

Cross Industry Standard Process for Data Mining (CRISP-DM)

CRISP-DM

Cross Industry Standard Process for Data Mining



CRISP-DM Phases

- **Business Understanding**
- **Data Understanding**
- **Data Preparation**
- **Modeling**
- **Evaluation**
- **Deployment**

Phase 1 – Business Understanding

- **Define problem or opportunity**
- **Assess situation**
- **Formulate goals**



Phase 2 – Data Understanding

- Data acquisition
- Data exploration



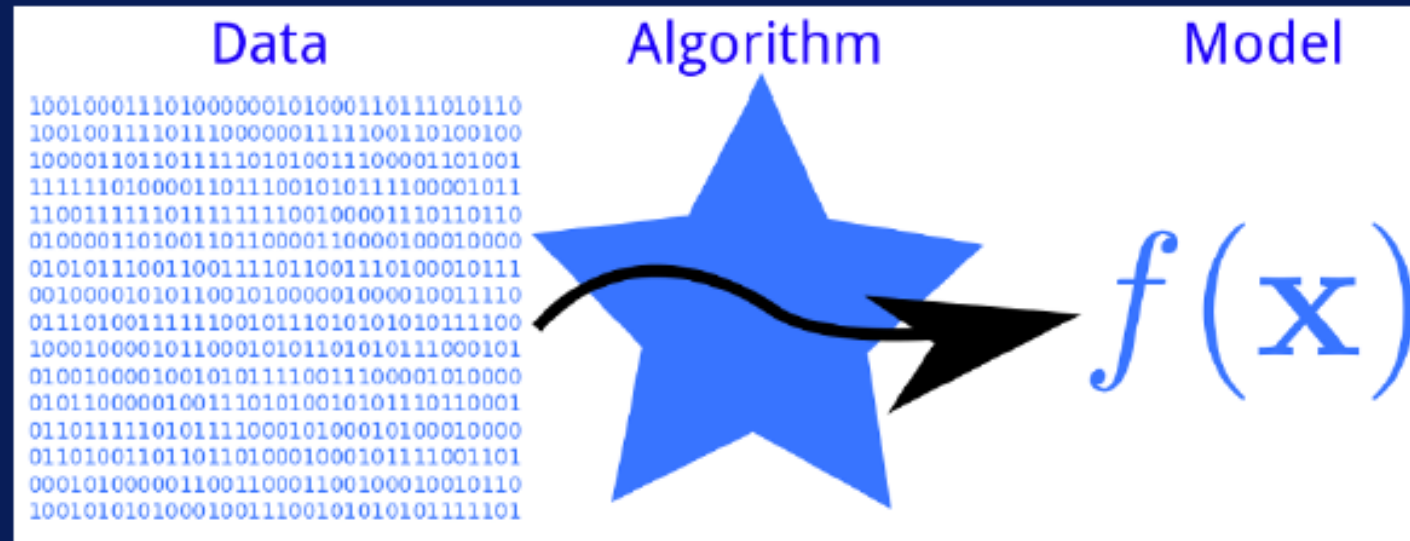
Phase 3 – Data Preparation

- Prepare data for modeling
- Address quality issues, select features to use, process data for modeling



Phase 4 – Modeling

- Determine type of problem
- Select modeling technique(s) to use
- Build model



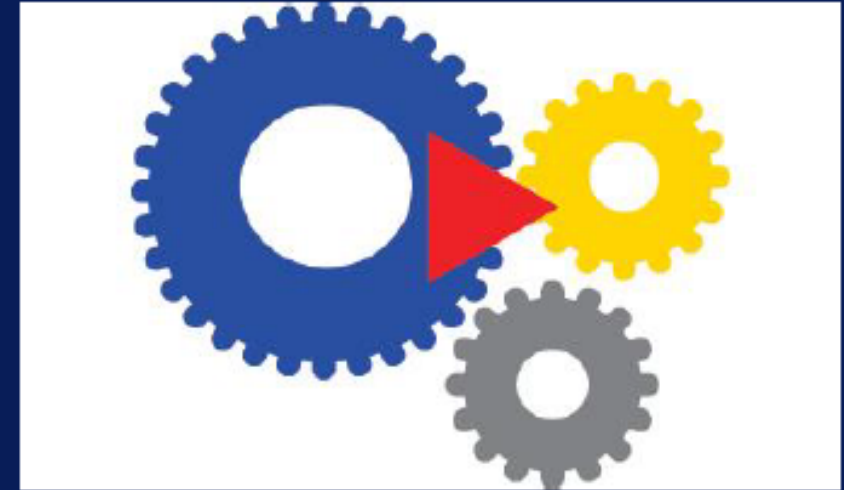
Phase 5 – Evaluation

- **Assess model performance**
- **Evaluate model results with respect to success criteria**



