

# Data Preparation

# Get to Know about data

**How would it be arranged for analytics?**

What is the data format?

# Data table

- Same terminology

- **Sample**

- Row
    - Instance
    - Record
    - Observation
    - Example

- **Variables**

- Features
    - Attribute
    - Field
    - Column
    - Dimension

A diagram illustrating a data table. The table has 5 columns and 4 rows. The columns are labeled 'ID', 'Date', 'MinTemp', 'MaxTemp', and 'Rainfall'. The rows are labeled with IDs 1, 2, 3, and 4. A blue bracket above the columns is labeled 'Variables'. A blue bracket to the left of the rows is labeled 'Samples'.

ID	Date	MinTemp	MaxTemp	Rainfall
1	2010-06-17	55	75	0.1
2	2010-06-18	52	78	0.0
3	2010-06-19	50	78	0.0
4	2010-06-20	54	77	0.0

# Variable types

- **DATE**
- **String**
- **Categorical data**
  - Color
  - Gender
  - Product Category
- **Number**
  - Temperature
  - Rainfall
  - Height
  - Age

A diagram illustrating a data table structure. A blue bracket labeled "Variables" spans the top of the table, grouping the column headers. Another blue bracket labeled "Samples" is positioned to the left of the table, grouping the rows. The table itself has a light blue header row and four data rows.

ID	Date	MinTemp	MaxTemp	Rainfall
1	2010-06-17	55	75	0.1
2	2010-06-18	52	78	0.0
3	2010-06-19	50	78	0.0
4	2010-06-20	54	77	0.0

# Get to Know about data

**Will data be in perfect form?**

Is there any missing or errorless data?

# What's wrong?

ID	Date	MinTemp	MaxTemp	Rainfall
1	2010-06-17	56	24	0.1
2	2016-06-18	52	26	3,678.9
3	2010-06-19	50	26	0.0
4	2010-06-20	54	25	0.0

ID	Date	MinTemp	MaxTemp	Rainfall
1	2010-06-17	56	75	--
2	2016-06-18	52	78	--
3	2010-06-19	--	78	0.1
4	2010-06-20	54	77	--

# What's wrong?

- Incomplete Data

- Missing Data
- Duplicate Data
- Different measure
- Different scaling
- Invalid
- Noisy
- Outliers

Name	Address
Angela	430 Park Drive
Sidney	7800 West View Street
Sid	7800 West View Street
Ratan	12442 Mountain Avenue
Kiril	45 East 5 <sup>th</sup> St
Kiril	1220 Mill Avenue
Zhou	4345 Apple Lane

# What's wrong?

- **Incomplete Data**

- Missing Data
- Duplicate Data
- Different measure
- Different scaling
- **Invalid**
- Noisy
- Outliers

Name	Zip Code
Angela	346412
Sidney	92618
Ratan	8033A
Kiril	11012
Zhou	59285



# What's wrong?

- Incomplete Data

- Missing Data
- Duplicate Data
- Different measure
- Different scaling
- Invalid
- Noisy
- Outliers

Name	Address
Angela	430 Park Drive
Sidney	780 ★❖©◆ View Street
Ratan	12443 Mountain Avenue
Kiril	1220 Mill Avenue
ZhČou	4345 Apple Lane

# Data Cleaning



# Data visualization

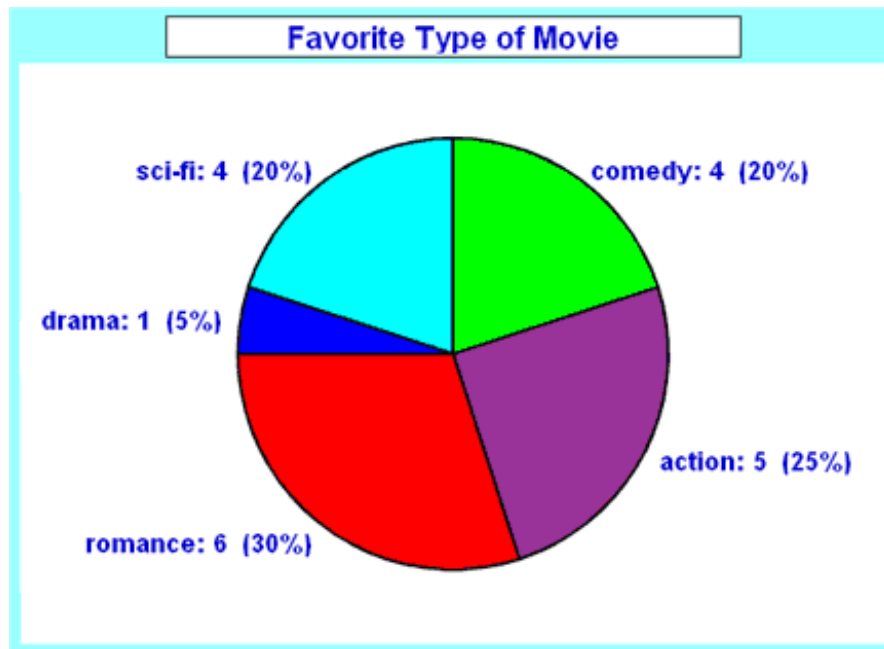
**How would we visualize the data?**

Is the data perfect or noisy?

