

# Data Engineering

MG-GY 8441

# Describing Data

- Agenda
  - Overview
  - Data Types
  - Numerical Summaries
  - Visualizations
- References
  - Han, Kamber, Pei, *Data Mining: Concepts and Techniques*
    - Chapter 1
    - Chapter 2.1 - 2.3

# Overview



# Overview

- Why Data Mining?
- What Is Data Mining?
- What Kind of Data Can Be Mined?
- What Technology Are Used?
- What Are the Applications?
- What Are the Challenges?

# Why Data Mining?

## Classification

- Classify credit applicants as low, medium, high risk

## Estimation

- Estimate the **click-through-rate** of an advertisement

## Prediction

- Predict which customers will leave within six months

# Why Data Mining?

## **Example from Marketing**

- You are in a meeting with your boss and a large publisher where you are negotiating to buy some advertising on their website.
- The publisher tells you the cost per thousand views of your advertisement (CPV) is \$10.
- Given your goal of collecting email addresses for potential new customers, you need to know the maximum CPV you can afford to effectively negotiate with the publisher.

# Why Data Mining?

## **Example from Marketing**

- You estimate that the click-through rate (CTR) of your advertisement has been around 1%.
- Your conversion rate (CR) has been averaging 10% in terms of email sign-ups.
- If you can afford to pay \$5 per email you acquire as your cost per acquisition (CPA), is \$10 CPV a good price from the publisher?

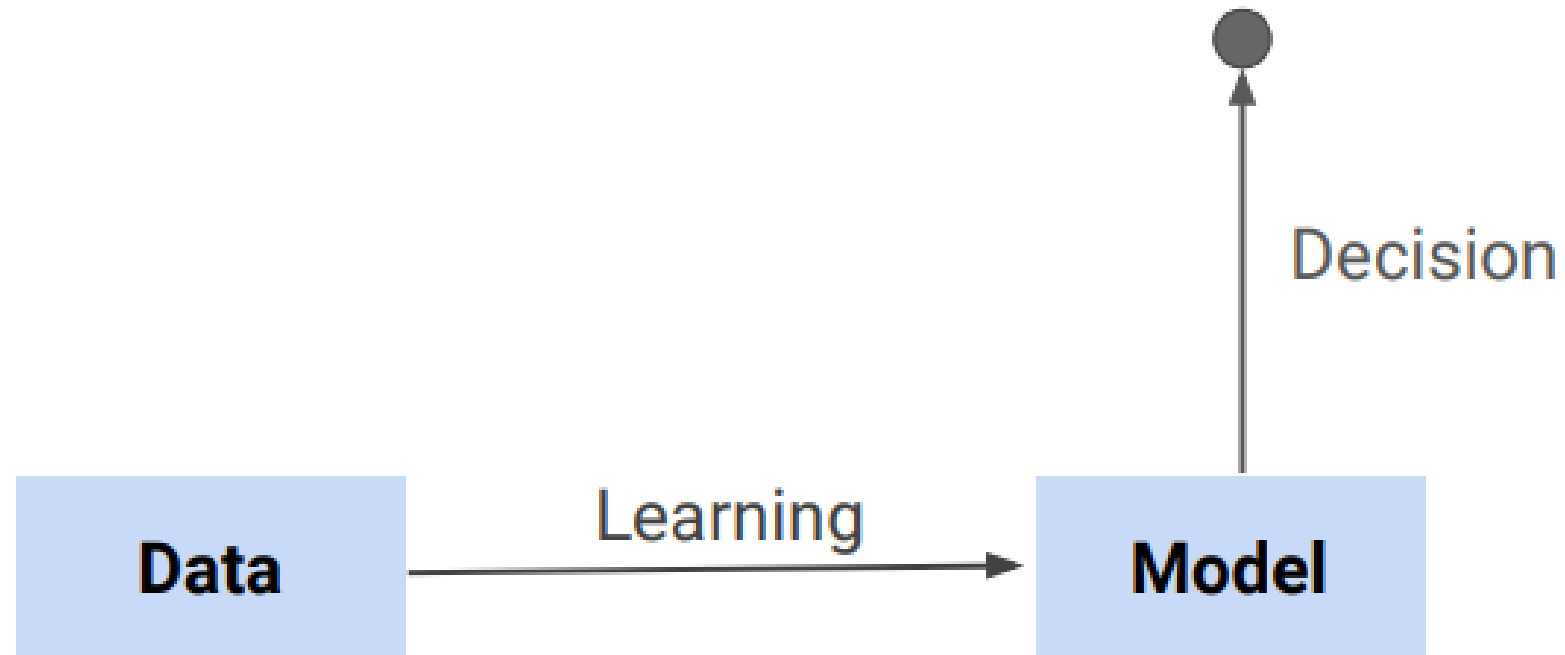
# Why Data Mining?

## Example from Marketing:

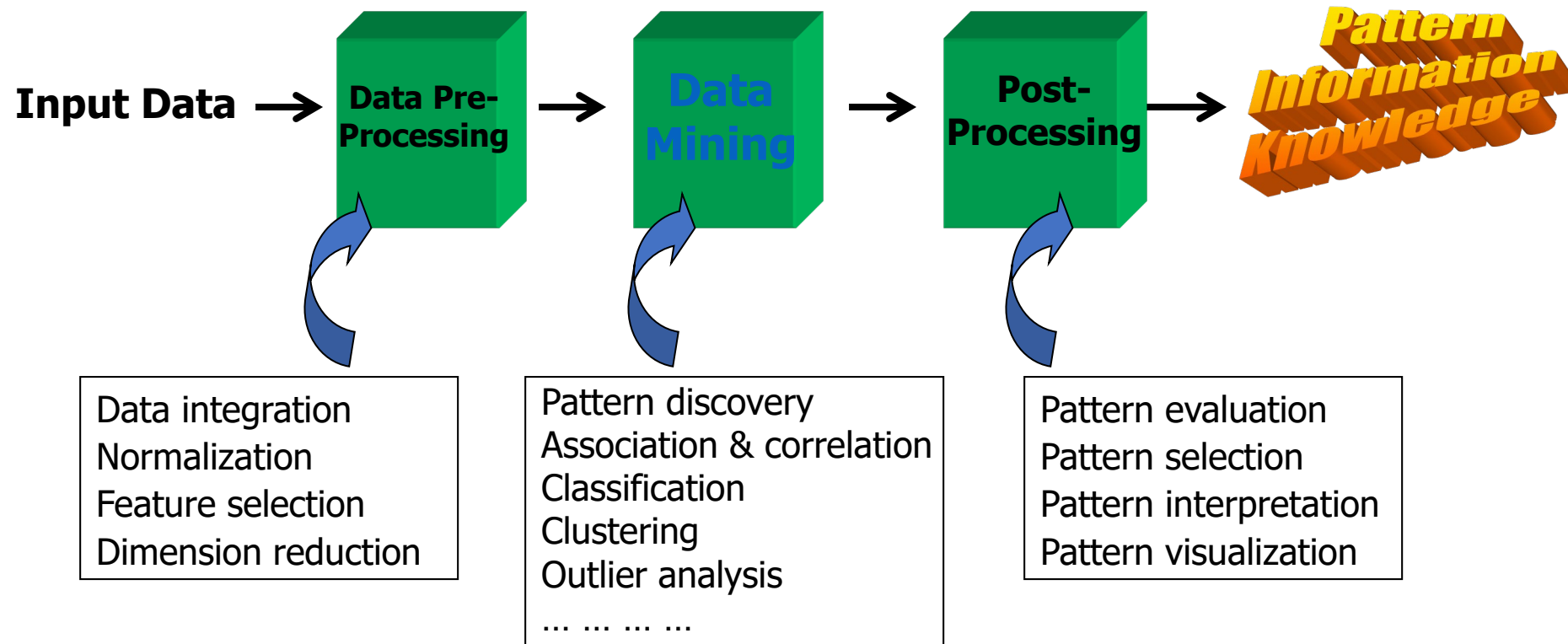
- CPV = \$10, CTR = 1%, CR = 10%, CPA goal = \$5
  - $CPA = CPV / ((CTR * 1000) * CR)$
  - $CPA = \$10 / ((0.01 * 1000) * 0.1)$
  - $CPA = \$10 / (10 * 0.1)$
  - $CPA = \$10 / 1$
- So the cost per acquisition goal would need to be doubled to match the publisher's price of \$10 for cost per thousand views



