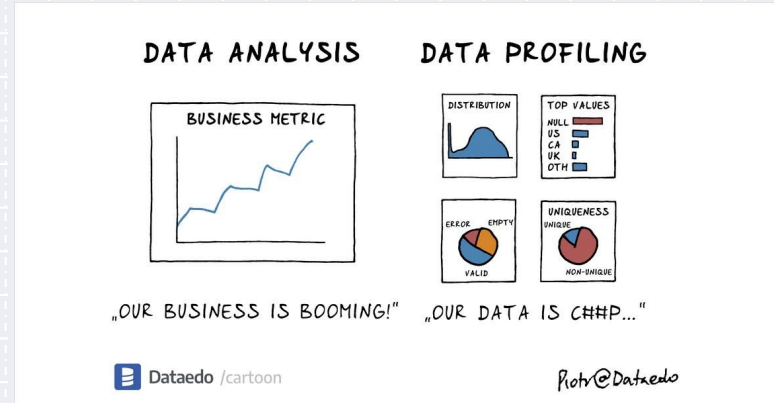
- Different sources → One View
- Timing varies
- Extract, Transform, Load
- Data Lakes/Warehouse



Data Preparation

Data Profiling

- Assess data quality
- Structure, content, relationship discovery
- Repeated throughout process

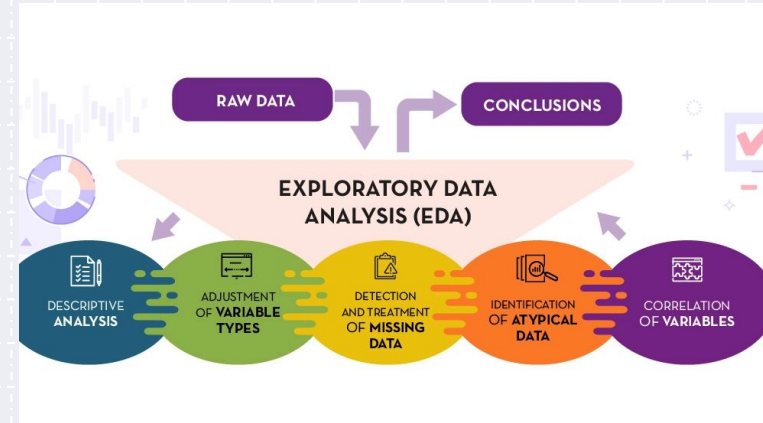
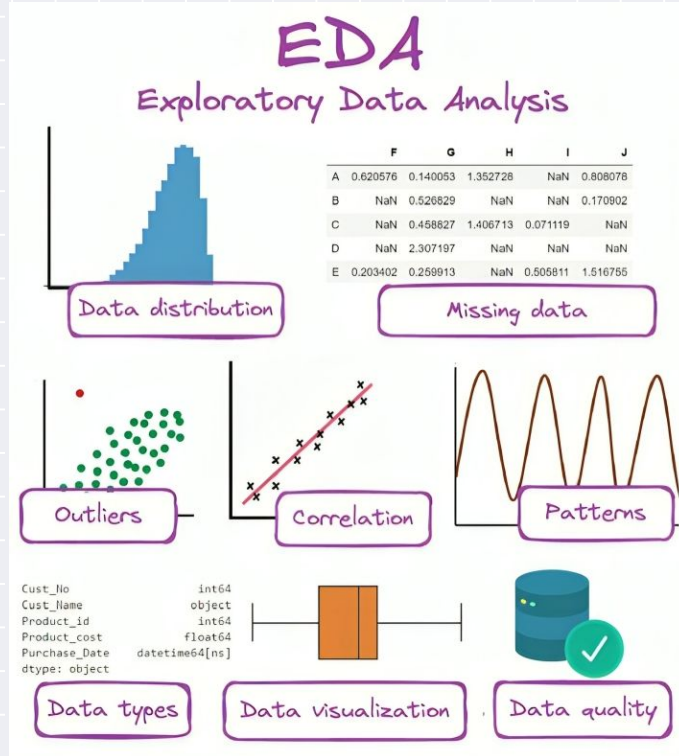


Data Cleaning

- Correcting inaccurate, incomplete, or irrelevant information from a dataset
- Nulls
- Noisy Data
- Inconsistent Data