# Visualize data

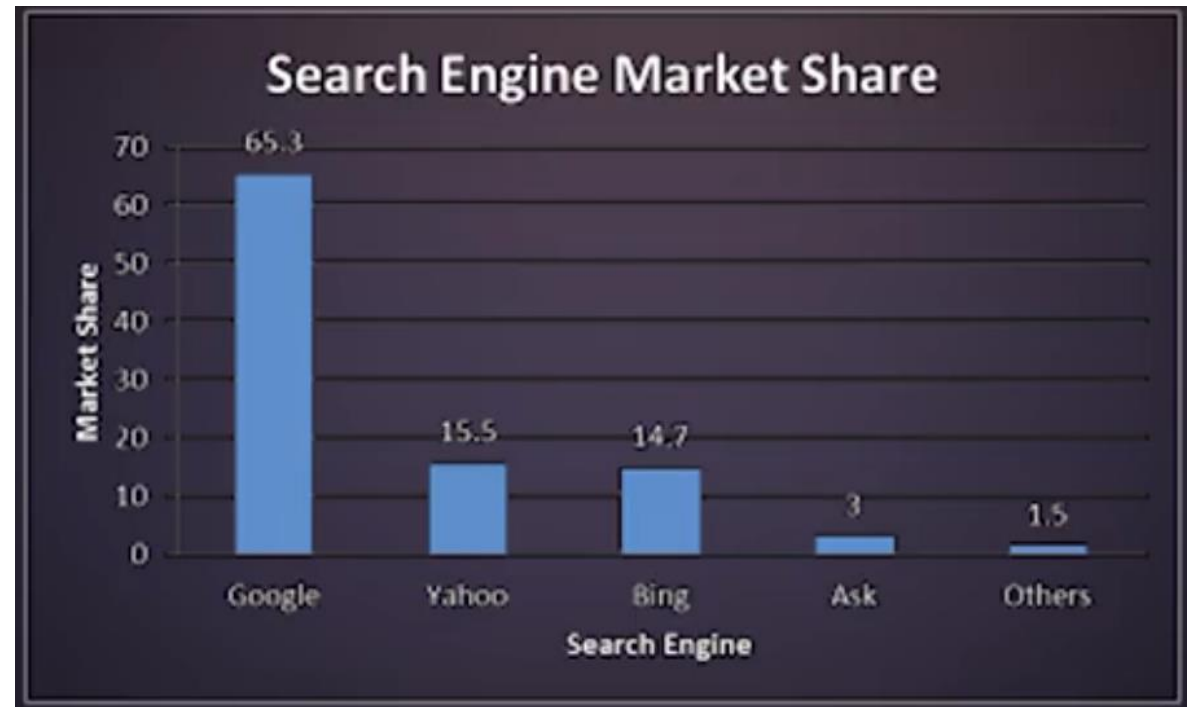
- **Pie Plot**

- represent the percentage share of different categories
  - implicit and assumed that the category is exhaustive



- **Bar Plot**

- category vs the count or percentage of each category
  - Show top category

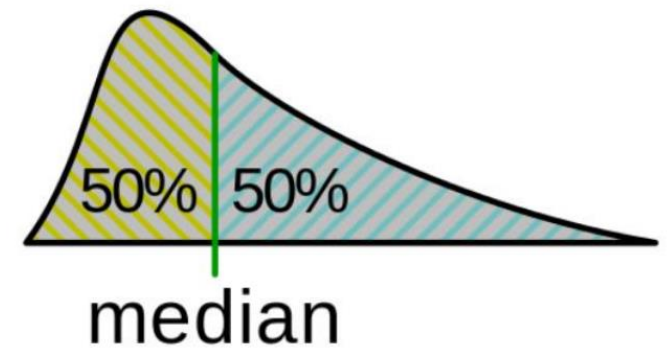
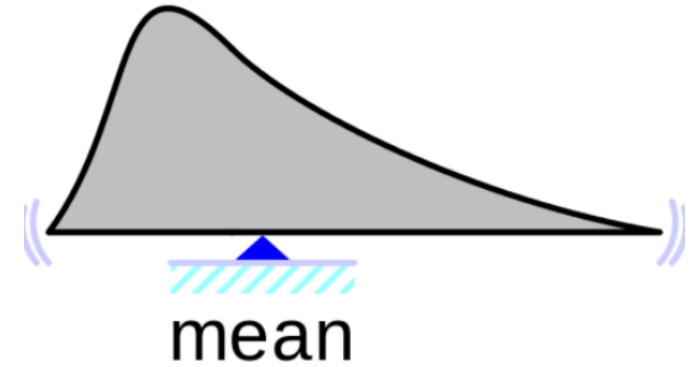
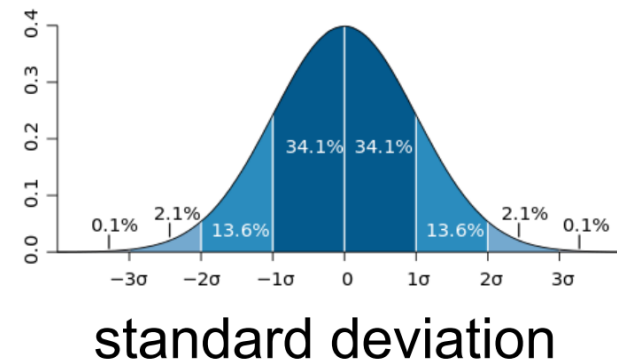
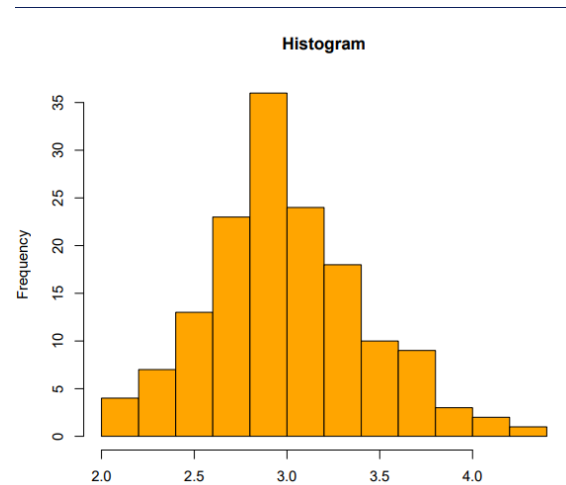


