# Machine Learning is Fun!

# Reference

1. https://towardsdatascience.com/types-of-machine-learning-algorithms-you-should-know-953a08248861
2. https://omdena.com/blog/supervised-and-unsupervised-machine-learning/
3. https://www.cs.huji.ac.il/~shais/UnderstandingMachineLearning/understanding-machine-learning-theory-algorithms.pdf
4.

# What is Machine Learning?

"Learning is any process by which a system improves performance from experience."

- Herbert Simon

**Machine Learning is the study of algorithms that**

- **improve their performance P**

- **at some task T**

- **with experience E.**

**A well-defined learning task is given by <P, T, E>.**

# Traditional Programming

Data →
Program →
Computer
→ Output

# Machine Learning

Data →
Output →
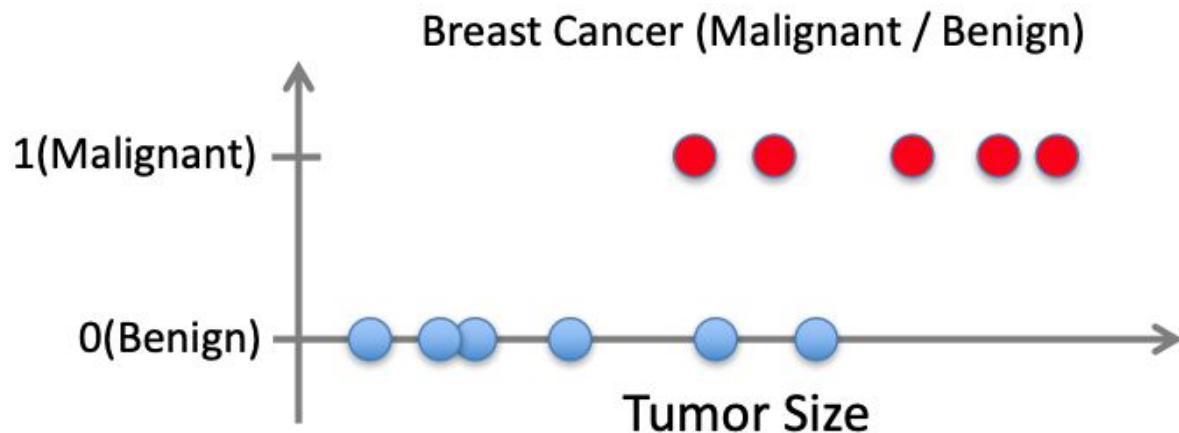Computer
→ Program

# Sample Applications

- Web search
- Computational biology
- Finance
- E-commerce
- Space exploration
- Robotics
- Information extraction
- Social networks
- Debugging software
- [Your favorite area]

# Types of Learning

- Supervised (inductive) learning
    Given: training data + desired outputs (labels)
- Unsupervised learning
    Given: training data (without desired outputs)
- Semi-supervised learning
    Given: training data + a few desired outputs
- Reinforcement learning
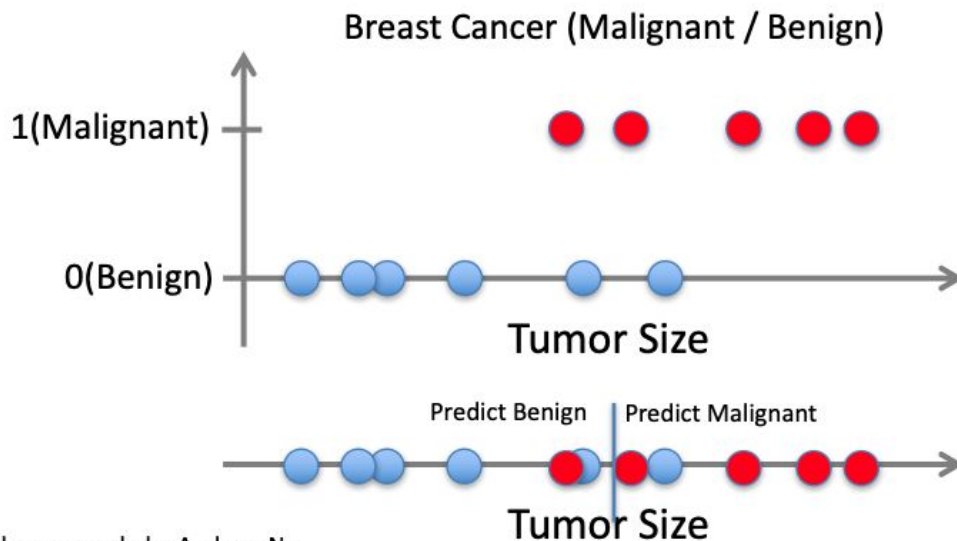    Rewards from sequence of actions