Phase 6 – Deployment

- **Produce final report**
- **Deploy model**
- **Monitor model**



Machine Learning Process

ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

Step 1: Acquire Data



Identify data sources

Collect data

Integrate data

ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

Step 2: Prepare Data

Step 2-A: Explore

Step 2-B: Pre-process

ACQUIRE

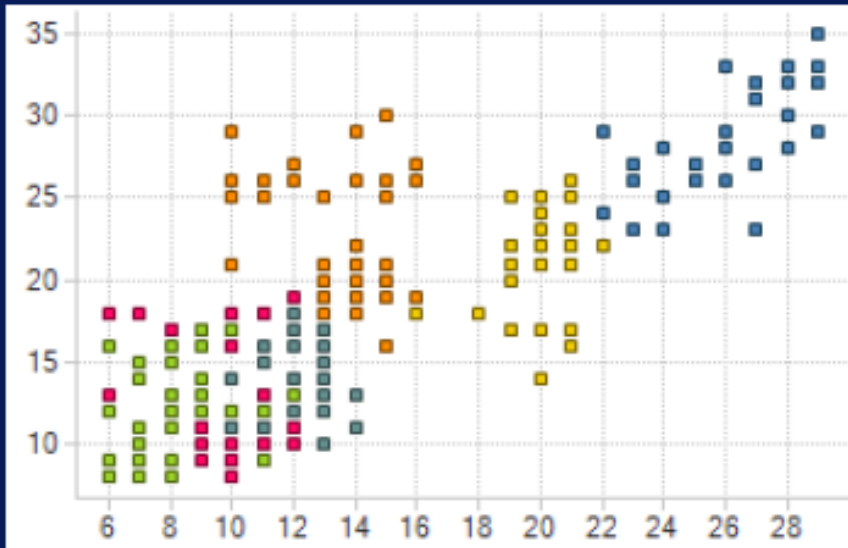
PREPARE

ANALYZE

REPORT

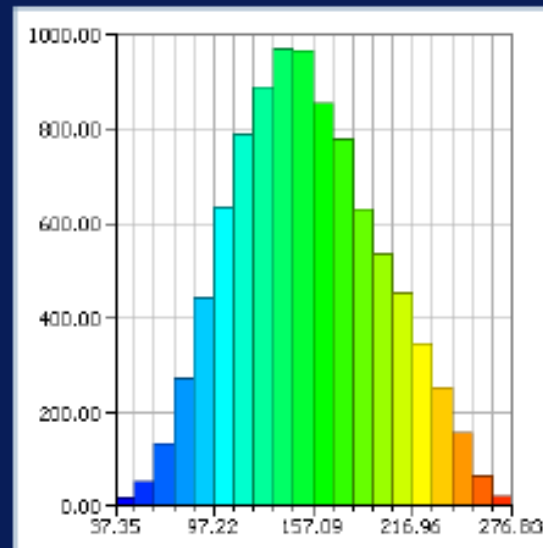
ACT

Step 2-A: Explore Data



Understand