# Visualize data

- **Statistical View**

- **Histogram**

- Distribution of Data
    - Probability
    - Statistical values
      - Mean
      - Standard deviation

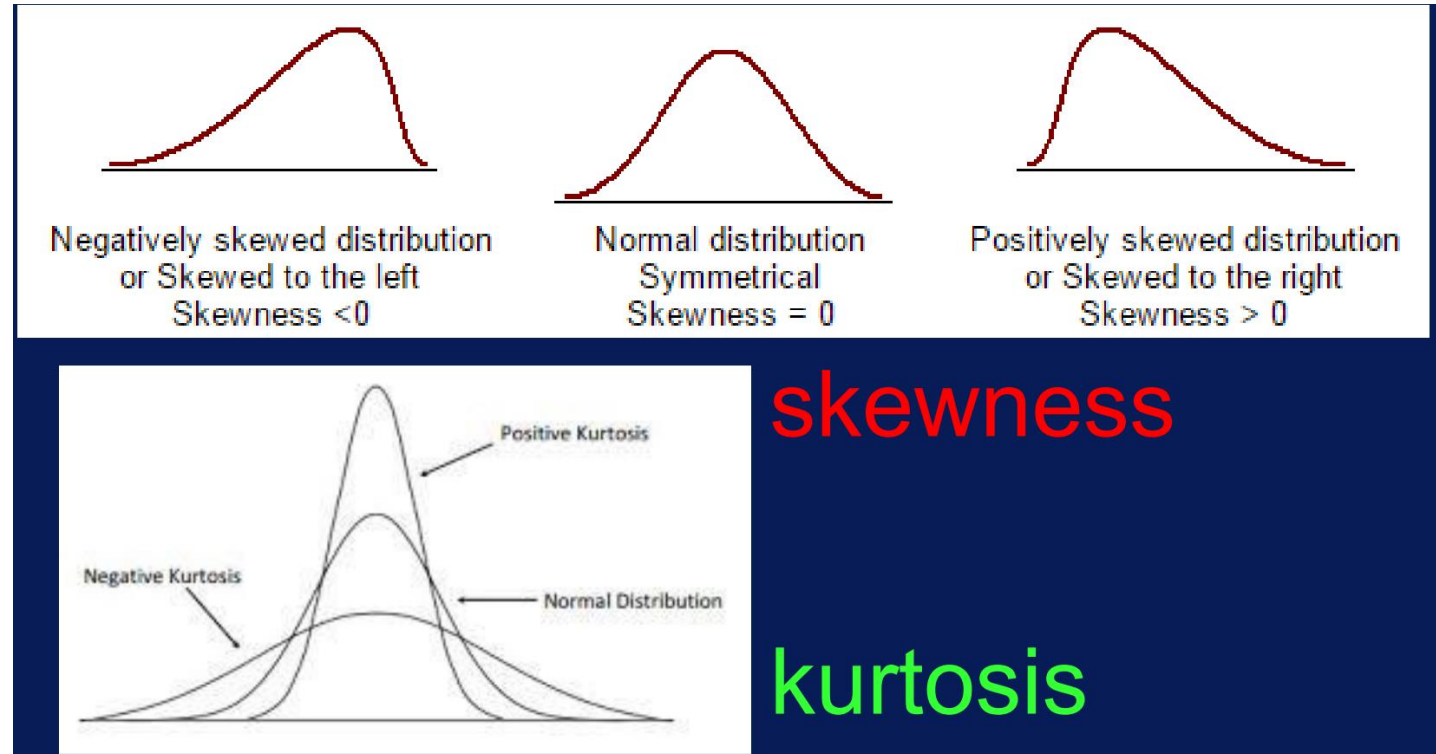