# Supervised Learning: Classification

- Given $(x_1, y_1), (x_2, y_2), ..., (x_n, y_n)$
- Learn a function $f(x)$ to predict $y$ given $x$
  - $y$ is categorical == classification
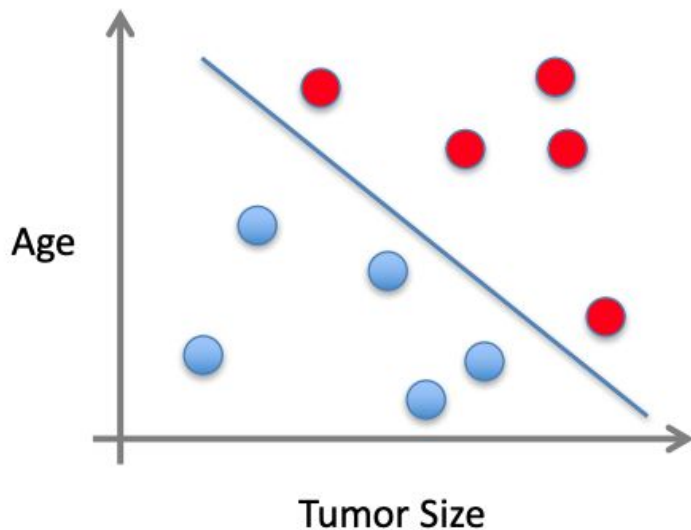
Breast Cancer (Malignant / Benign)

# Supervised Learning: Classification

- Given $(x_1, y_1), (x_2, y_2), ..., (x_n, y_n)$
- Learn a function $f(x)$ to predict $y$ given $x$
  - $y$ is categorical == classification



Breast Cancer (Malignant / Benign)