nature of data

Preliminary
analysis



ACQUIRE

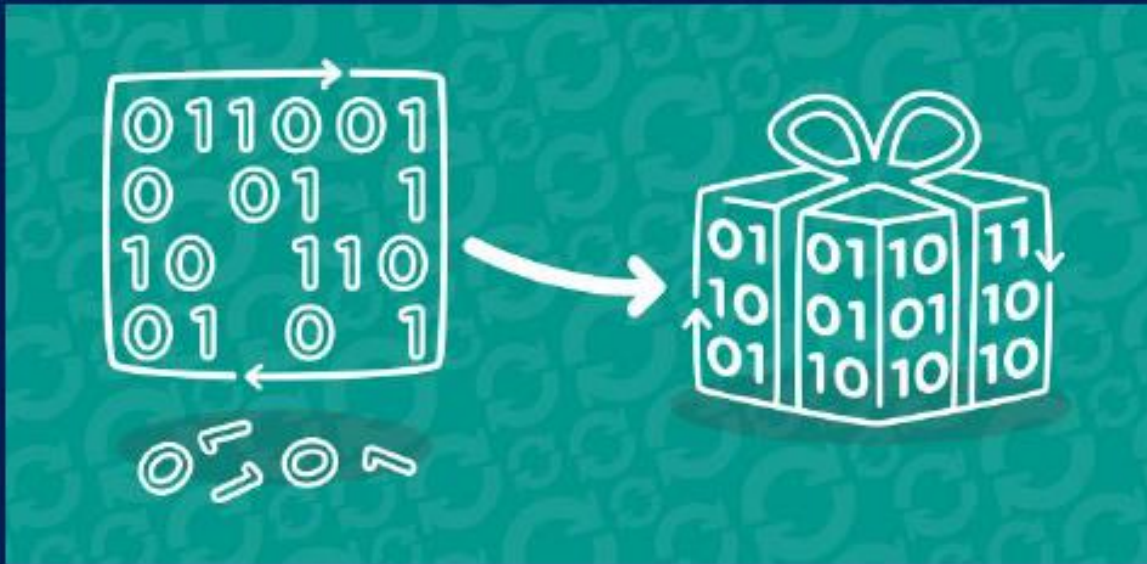
PREPARE

ANALYZE

REPORT

ACT

Step 2-B: Pre-process Data



Clean

Select

Transform

ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

Step 3: Analyze Data



Select analytical techniques

Build models

Assess results

ACQUIRE

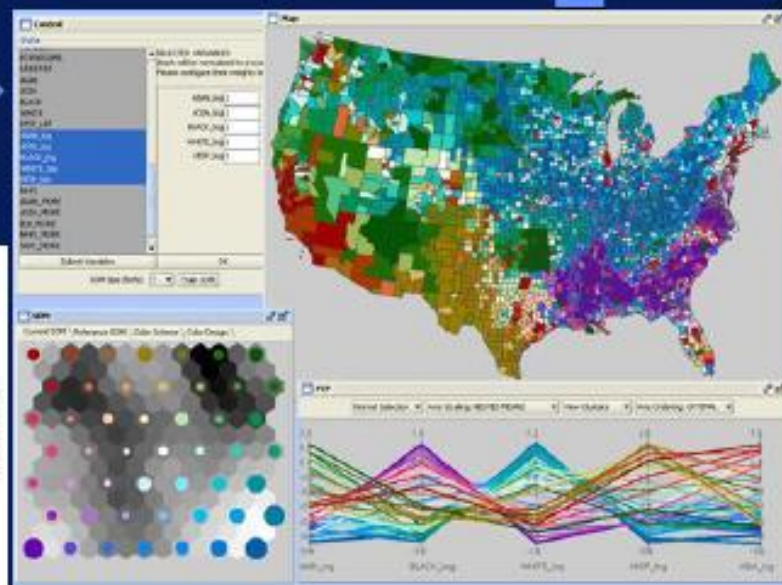
PREPARE

ANALYZE

REPORT

ACT

Step 4: Communicate Results



ACQUIRE

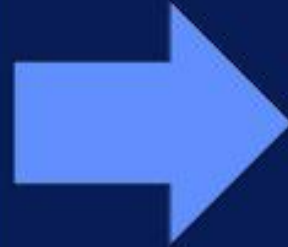
PREPARE

ANALYZE

REPORT

ACT