# Visualize data

- **Statistical View**

- **Histogram**

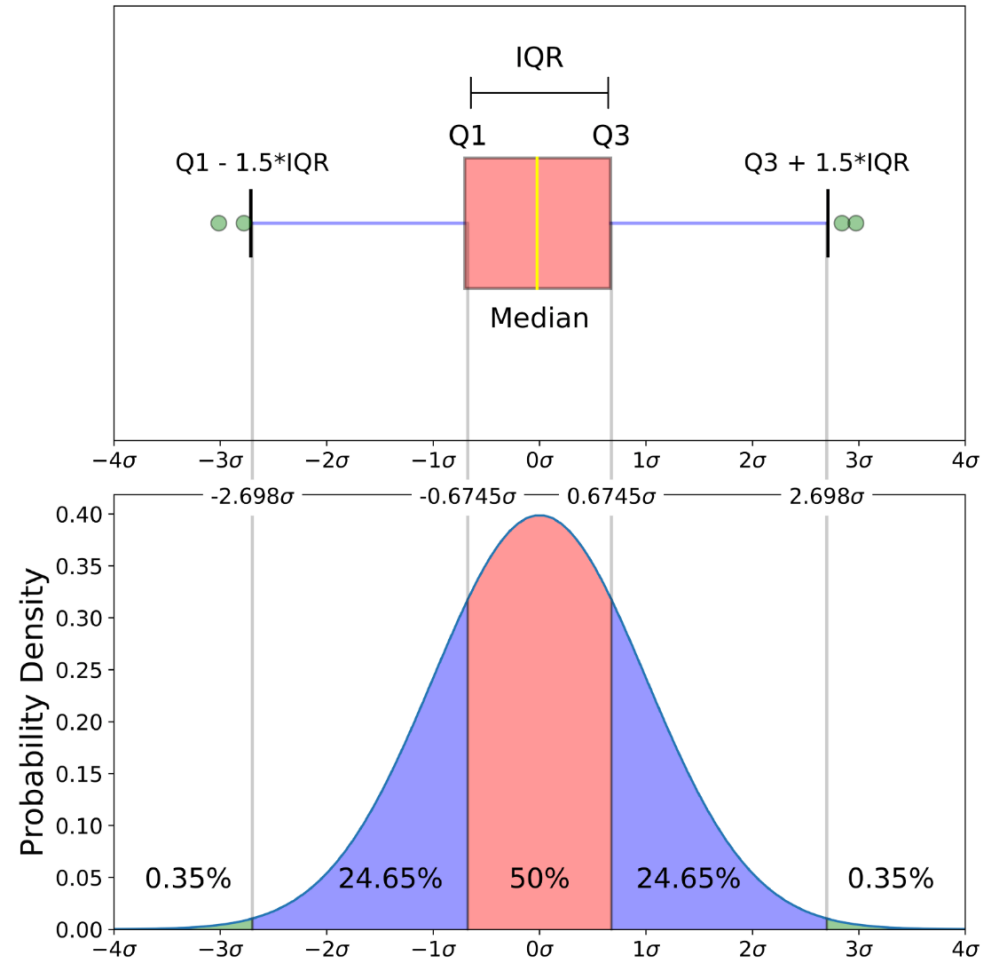
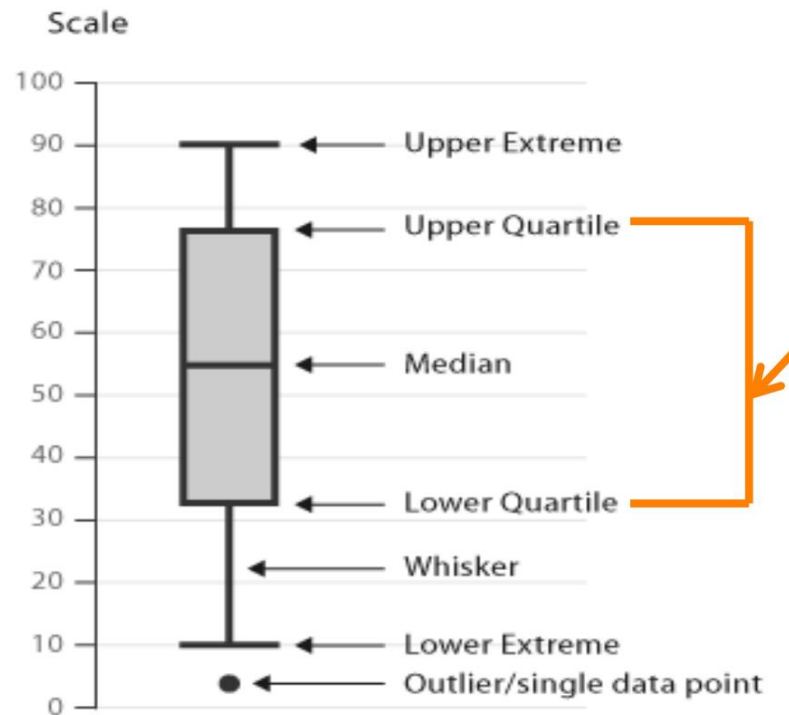
- Distribution of Data
    - Probability
    - Statistical values
      - Skewness
      - kurtosis