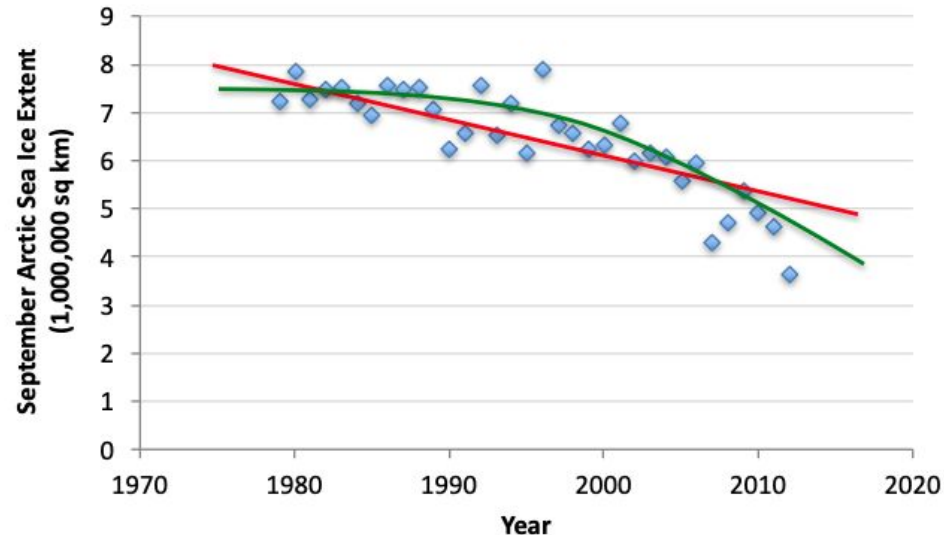Based on example by Andrew Ng

# Supervised Learning

- $x$ can be multi-dimensional
  - Each dimension corresponds to an attribute



- Clump Thickness
- Uniformity of Cell Size
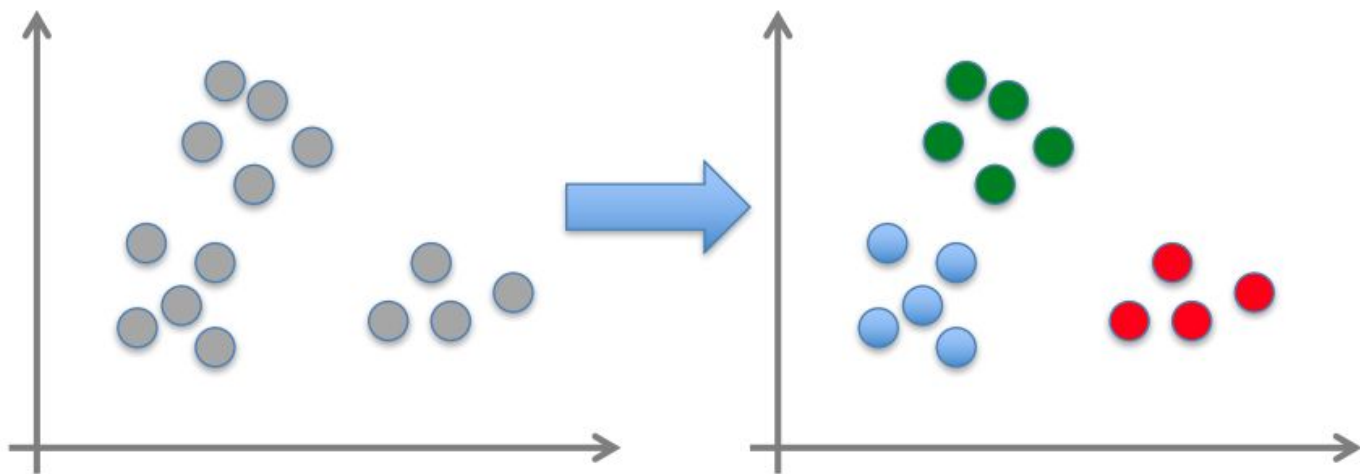- Uniformity of Cell Shape

...

# Supervised Learning: Regression

- Given $(x_1, y_1), (x_2, y_2), ..., (x_n, y_n)$
- Learn a function $f(x)$ to predict $y$ given $x$
  - $y$ is real-valued == regression
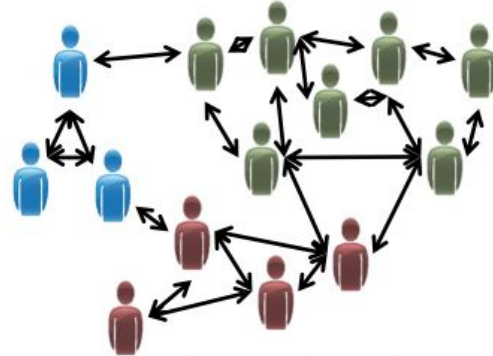
# Unsupervised Learning

- Given $x_1, x_2, ..., x_n$ (without labels)
- Output hidden structure behind the $x$'s
  - E.g., clustering

# Unsupervised Learning



Organize computing clusters

Social network analysis

Market segmentation

Astronomical data analysis

# Metrics to Evaluate your Machine Learning Algorithm

*Classification Accuracy*

*Logarithmic Loss*

*Confusion Matrix*

*Area under Curve*

*F1 Score*

*Mean Absolute Error*

*Mean Squared Error*

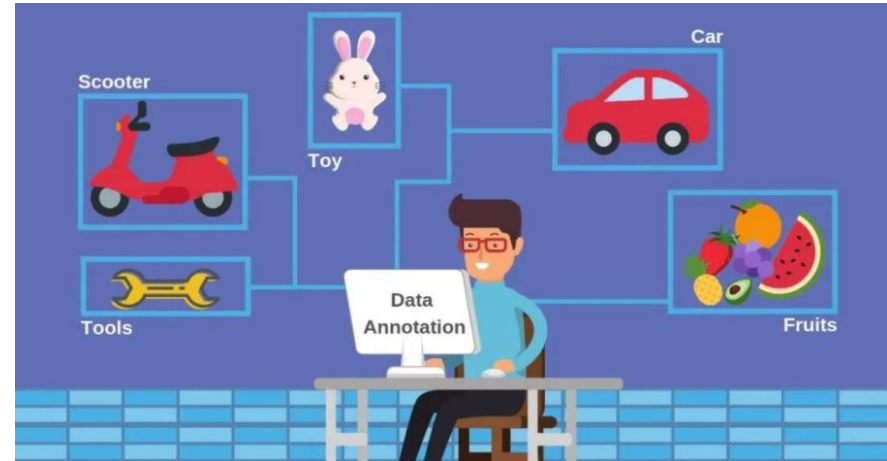| n=165 | Predicted: NO | Predicted: YES |
|---|---|---|
| Actual: NO | 50 | 10 |
| Actual: YES | 5 | 100 |