Step 5: Apply Results



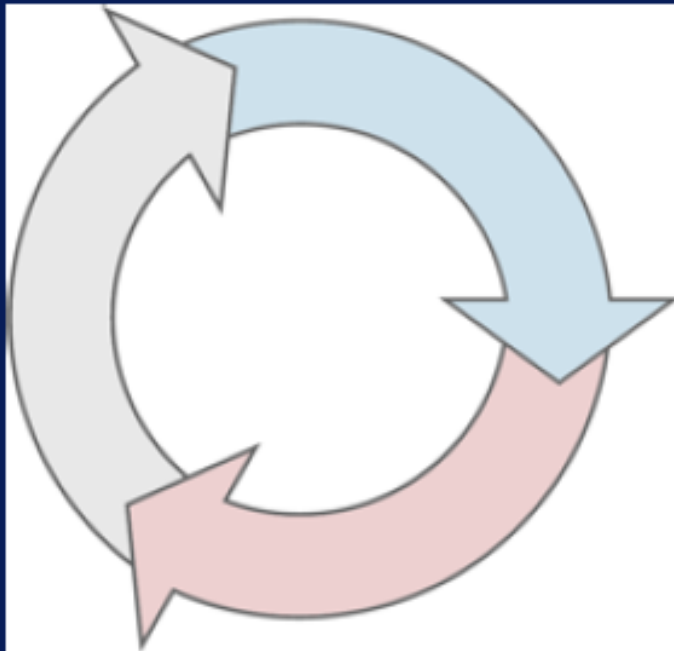
ACQUIRE

PREPARE

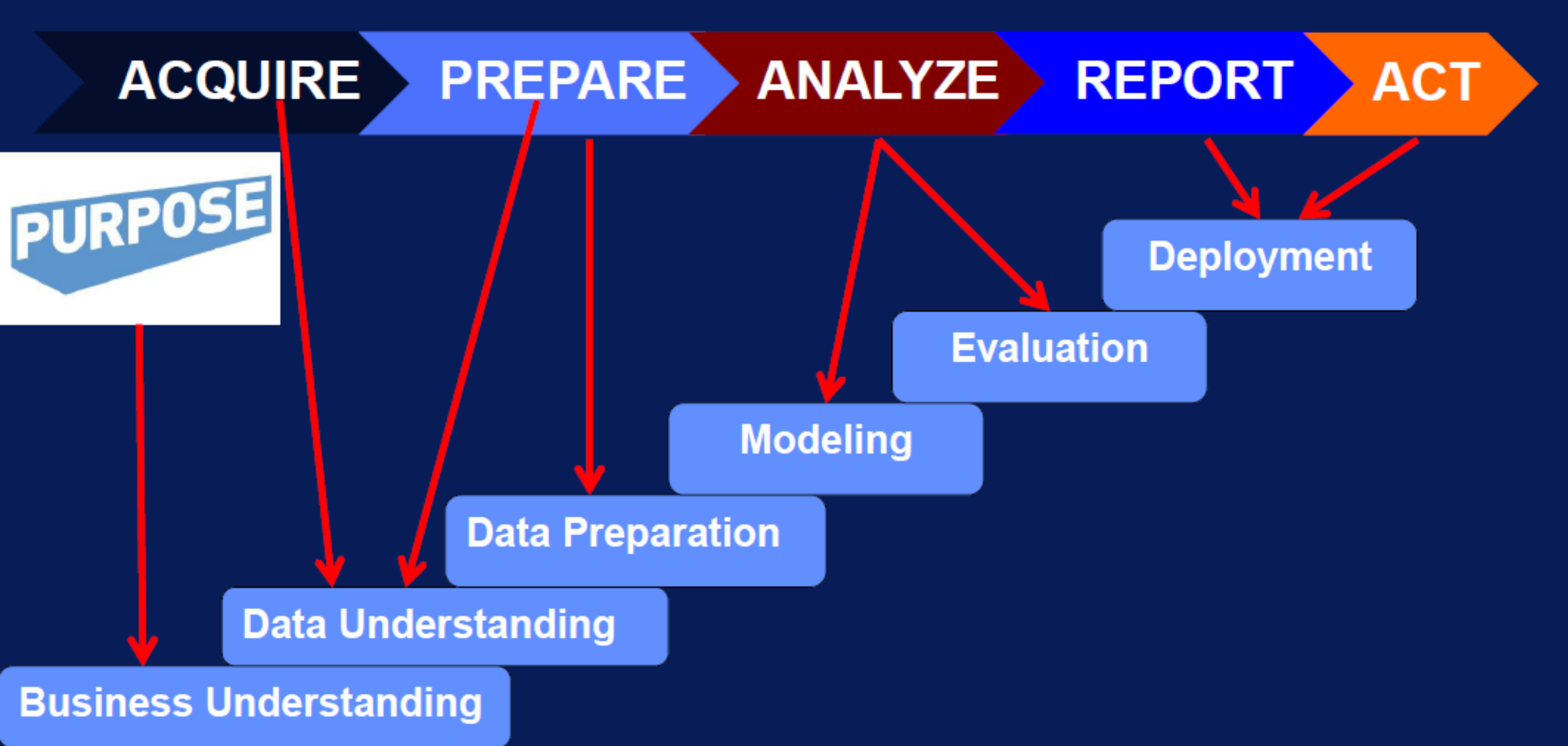
ANALYZE

REPORT


ACT



Iterative process



Goals and Activities in the Machine Learning Process



ACQUIRE

PREPARE

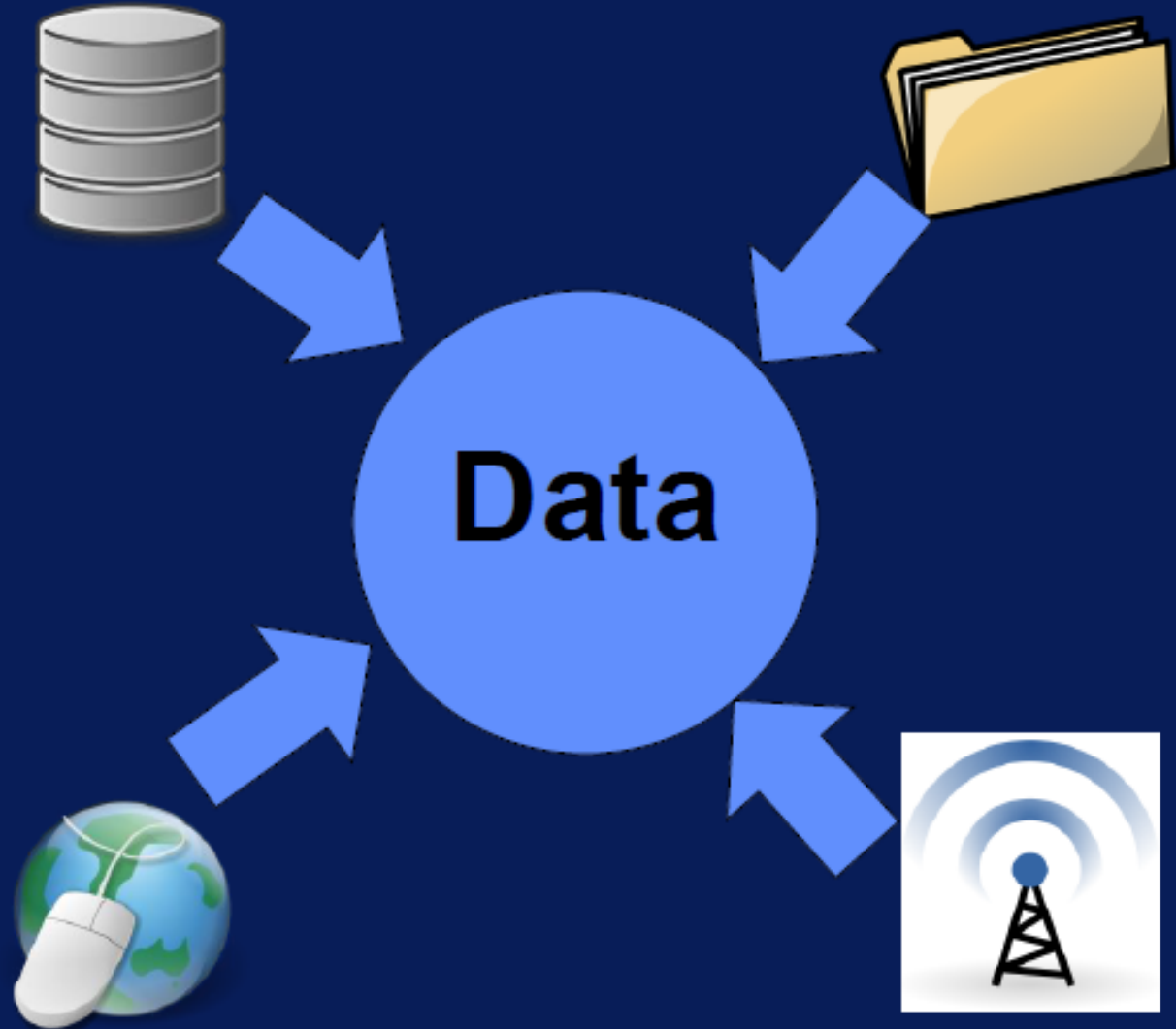
ANALYZE

REPORT

ACT

Goal: Identify and obtain all data
related to problem

Acquire Data



Identify data sources
Collect data
Integrate data



Step 2-A: Explore

Step 2-B: Pre-process



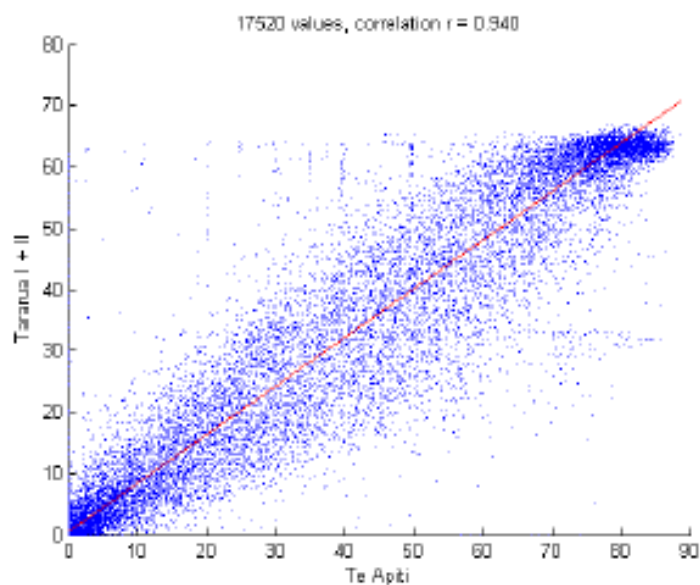
Why Explore?