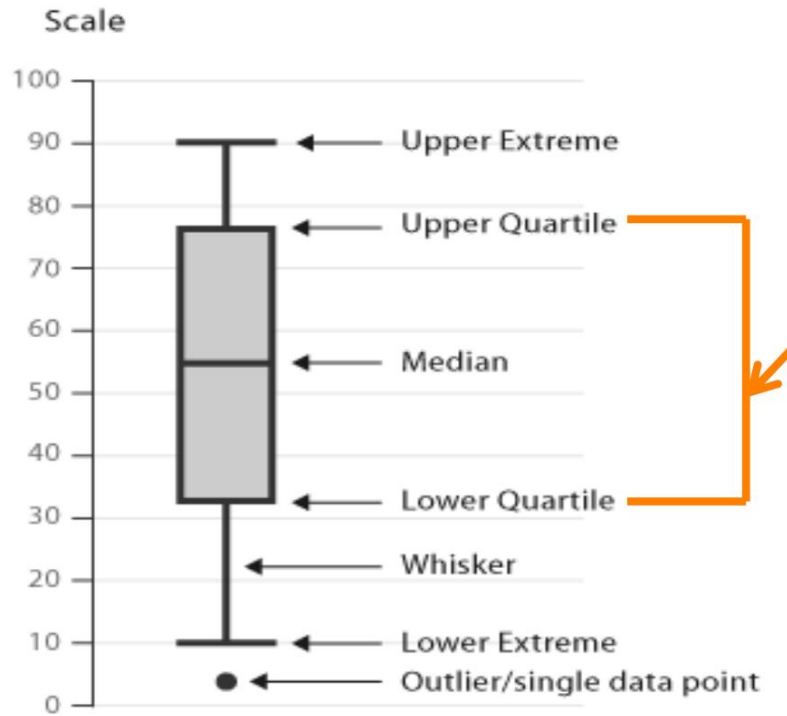
# Visualize data

- **Box Plot**

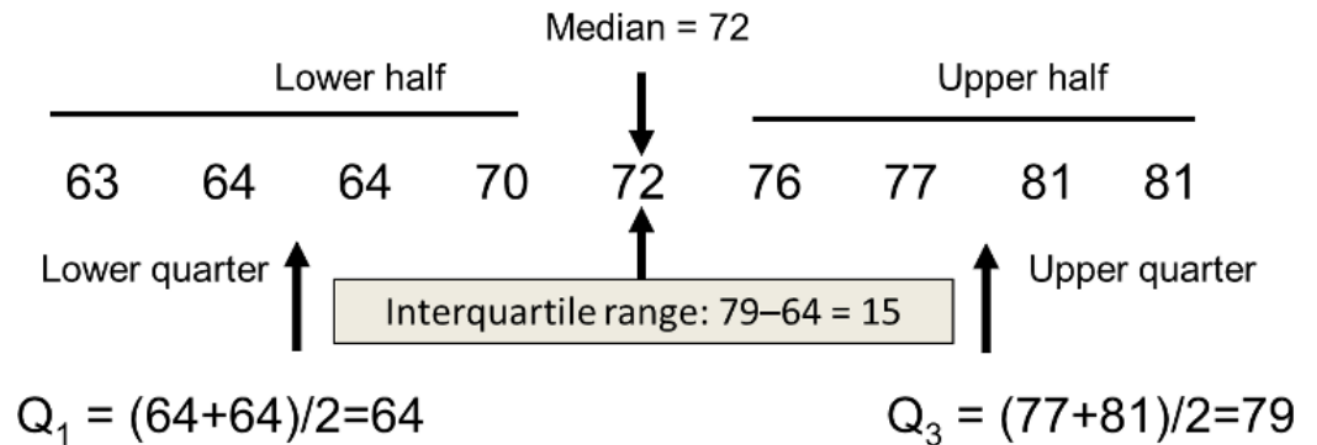
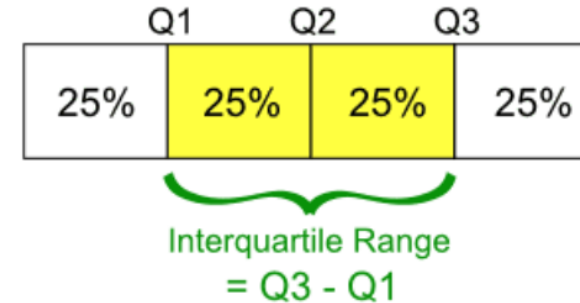
- Distribution of data



# alize data

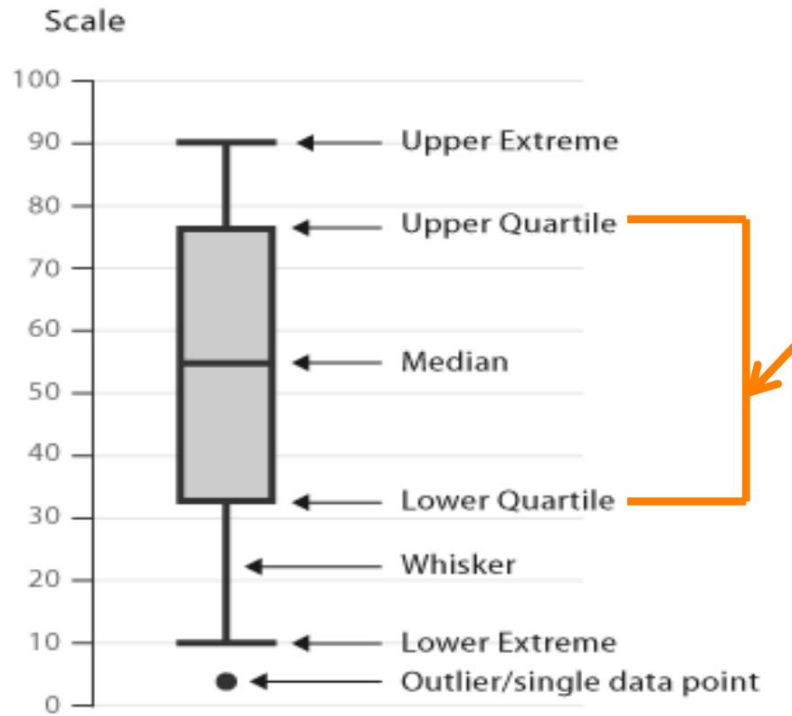


- Relationship of box plot vs distribution of data





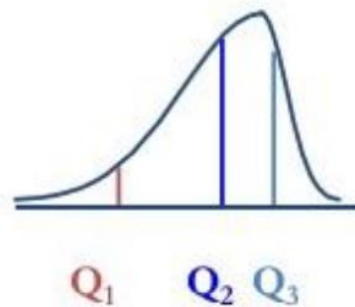
# alyze data



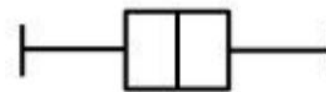
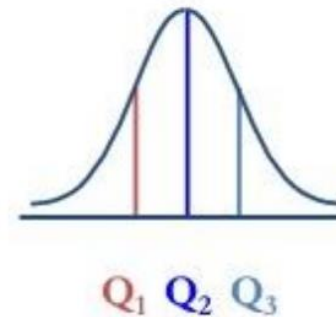
- Relationship of box plot vs distribution of data