Confusion Matrix

There are 4 important terms :

- **True Positives** : The cases in which we predicted YES and the actual output was also YES.
- **True Negatives** : The cases in which we predicted NO and the actual output was NO.
- **False Positives** : The cases in which we predicted YES and the actual output was NO.
- **False Negatives** : The cases in which we predicted NO and the actual output was YES.
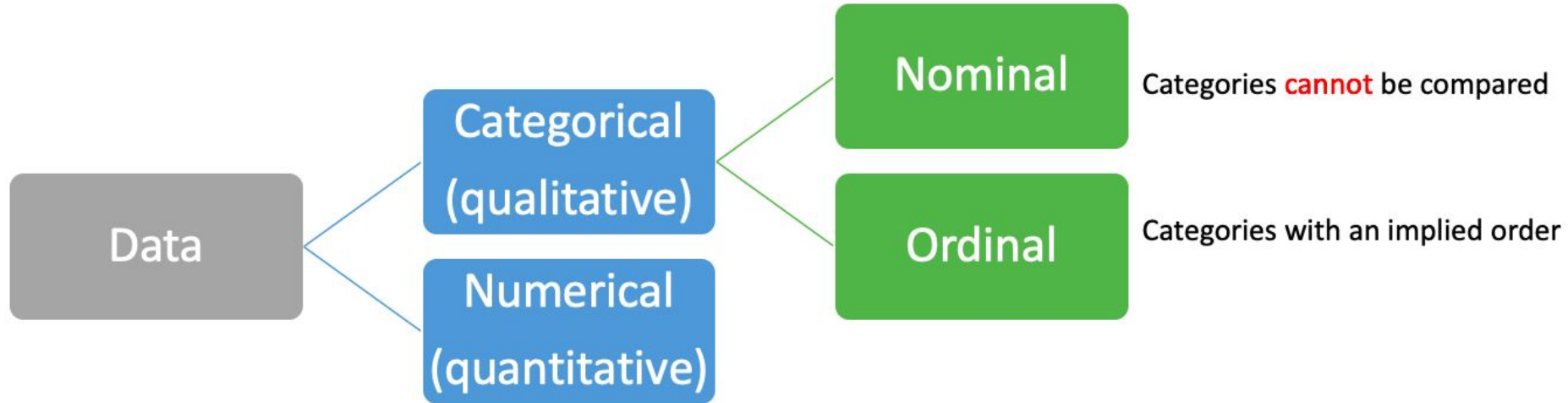
# Terminology: Data table