Goal: Understand your data

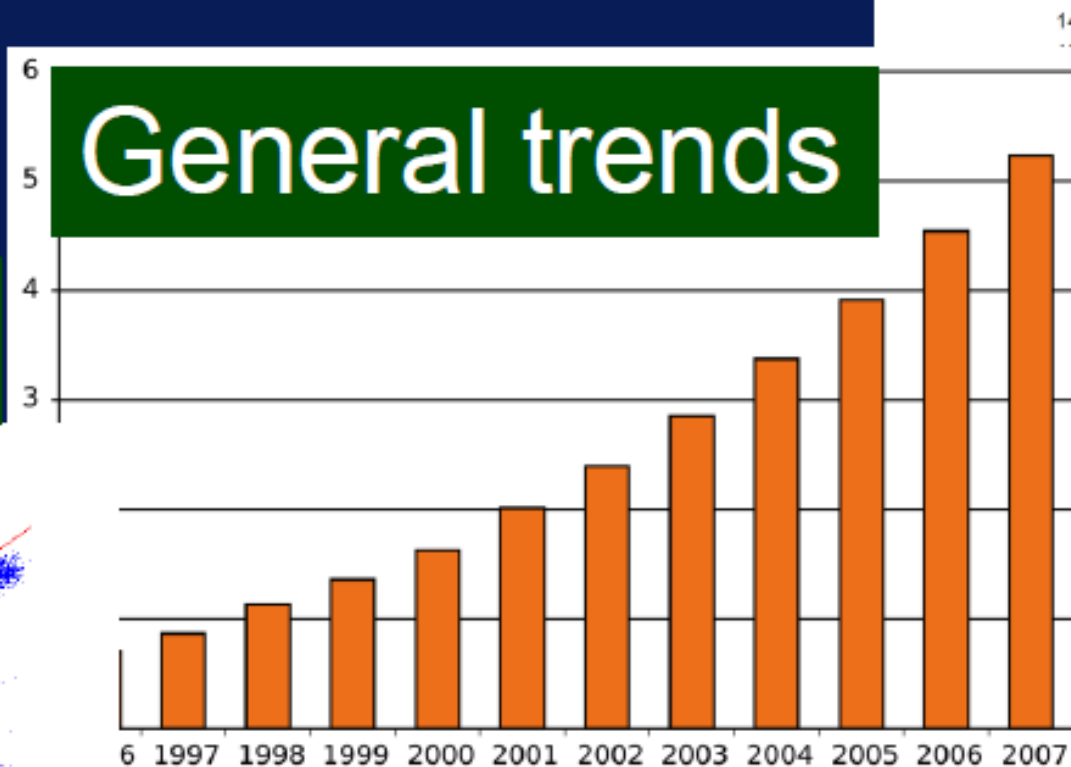


Why Explore?