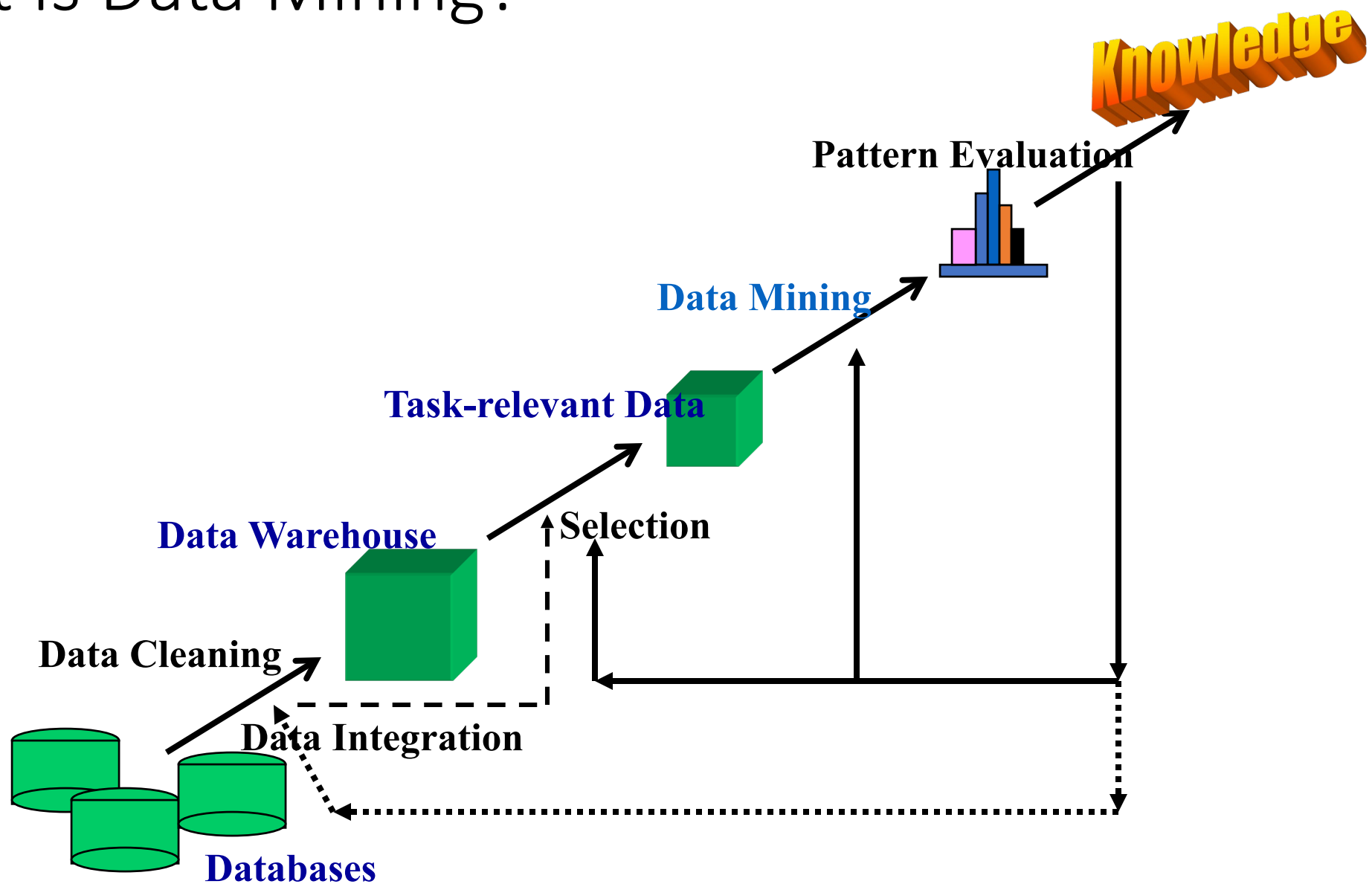
# What Is Data Mining?



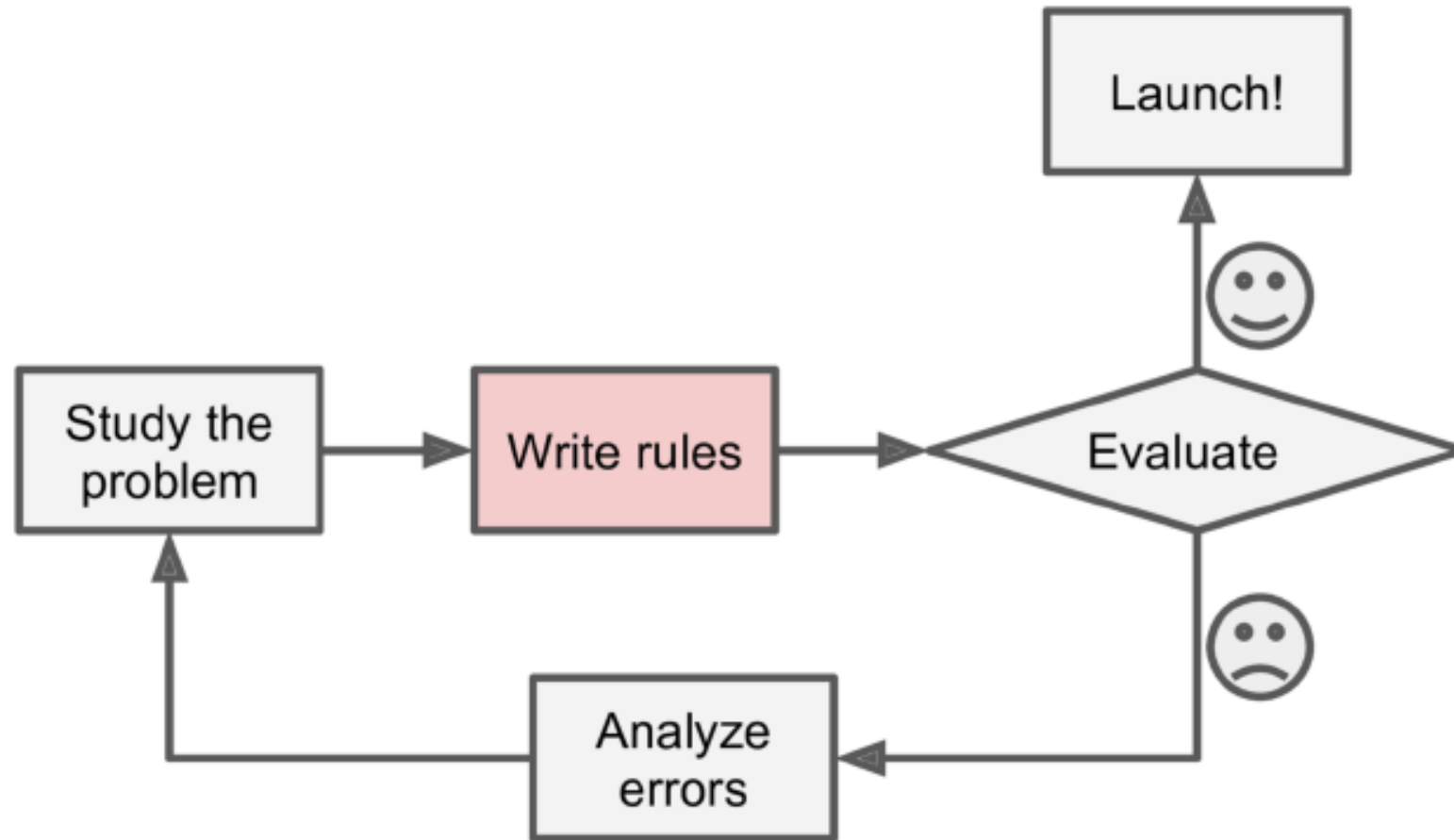
# What Is Data Mining?