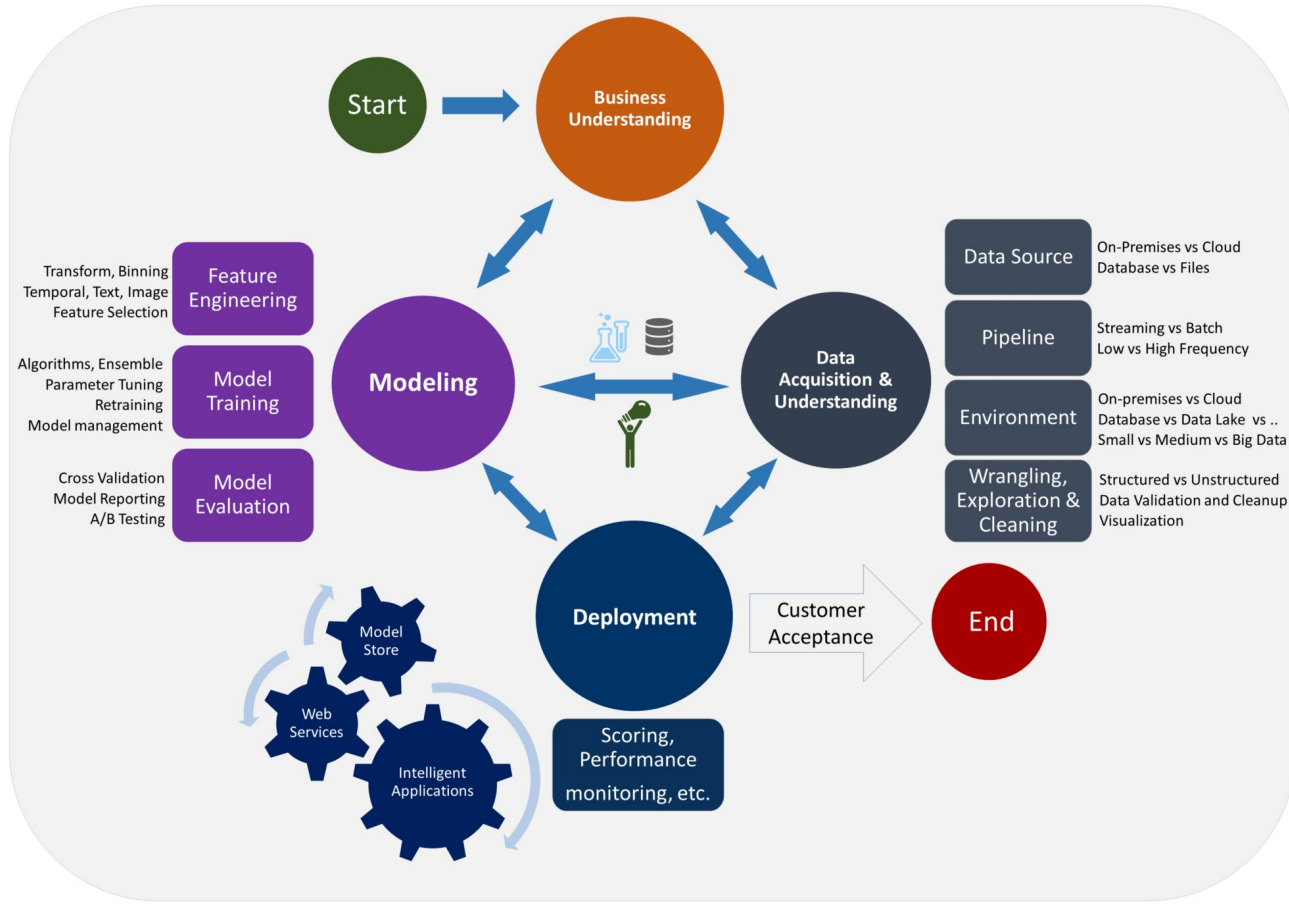| ID | INPUT | | | | LABEL |
|---|---|---|---|---|---|
| Day | Outlook | Temp. | Humidity | Wind | Decision |
| 1 | Sunny | Hot | High | Weak | No |
| 2 | Sunny | Hot | High | Strong | No |
| 3 | Overcast | Hot | High | Weak | Yes |
| 4 | Rain | Mild | High | Weak | Yes |

# What is the difference between supervised and unsupervised learning techniques?
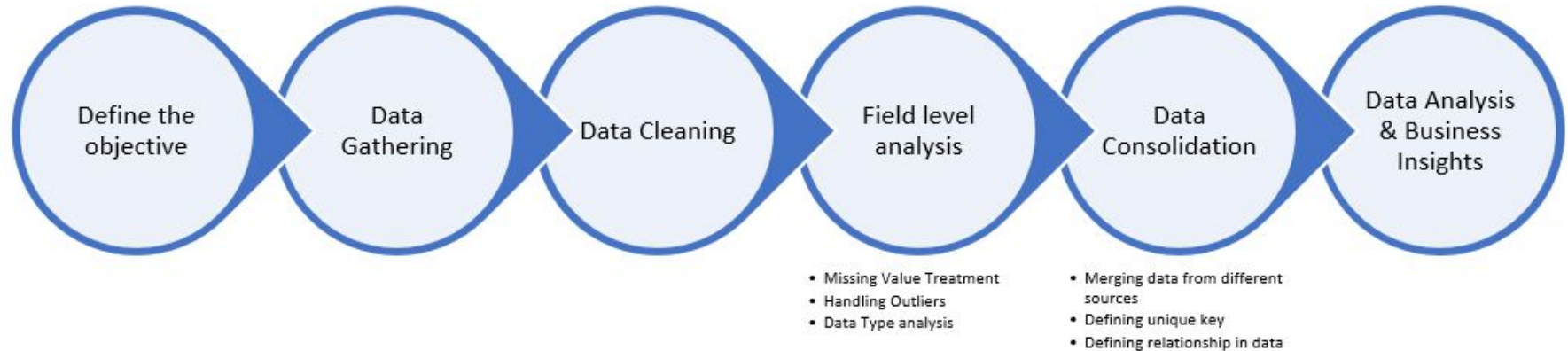
# Terminology: Data table