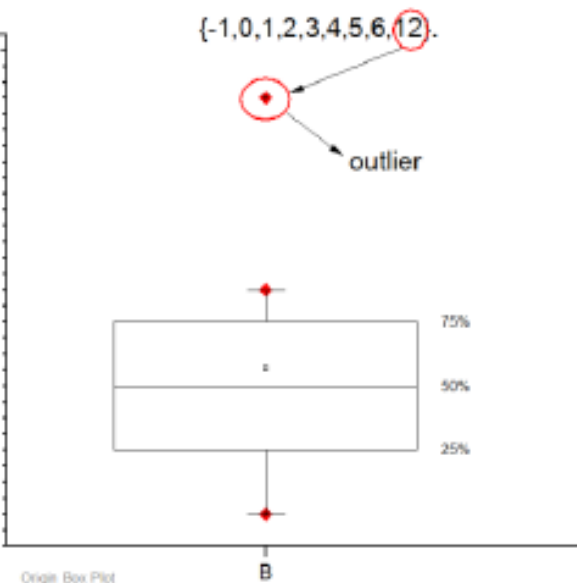
Correlations



General trends



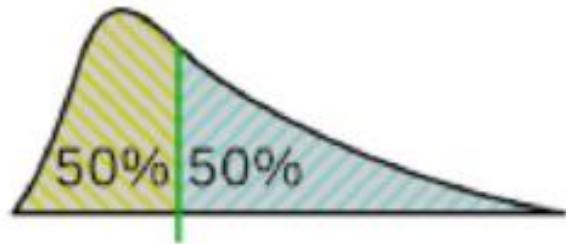
Outliers



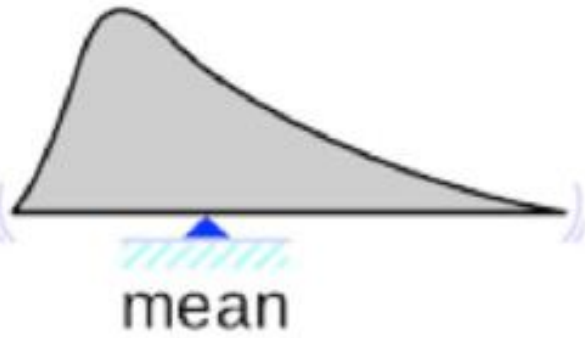
Describe Your Data



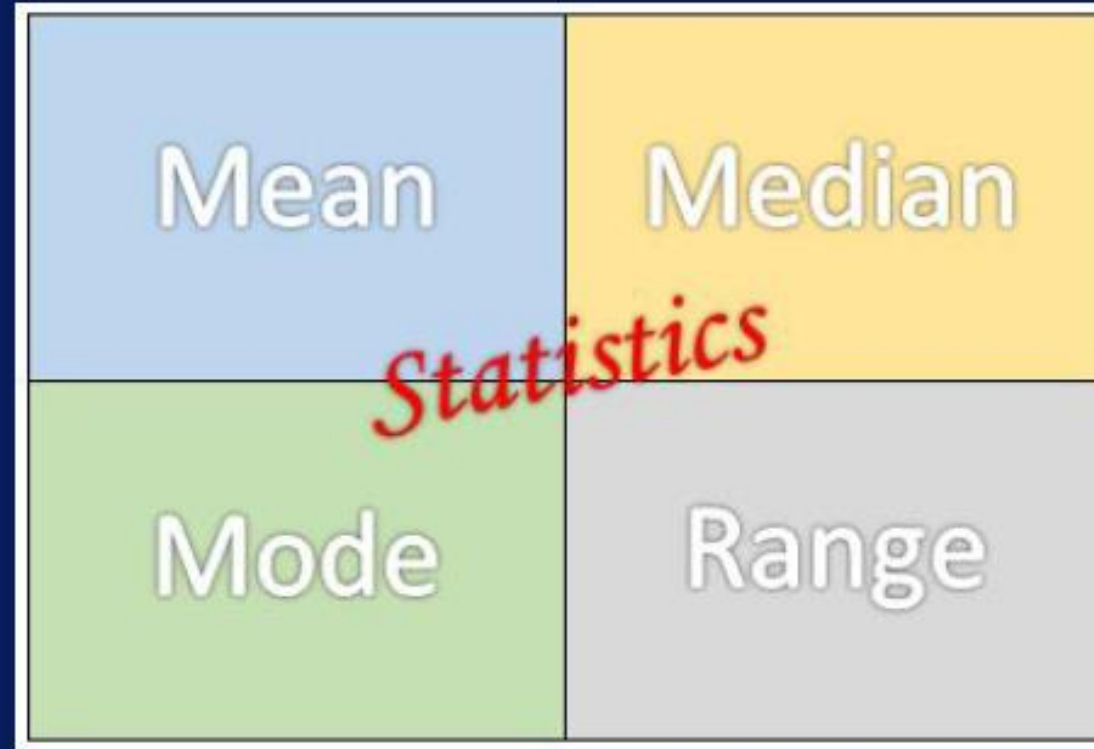
mode



median

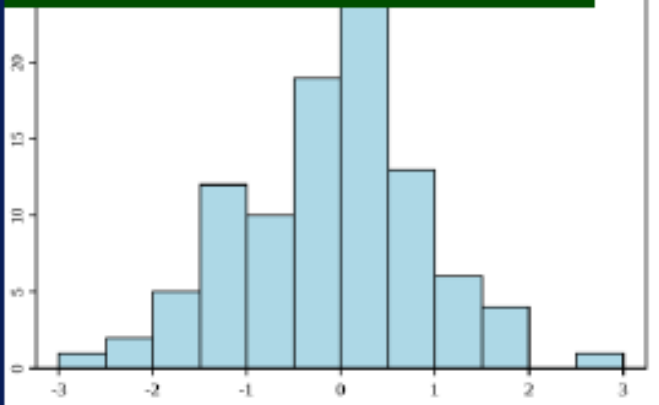


mean

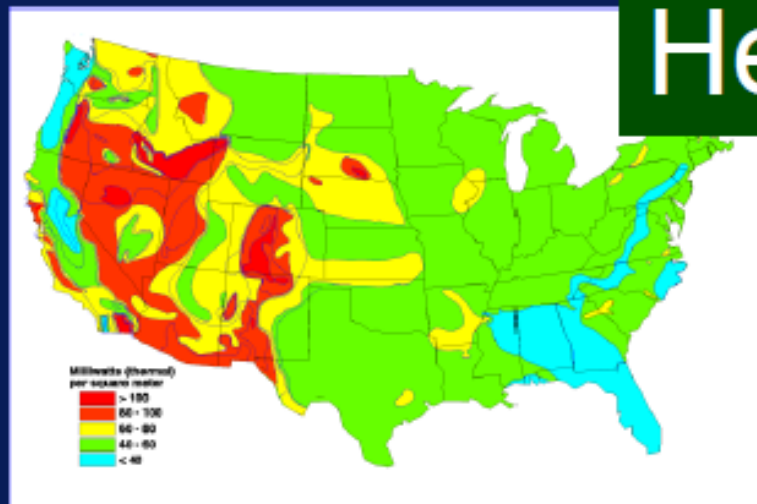


Visualize Your Data

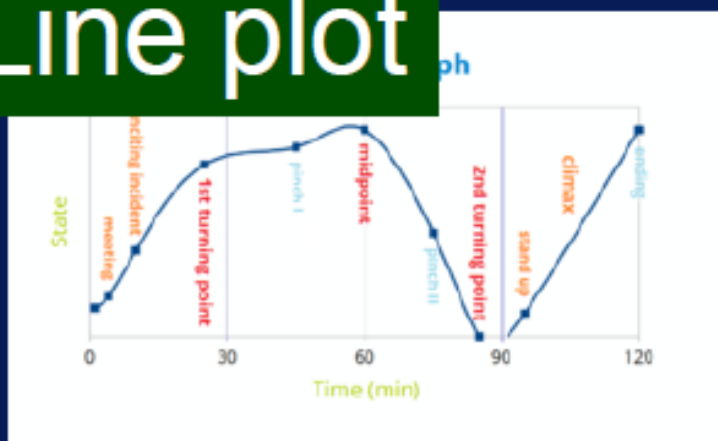
Histogram



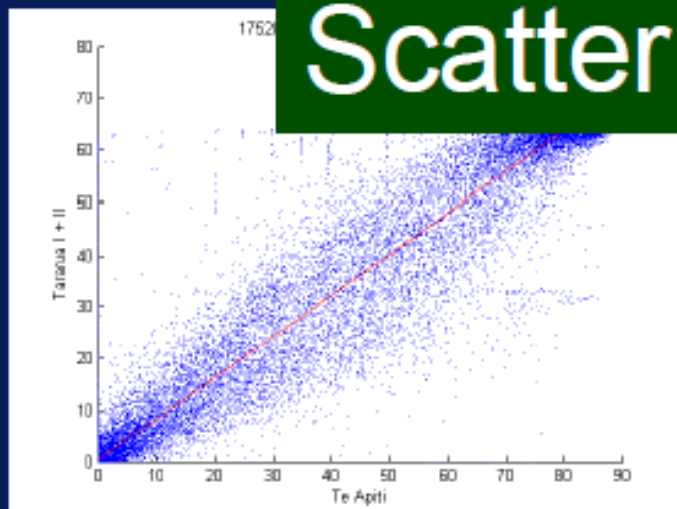