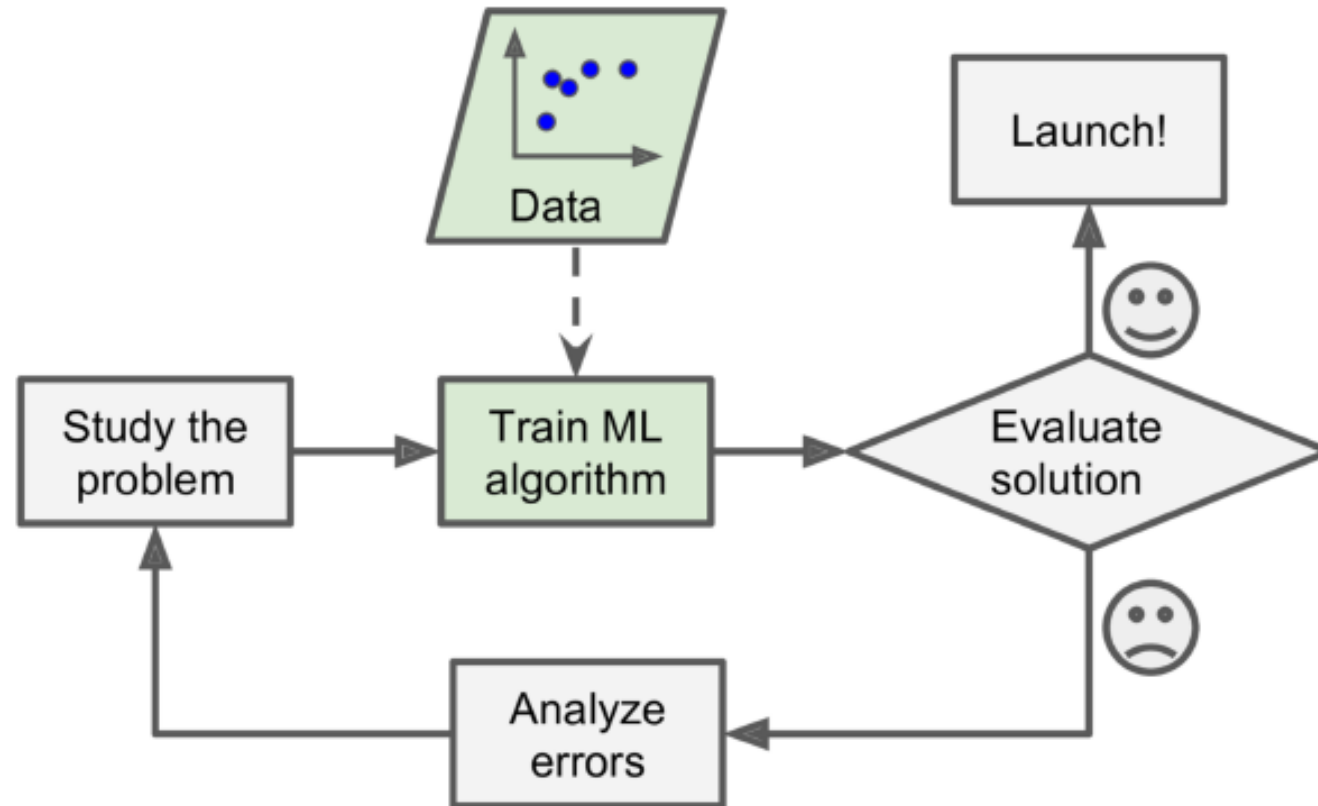
# What Is Data Mining?



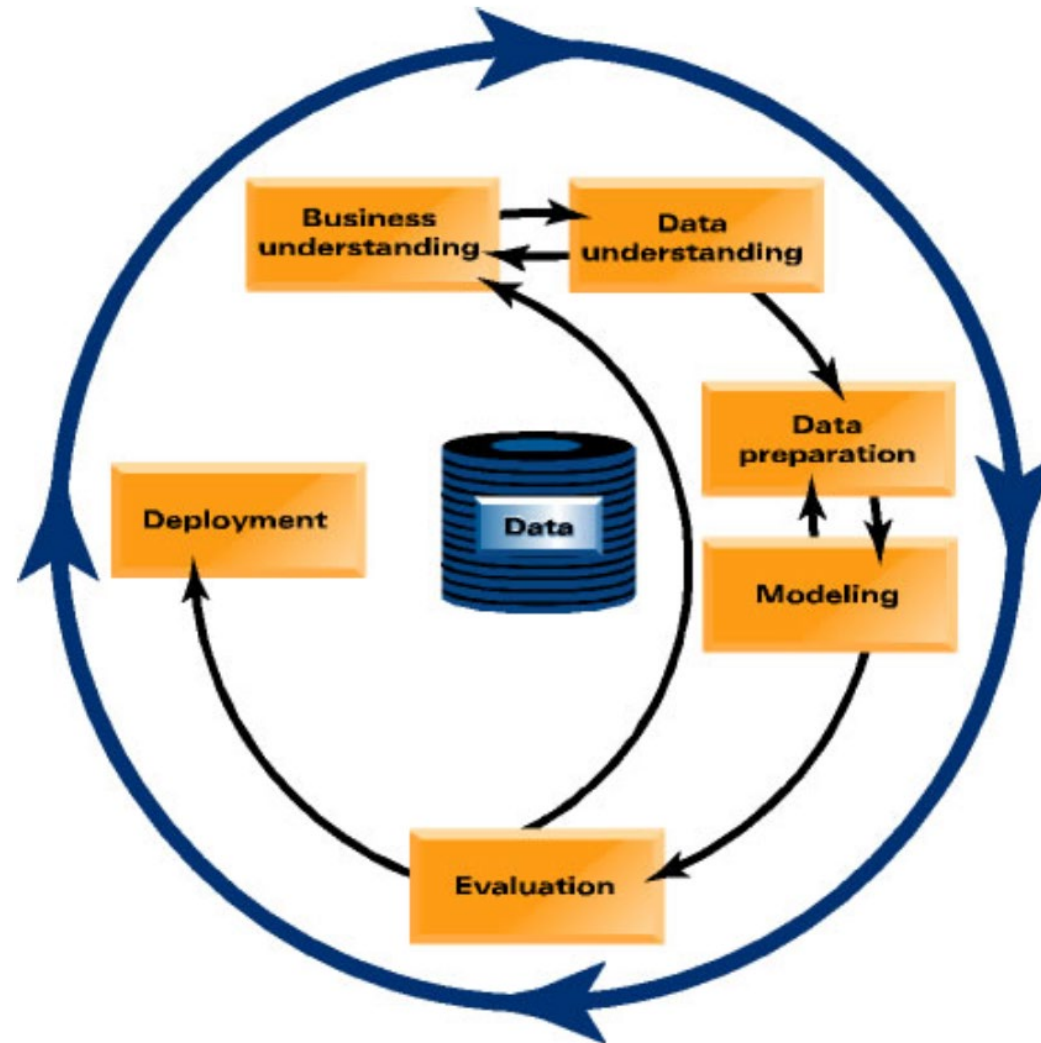
# What Is Data Mining?