Data → Categorical (qualitative) → Nominal → Categories **cannot** be compared
Categorical (qualitative) → Ordinal → Categories with an implied order

Data → Numerical (quantitative)

# Data Science Lifecycle

Ref: https://medium.com/@rathi.ankit/data-science-introduction-e03773919c6

# Data Analytics: Step by Step Approach

**The Most Important Thing: Define Your Questions !!!!**



- Define the objective
- Data Gathering
- Data Cleaning
- Field level analysis
  - Missing Value Treatment
  - Handling Outliers
  - Data Type analysis
- Data Consolidation
  - Merging data from different sources
  - Defining unique key
  - Defining relationship in data
- Data Analysis & Business Insights

https://medium.com/datadriveninvestor/data-analytics-step-by-step-approach-757c6a0bd8a2

Let's explore the exercises!

# + Data Mining Tools



CHALLENGERS | LEADERS

Alteryx
SAS
RapidMiner
H2O.ai
TIBCO Software
MathWorks
Domino  IBM
Microsoft
SAP
Databricks
Anaconda  Angoss
Dataiku
Teradata

NICHE PLAYERS | VISIONARIES

ABILITY TO EXECUTE
COMPLETENESS OF VISION
As of January 2018      © Gartner, Inc

Code

Ssas

R

python™

GU

SAS® Enterprise Miner™ 12.3

WEKA
The University of Waikato

rapidminer

SPSS Modeler

# Rapid Miner

# + What is Rapid Miner?



Real Data Science, Fast and Simple

# Platform

## One Platform. Does *Everything*.

RapidMiner's unified data science platform accelerates the building of complete analytical workflows – **from data prep to machine learning to model validation to deployment** – in a single environment, dramatically improving efficiency and shortening the time to value for data science projects.

### RapidMiner Studio

Visual workflow designer for data science teams

### RapidMiner Server

Share, reuse, and deploy predictive models from RapidMiner Studio

### RapidMiner Radoop

Run data science workflows directly inside Hadoop

# Process and Operators

- An analytical workflow is called "Process"

- Each process consists of one or more "Operators"

- Connect output of an operator to input of the next operator

# Tools in Rapid Miner

# Examples of Operators

- Data accessing
  - e.g. file, cassandra, mongoDB, Amazon S3

- Data blending
  - e.g. mapping, filter, select, aggregate, split

- Data cleansing
  - e.g. normalize, deduplication, outlier

- Data modeling
  - e.g. Bayesian, decision tree, neural net

# **+** Anatomy of Operators in RapidMiner

- **"inp" – Input**
- **"out" – Output**
- **"thr" – Through**
- **"fil" – File**
- **"exp" – Example**
- **"ori" - Original**
- **"lab" – Label**
- **"tra" – Training**
- **"mod" – Model**
- **"wei" -Weight**

# Extensions in Rapid Miner: Marketplace

# Machine Learning Exercises using RapidMiner

# Lab1: Classification Lab | Pima Indians Diabetes

https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database



| Name | | Type | Missing |
|------|---|------|---------|
| Id **PatientID** | | Polynominal | 0 |
| Label **HasDiabetes** | | Binominal | 0 |
| **NumberOfPregnant** | | Integer | 0 |
| **Glucose** | | Integer | 0 |
| **BloodPressure** | | Integer | 0 |
| **SkinThickness** | | Integer | 0 |
| **Insulin** | | Integer | 0 |

| | | | |
|---|---|---|---|
| BMI | | Real | 0 |
| DiabetesPedigree | | Real | 0 |
| Age | | Integer | 0 |

# Lab2: Customer Churn



| Name | | Type |
|---|---|---|
| Label<br>**Churn** | | Polynominal |
| ⚠ **Gender** | | Polynominal |
| ⚠ **Age** | | Integer |
| **Payment Method** | | Polynominal |
| **LastTransaction** | | Integer |

# Exercise: Student Grade

| Name | | Type | Missing |
|------|---|------|---------|
| Id<br>**ID** | | Integer | 0 |
| Label<br>**IsFail** | | Binominal | 3 |
| **GPAX** | | Real | 3 |
| **Gender** | | Binominal | 0 |
| **Department** | | Polynominal | 0 |
| **AttendScore** | | Real | 5 |

34