Heat map



Line plot



Scatter plot





Step 2-A: Explore

Step 2-B: Pre-process

ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

Step 2-A: Explore

Goal: Create data
for analysis

Step 2-B: Pre-process

Clean

Select

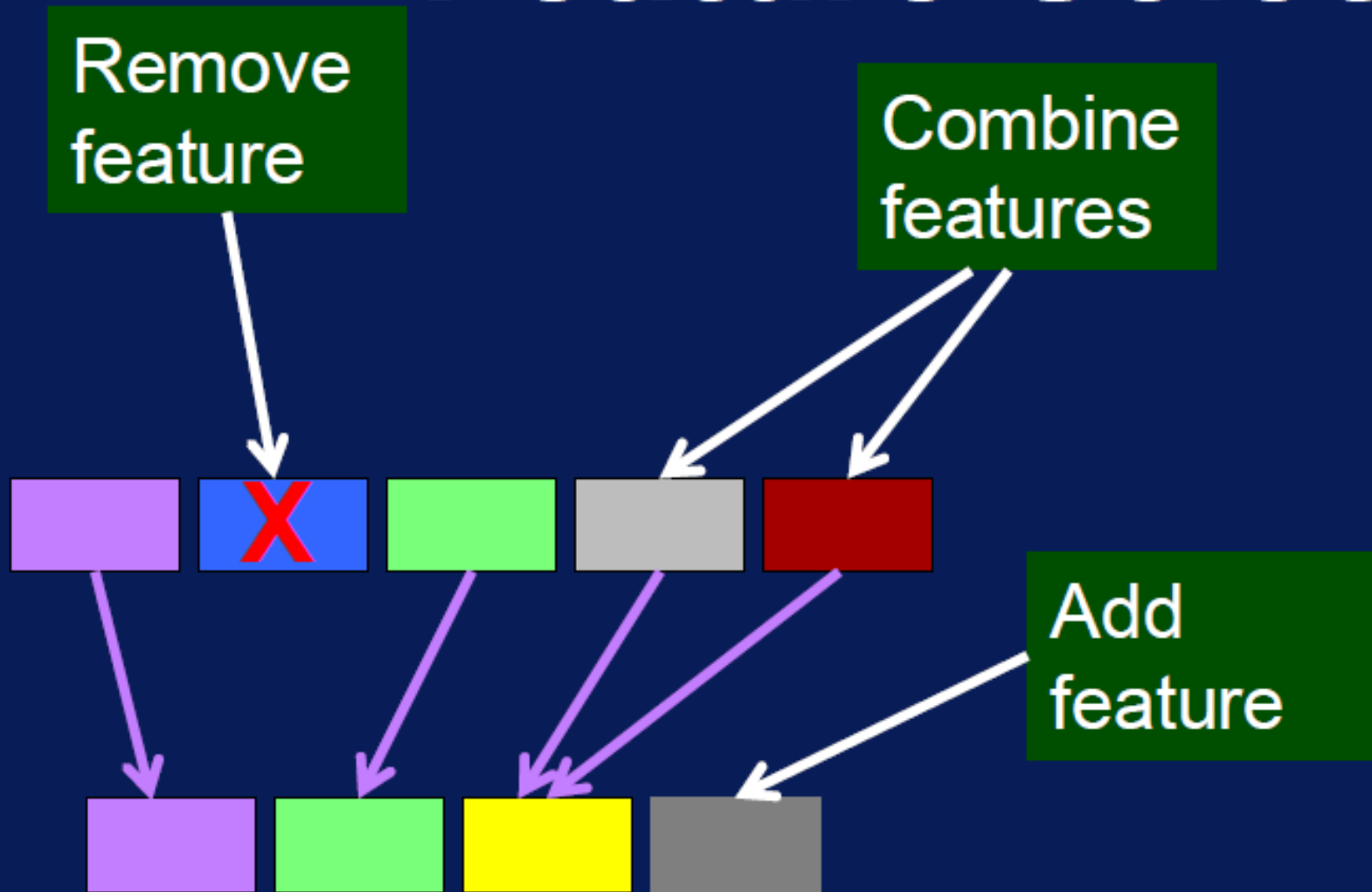
Transform

Data Cleaning

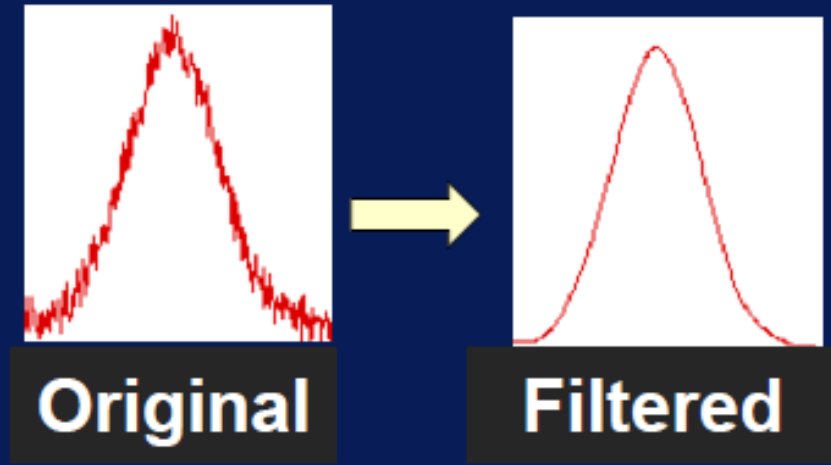
- **Missing values**
- **Duplicate data**
- **Inconsistent data**
- **Noise**
- **Outliers**