## Distribution Shape and The Boxplot

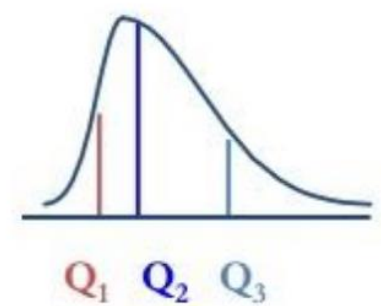
Negative Skew



Symmetric



Positive Skew



# Visualize data

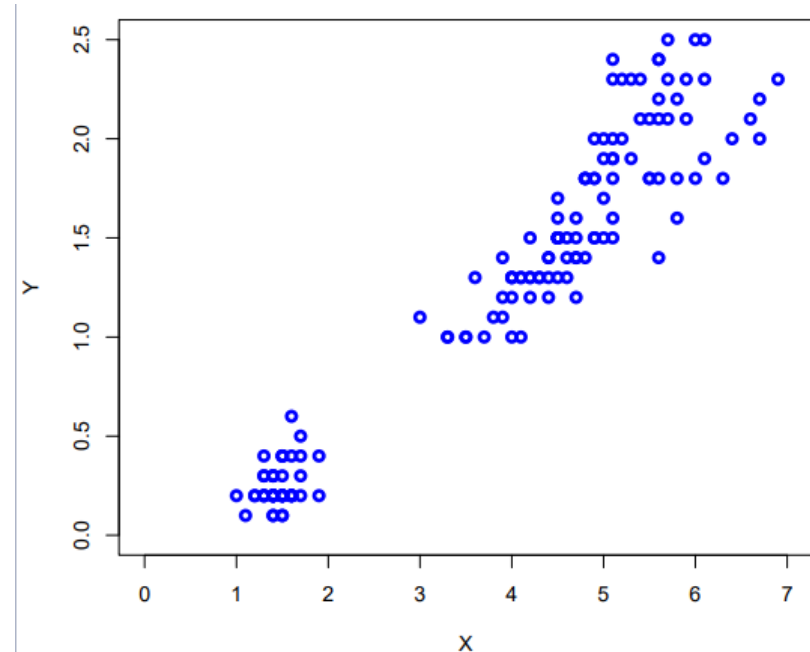
- **Line Plot**

- Sequential data (Time series data)
  - Trends Prediction



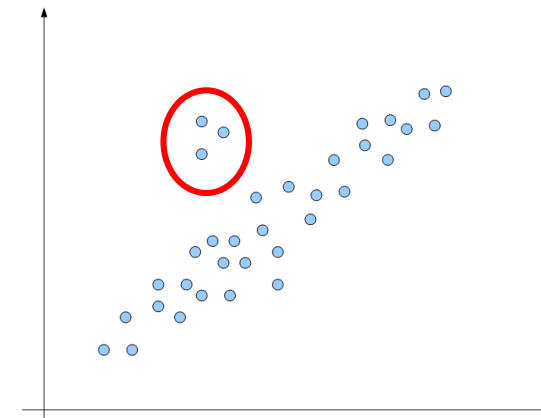
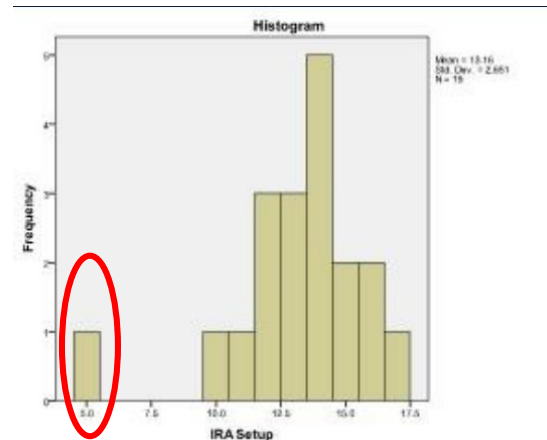
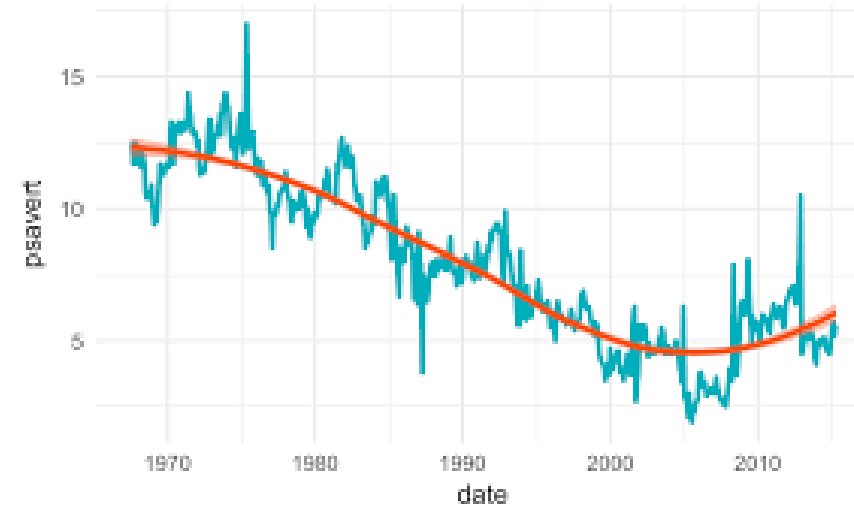
- **Scatter Plot**