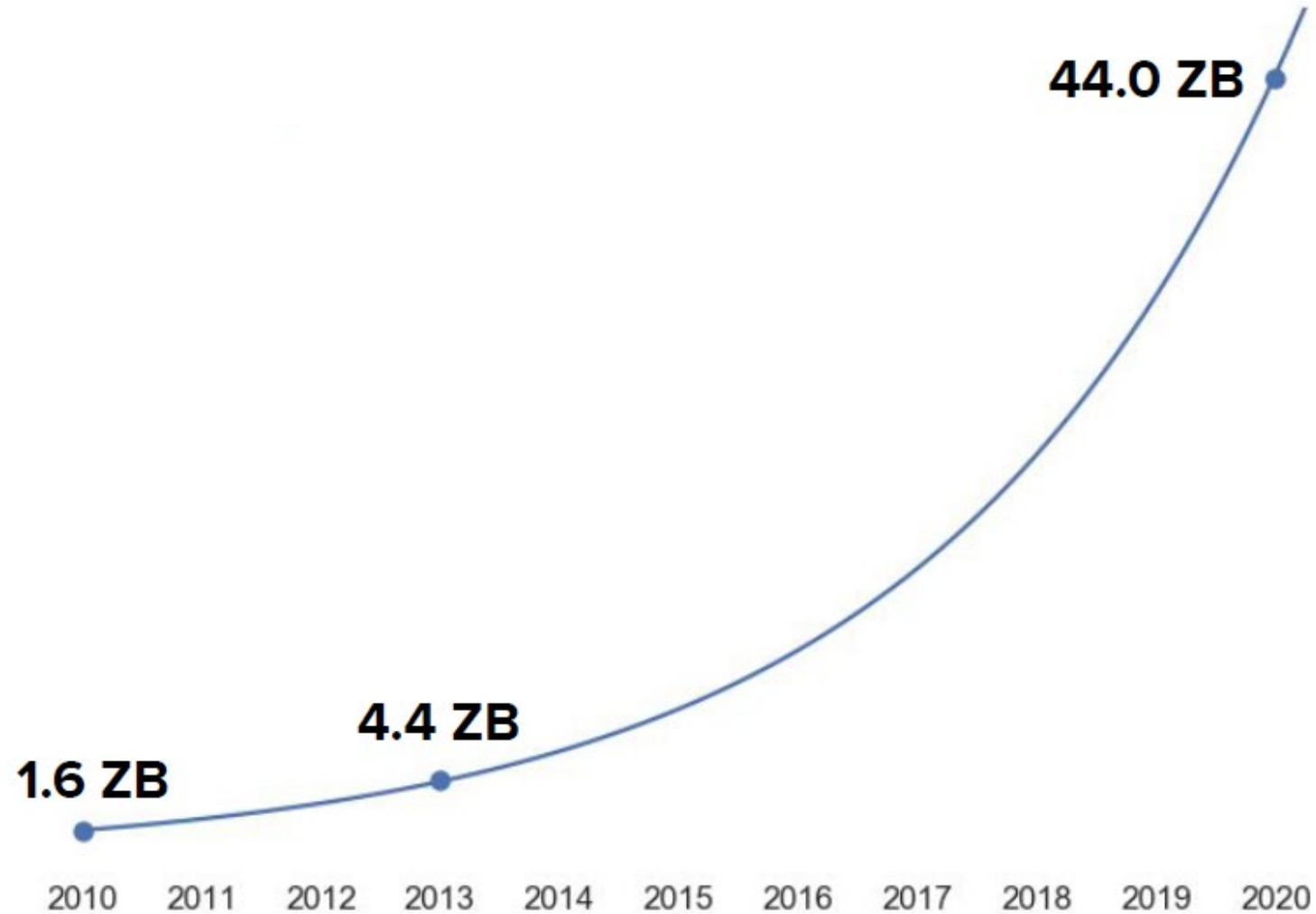
# What Is Data Mining?



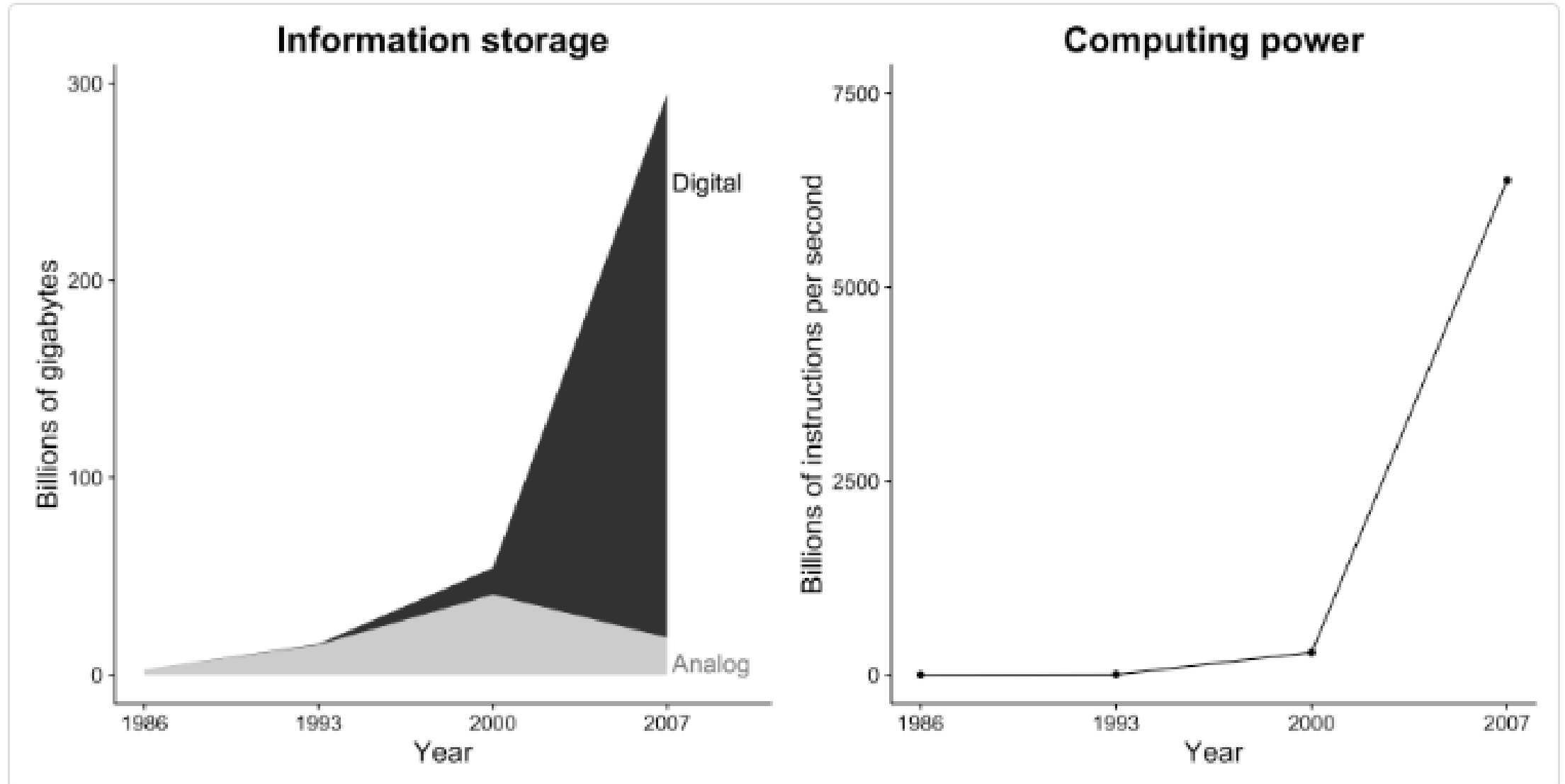
# What Is Data Mining?