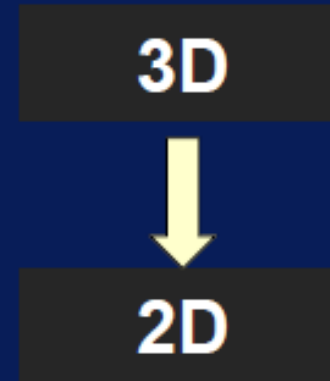
Feature Selection



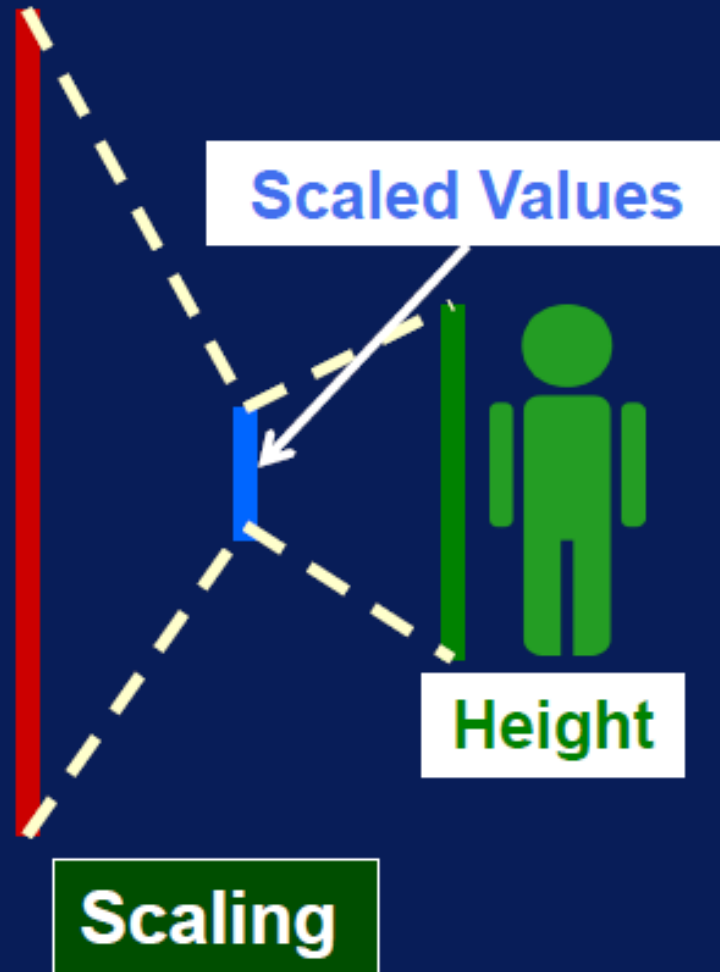
Feature Transformation



Filtering




**Dimensionality
Reduction**



Weight

Height



ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

Goals:

- Build model
- Evaluate results

Analyze

Select technique

Build model

Evaluate

Classification
Regression
Cluster Analysis
Association
Analysis



Analyze


Select technique



Build model



Evaluate results



ACQUIRE

PREPARE

ANALYZE

REPORT

ACT

Goal: Communicate results and
recommend actions

Present

Report



with



using



by University of California San Diego



Goal: Determine actions based on insights

Act

Determine action



Implement



Assess impact

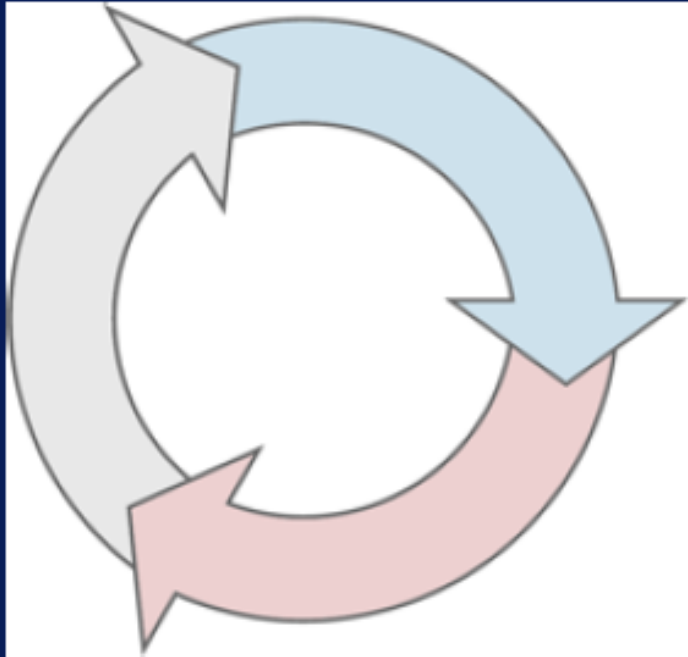
ACQUIRE

PREPARE

ANALYZE

REPORT

ACT



Iterative process