- Relationship between 2 variables
  - Correlation



# Visualize data

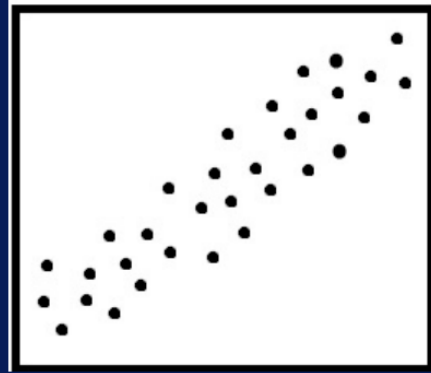
- Trend
- Outlier



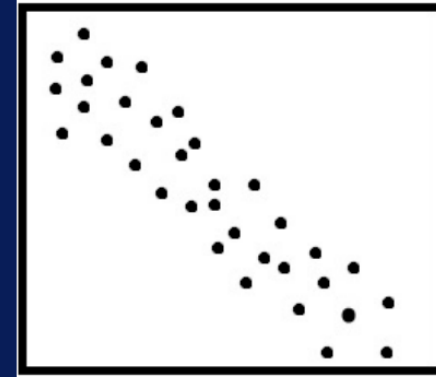
# Visualize data

- Scatter plots of two variables (features)
- Show Correlation between them

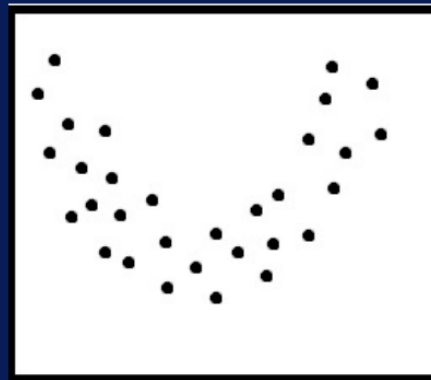
Positive  
Correlation



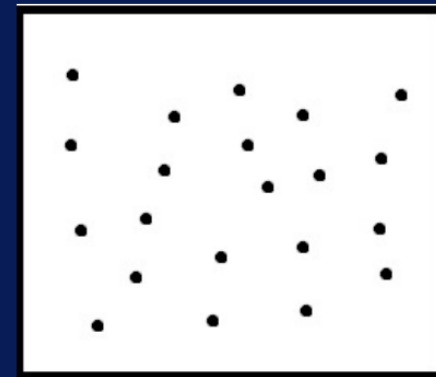
Negative  
Correlation



Non-  
Linear  
Correlation

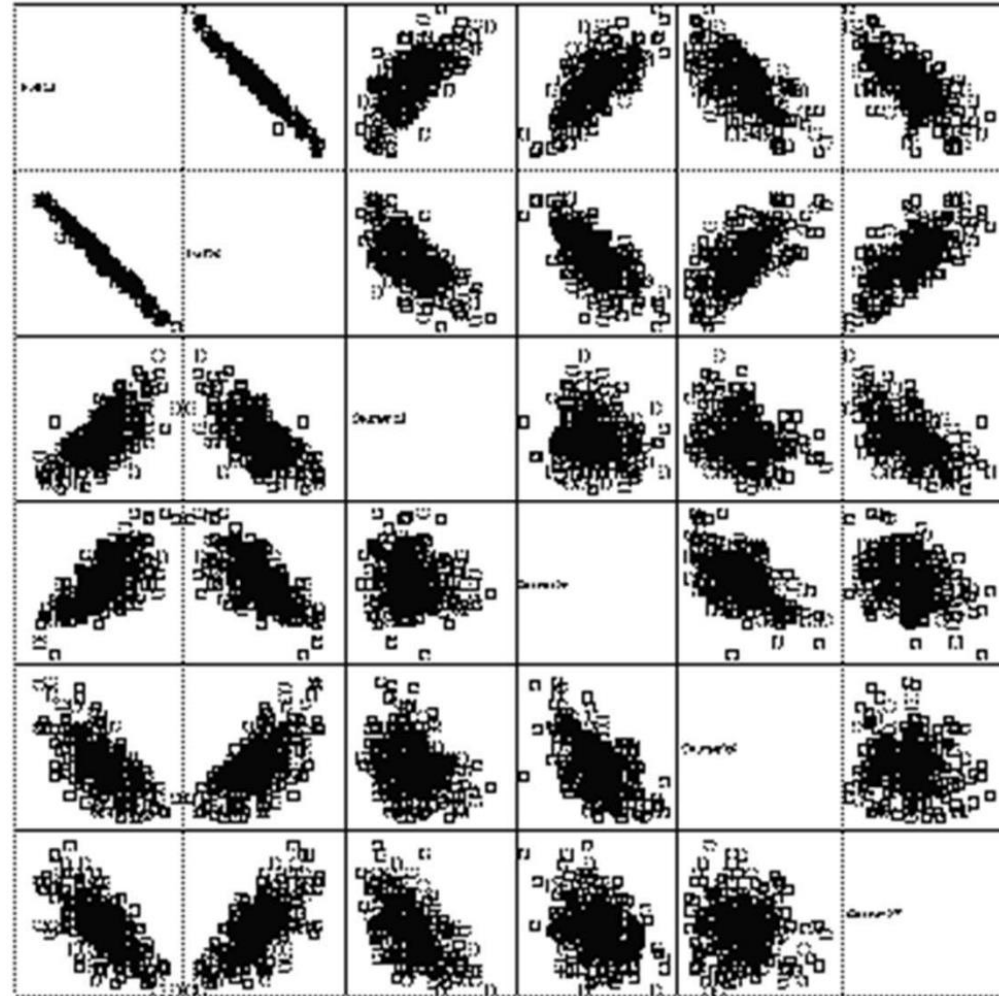


No Correlation



# Visualize data

- Scatter plots of multi-variables (features)
- Show Correlation between a pair of them



# How to make data valid

**Should data be changed or transformed?**

How to correct any missing or errorless data?