# What Kind of Data Can Be Mined?

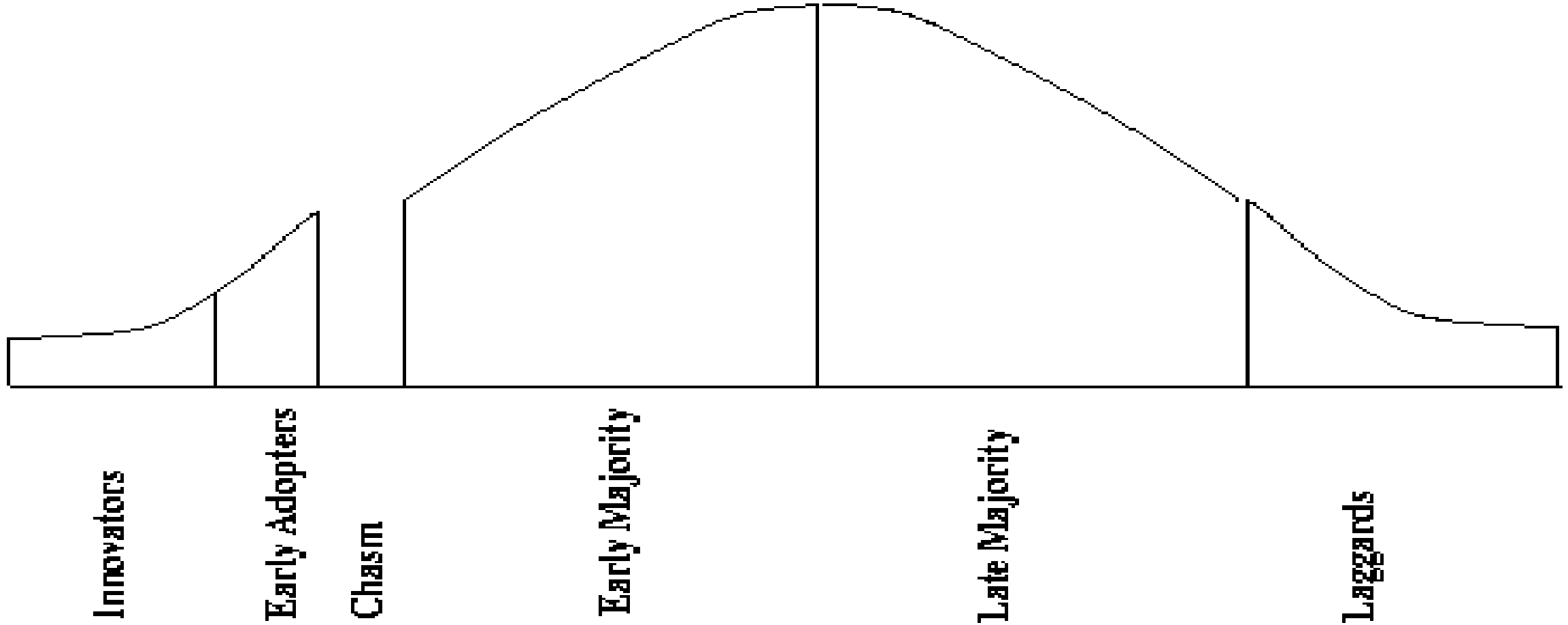


# What Kind of Data Can Be Mined?





# What Technology Are Used?



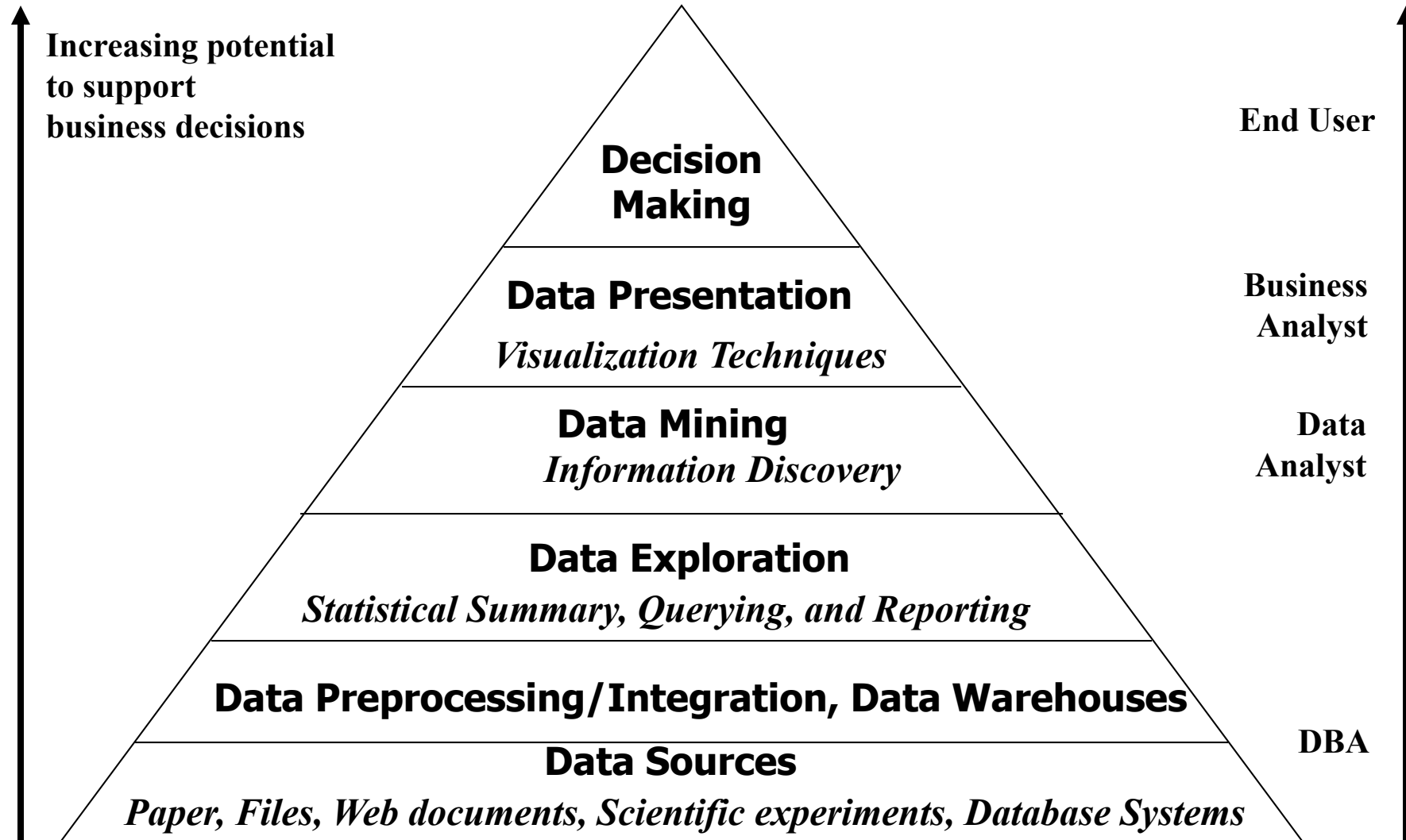
# What Technology Are Used?



# What Are the Applications?

- Market Basket Analysis
  - Identify what products are likely to be bought together
- Entity Resolution
  - Disambiguate records by linking various data sources
- Market segmentation
  - Identify common characteristics of customers who buy same products
- Collaborative Filtering
  - Recommend products to customers based on preferences

# What Are the Applications?



# What Are the Challenges?

- Privacy
  - Right to be unknown or forgotten
- Transparency
  - Redistribution of data
- Accountability
  - Oversight of companies and government agencies
- Fairness
  - Social impact of data driven decision making

# Data Types



# Data Types

- Categorical
  - Ordinal
  - Nominative
- Numerical
  - Continuous
  - Discrete

# Properties of Data

|                  |   |
|------------------|---|
| <b>Volume</b>    | The quantity of data                    |
| <b>Velocity</b>  | Speed at which data is collected        |
| <b>Variety</b>   | Data may be structured or heterogeneous |
| <b>*Veracity</b> | Data can be noisy, incomplete, or wrong |



# Properties of Data

- Dimensionality
  - High dimensional data needs to transform to low dimension data for processing
- Sparsity
  - If many entries lack information, then we need to compress the data to improve storage and computation
- Resolution
  - The granularity of a dataset refers to the level of detail. Sometimes we need to adjust the granularity by grouping records.
- Distribution
  - Aggregations of numbers allow use to track trends in data like dispersion around a common center

# Data Object

- Datasets are made up of data objects. A **data object** represents an entity amongst the records.
- Example:
  - sales database: customers, store items, sales...
- Data objects are described by **attributes**.
  - We tend to organize data into tables with rows and columns
    - rows -> data objects
    - columns -> attributes.