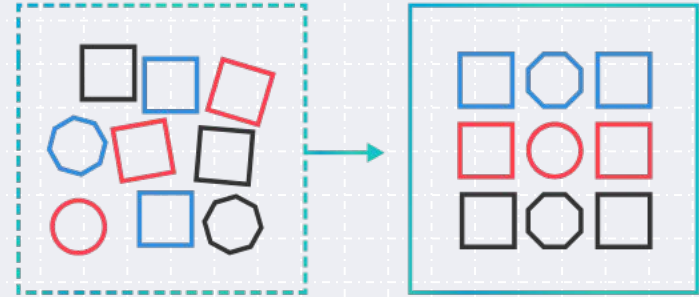
Exploratory Data Analysis



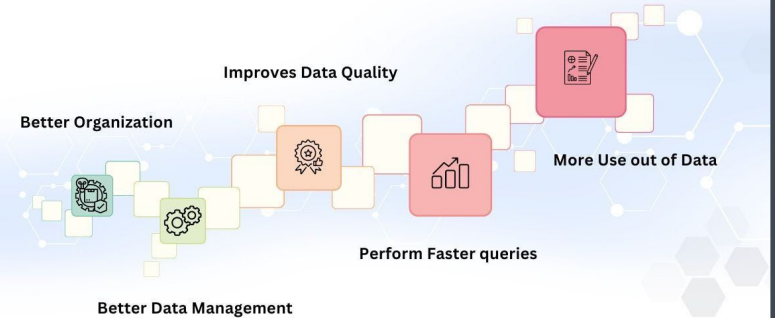
- Visualizations to analyze data
- Reveals patterns, checks assumptions
- More detailed than profiling
- For scientists and stakeholders

Data Transformation

- Changing structure or format
- Goal to improve usability
- Various techniques
- Feature extraction
- Encoding and normalization



Advantages of Data Transformation

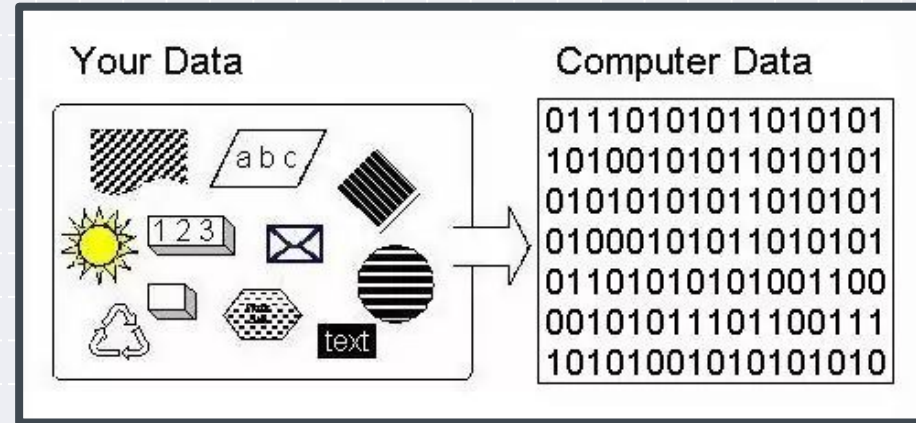


Data Transformation

Encoding

- Categorical data → Quantitative
- Label encoding vs. one hot encoding
- One number per unique entry
- Columns of unique entries

Label Encoding			One Hot Encoding			
Food Name	Categorical #	Calories	Apple	Chicken	Broccoli	Calories
Apple	1	95	1	0	0	95
Chicken	2	231	0	1	0	231
Broccoli	3	50	0	0	1	50



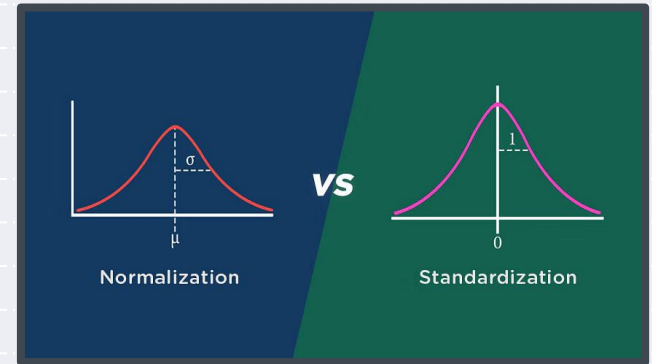
Data Transformation

Normalization/Standardization

- Quantitative → Quantitative
- Protection against outliers
- Plug data into formulas
- Normalization vs. Standardization

$$\text{new value} \rightarrow X' = \frac{\text{original value} - \min(x)}{\max(x) - \min(x)}$$

$$Z = \frac{X - \mu}{\sigma}$$



Data Source



ETL Tools



Data Warehouse



Data Transformation



BI Tools





Thanks NDL!

Feel free to stay and ask questions!