# Data Cleaning

- Incomplete Data

- Missing Data
- Duplicate Data
- Different measure
- Different scaling
- Invalid
- Noisy
- Outliers

Name	Age	Income
Angela	34	80
<del><i>Sidney</i></del>	<del>--</del>	<del><i>56</i></del>
<del><i>Ratan</i></del>	<del><i>10</i></del>	<del>--</del>
<del><i>Kiril</i></del>	<del><i>68</i></del>	<del>--</del>
Zhou	45	120

Problem Solving: Remove all records containing missing data

Good point: Simple

Critical point: may not have enough data

# Data Cleaning

- Incomplete Data

- Missing Data
- Duplicate Data
- Different measure
- Different scaling
- Invalid
- Noisy
- Outliers

Name	Age	Income
Angela	34	80
Sidney	0	56
Ratan	10	0
Kiril	68	0
Zhou	45	120

Problem Solving: Replace missing data with 0

Good point: Simple

Critical point: Lead to mislead and become outlier

May not be allowed especially for regression application



# Data Cleaning

- Incomplete Data

- Missing Data
- Duplicate Data
- Different measure
- Different scaling
- Invalid
- Noisy
- Outliers

Name	Age	Income
Angela	34	80
Sidney	39	56
Ratan	10	80
Kiril	68	80
Zhou	45	120

Problem Solving: Replace missing data with some values

Ex. Mean / Median / Most occurrence / History value /  
Regression prediction / center of the cluster

Good point: Better representation

Critical point: Need Historical Data and Higher Computation

# Data Cleaning

- Incomplete Data

- Missing Data
- Duplicate Data
- Different measure
- Different scaling
- Invalid
- Noisy
- Outliers

Name	Address
Sidney	7800 West View Street
<del>Sid</del>	<del>7800 West View Street</del>
<del>Kiril</del>	<del>45 East 5<sup>th</sup> St</del>
Kiril	1220 Mill Avenue

Problem Solving: Delete Old record -> Replace with New one  
or by defined rules

# Data Cleaning

- Incomplete Data
  - Missing Data
  - Duplicate Data
  - Different measure
  - Different scaling
- Invalid
- Noisy
- Outliers



**Problem Solving:**

Ex: Distribution Shifting and Scaling  
Normalization / Log transform

**Note:** scaling method selection depends on the domain problem.

# Data Cleaning

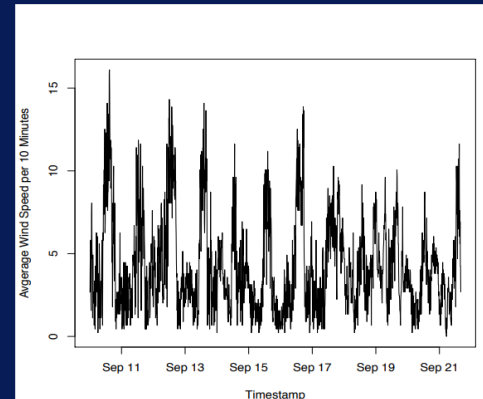
- Incomplete Data
  - Missing Data
  - Duplicate Data
  - Different measure
  - Different scaling
  - Invalid
  - Noisy
  - Outliers

## Problem Solving:

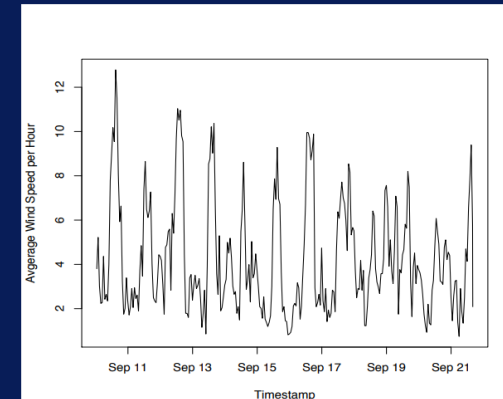
Ex: Removing / Filtering

Note: Removing carefully.

Avg Wind Speed  
(every 10 minutes)



Avg Wind Speed  
(every 60 minutes)



# Data Cleaning

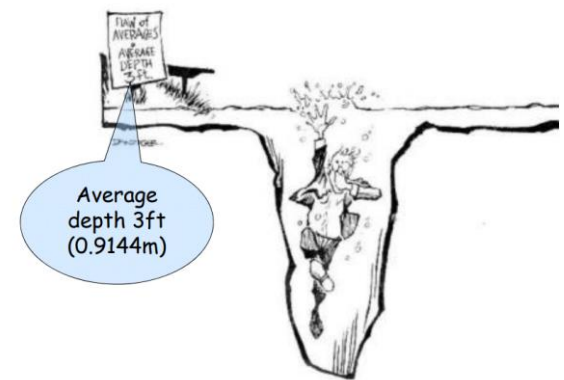
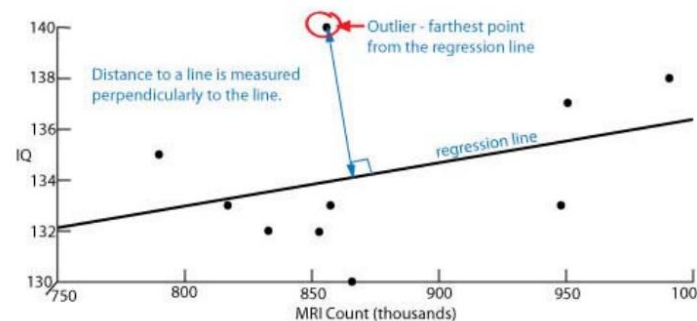
- Incomplete Data
  - Missing Data
  - Duplicate Data
  - Different measure
  - Different scaling
  - Invalid
  - Noisy
  - Outliers

## Problem Solving:

Ex: Removing

Replacing with specify value  
such as most frequent value or  
trend of its group

Note: Removing carefully.



[https://paginas.fe.up.pt/~ec/files\\_1112/week\\_02\\_descriptive\\_statistics.pdf](https://paginas.fe.up.pt/~ec/files_1112/week_02_descriptive_statistics.pdf)

# Exercise

**[13, -, 15, 16, 16, -, 19, 20, 21, 0, 22, 25, 0, 33, 1000, 35, 36, 40, -, 45]**

**Perform Data Cleaning**

**Calculate**

- sorting from min to max**
- mean / median / min / max**
- normalized value in range of 0 to 1 ( [0,1] )**
- z-normalized value**