# Qualitative Data

- **Nominal:** categories, states, or “names of things”
  - Example: marital status, occupation, zip codes
- **Ordinal**
  - Values have a meaningful order (ranking) but magnitude between successive values is not known.
  - Example: *small, medium, large*}

# Quantitative Data

- **Discrete Attribute**

- We can count the distinct numbers.
  - Commonly the numbers have integer values like 0,1,2,3,..
- Binary attributes are a special case of discrete attributes containing the value 0 or 1

- **Continuous Attribute**

- We cannot count the range of numbers
  - Commonly the numbers have floating point values containing fractions with many digits

# Quantitative Data

- **Interval**

- Range of numbers differing by increments
- Example: 1,2,3,4,5 contains numbers differing by increment of 1

- **Ratio**

- Range of numbers differing by factor
- Example:
  - $10^1, 10^2, 10^3, 10^4, 10^5$  contains numbers differing by factor of 10
  - Factors of 10 are called order of magnitude. This type of range is a logarithmic scale. It can be use for working with large ranges of values in visualizations

# Numerical Summaries



# Numerical Summaries

- Measuring Central Tendencies
  - Mean
  - Median
  - Mode
- Ranking Numbers
  - Quantiles

# Measuring the Central Tendency

- Mean
  - Measures the average of a collection of numbers.
  - If we want to put different emphasis on the numbers, then we can use weighted mean
  - If we want to remove outliers, then we could discard or round some values

$$\bar{x} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$



# Measuring the Central Tendency

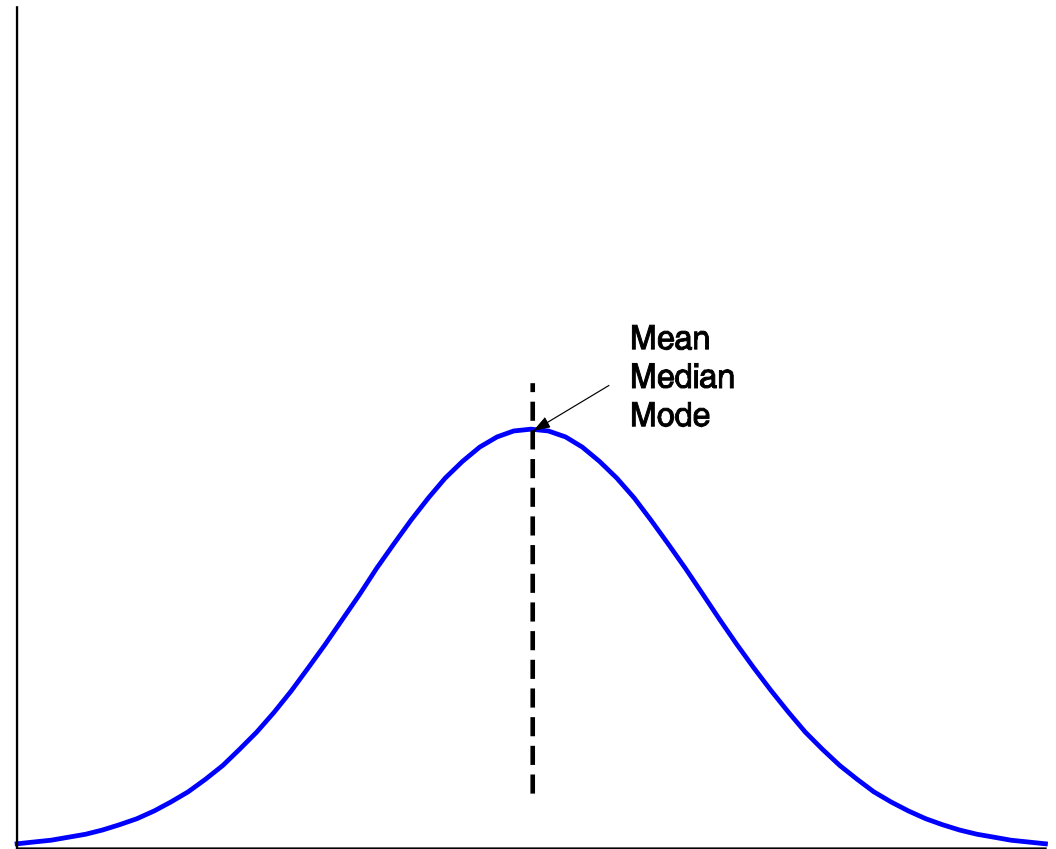
- Median:
  - We can rank the numbers from smallest to largest.
  - The median is the smallest number such that at least half of the other numbers of lesser value



# Measuring the Central Tendency

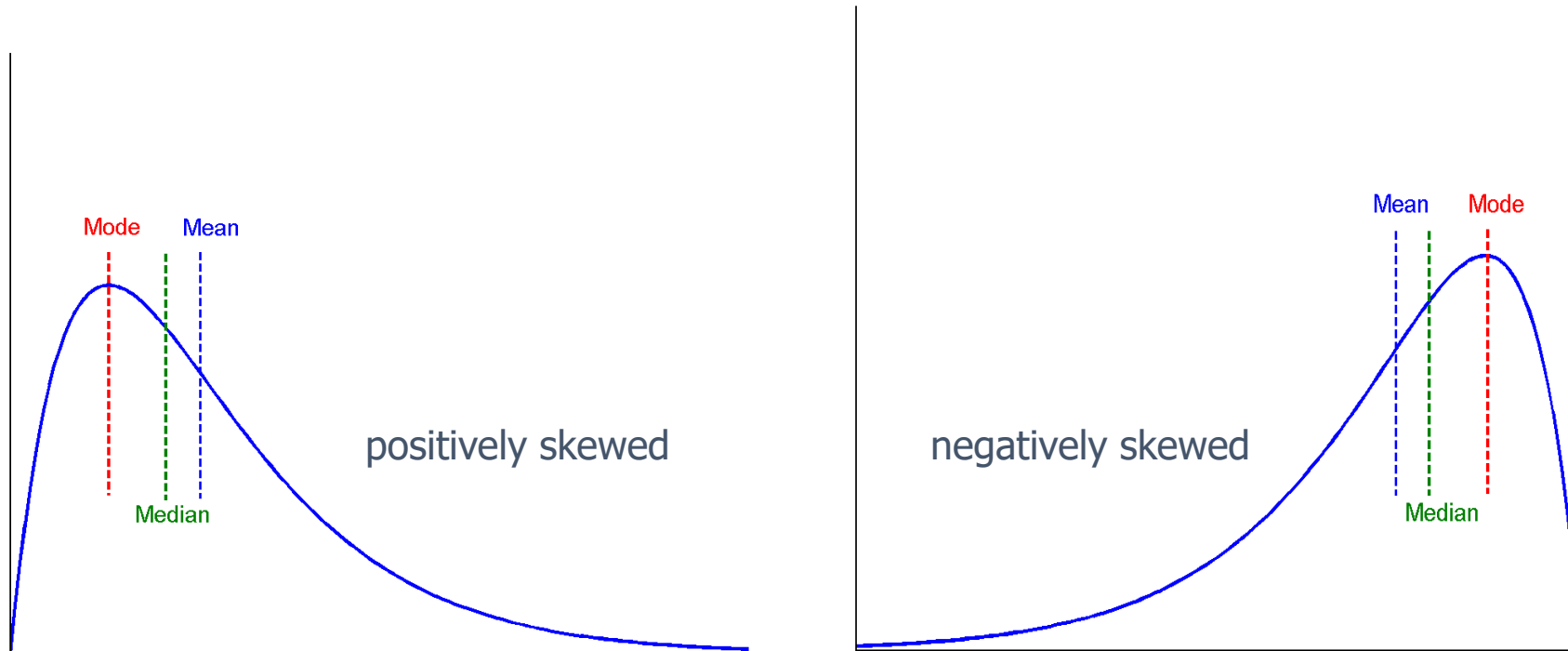
## Mode

- Value that occurs most frequently in the data
- Unimodal, bimodal, trimodal,...



# Symmetric vs. Skewed Data

- Median, mean and mode of symmetric, positively and negatively skewed data



# Dispersion of Data

- Quantile
  - Quantiles generalize median. Instead of splitting the data in half, we can split into any fractions.
    - **Quartiles:**  $Q_1$  (25<sup>th</sup> percentile),  $Q_3$  (75<sup>th</sup> percentile)
  - **Five number summary:** minimum,  $Q_1$ , median,  $Q_3$ , maximum
  - **Inter-quartile range:**  $IQR = Q_3 - Q_1$ 
    - **Outlier:** usually, a value higher/lower than  $1.5 \times IQR$

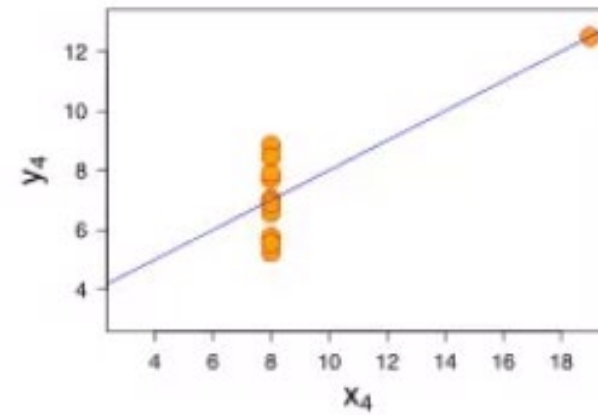
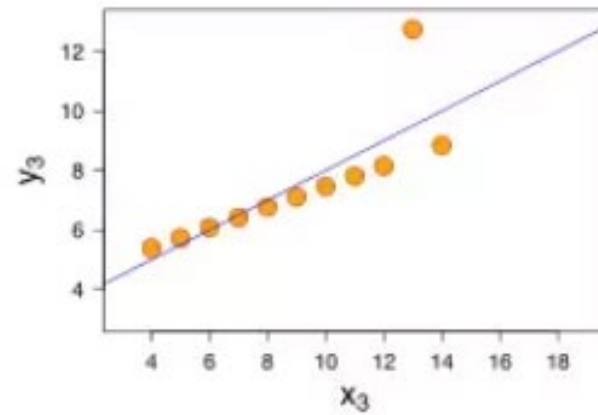
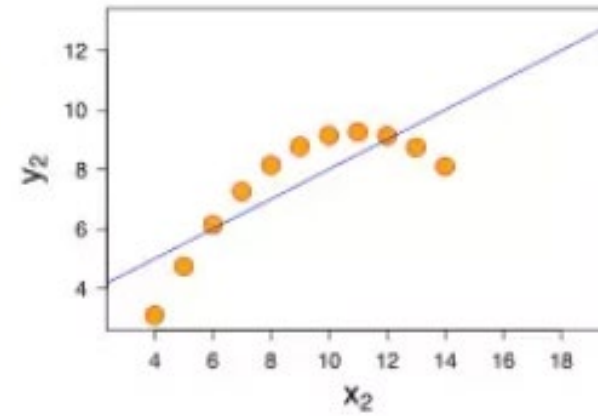
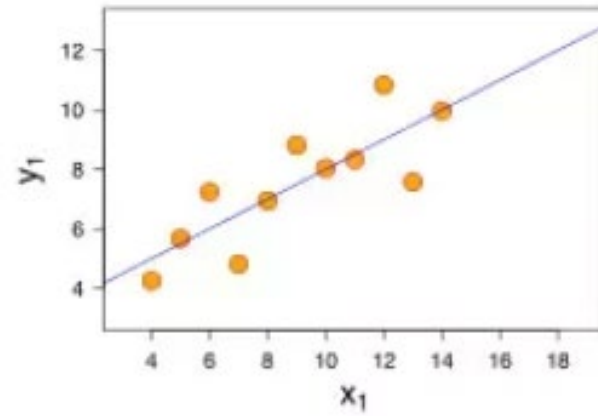
# Dispersion of Data

- **Variance** is average square distance around the mean

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left[ \sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2 \right]$$

- **Standard deviation** is the square root of variance  $s^2$

# Visualization

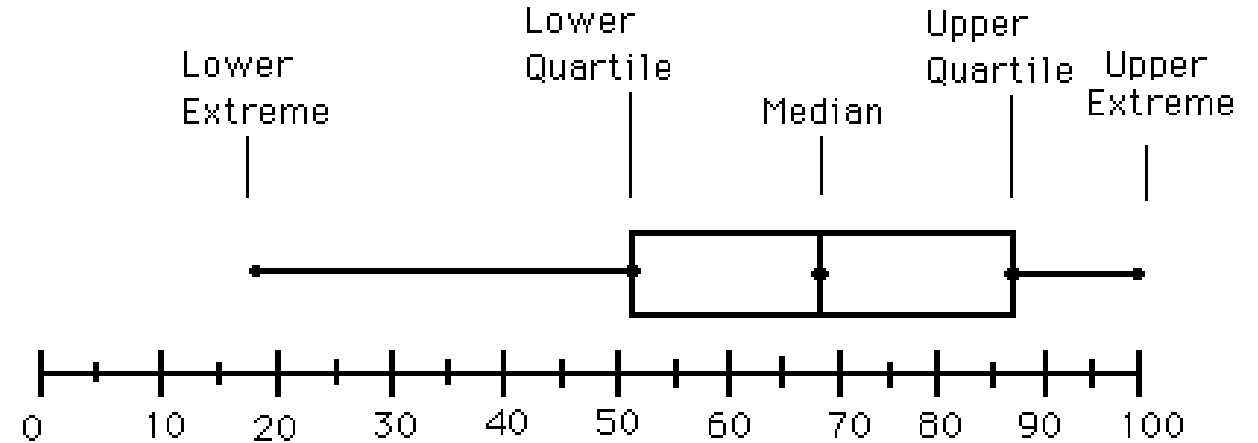


# Visualization

- Box-plot
- Histogram
- Scatter-plot
- Quantile and Quantile-Quantile Plot

# Boxplot

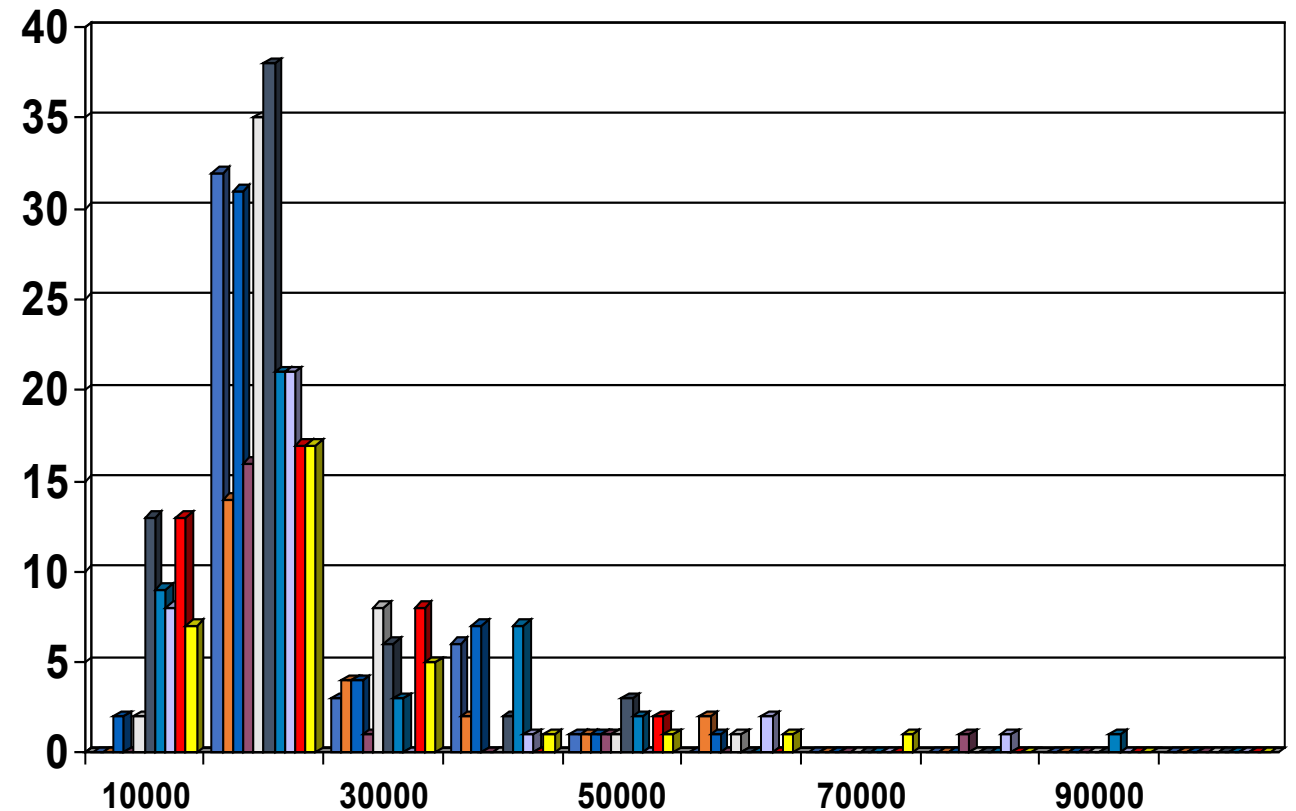
- **Five-number summary** of a distribution
  - Minimum, Q1, Median, Q3, Maximum
- **Boxplot**
  - Data is represented with a box
  - The ends of the box are at the first and third quartiles, i.e., the height of the box is IQR
  - The median is marked by a line within the box
  - Whiskers: two lines outside the box extended to Minimum and Maximum
  - Outliers: points beyond a specified outlier threshold, plotted individually



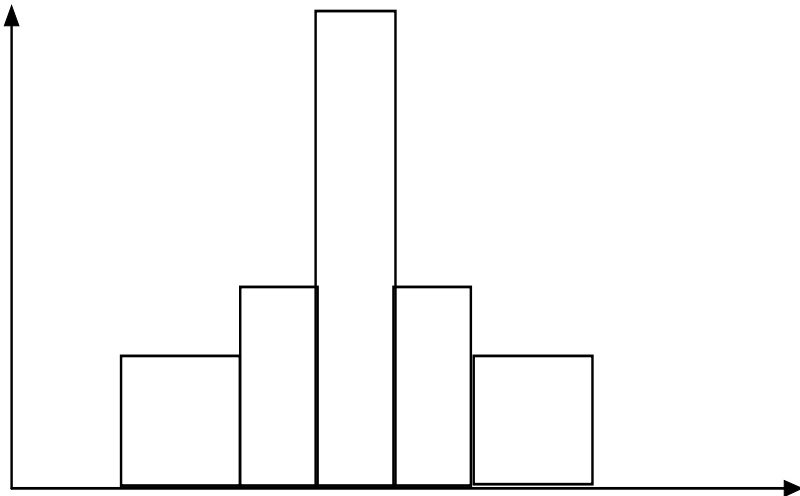
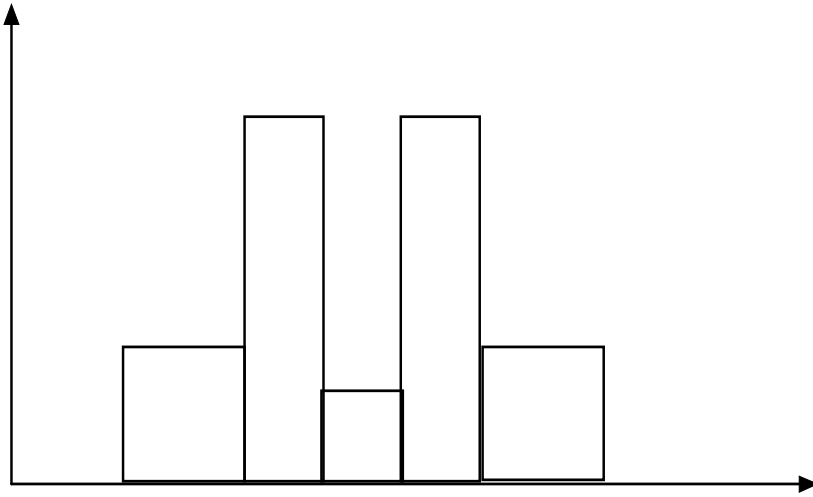


# Histogram

- Histogram: Graph display of tabulated frequencies, shown as bars
- It shows what proportion of cases fall into each of several categories
- Differs from a bar chart in that it is the *area* of the bar that denotes the value, not the height as in bar charts, a crucial distinction when the categories are not of uniform width
- The categories are usually specified as non-overlapping intervals of some variable. The categories (bars) must be adjacent



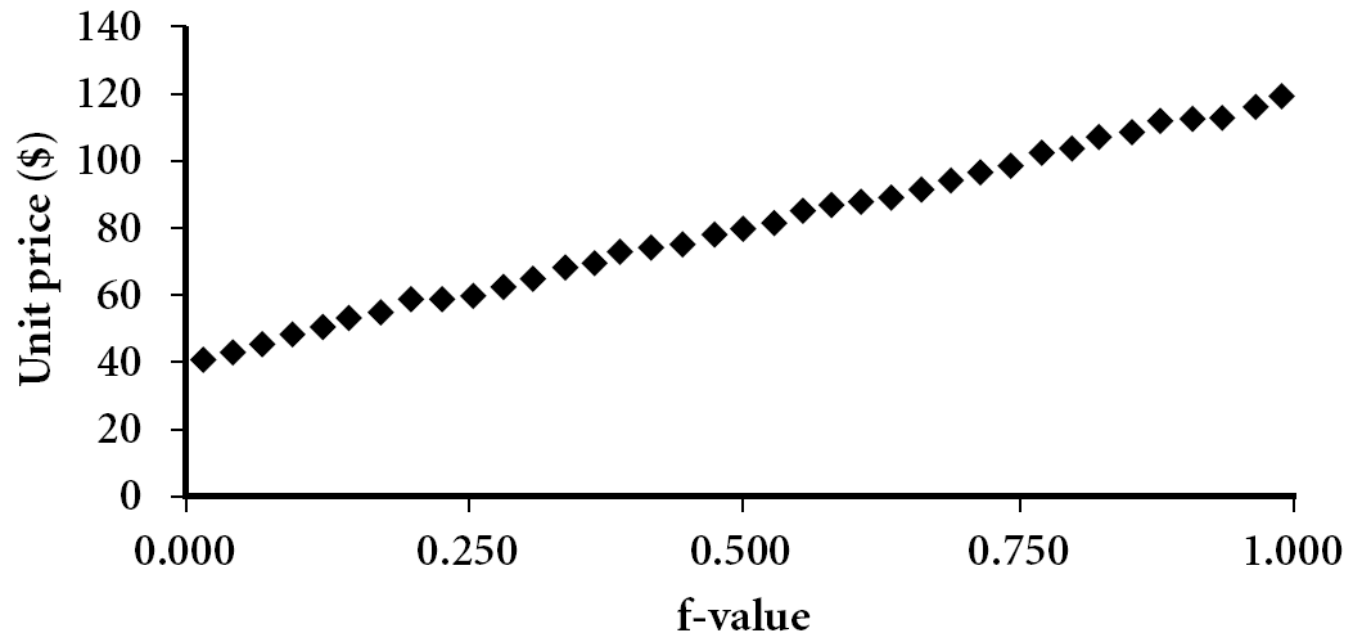
# Histograms Often Tell More than Boxplots



- The two histograms shown in the left may have the same boxplot representation
  - The same values for: min, Q1, median, Q3, max
- But they have rather different data distributions

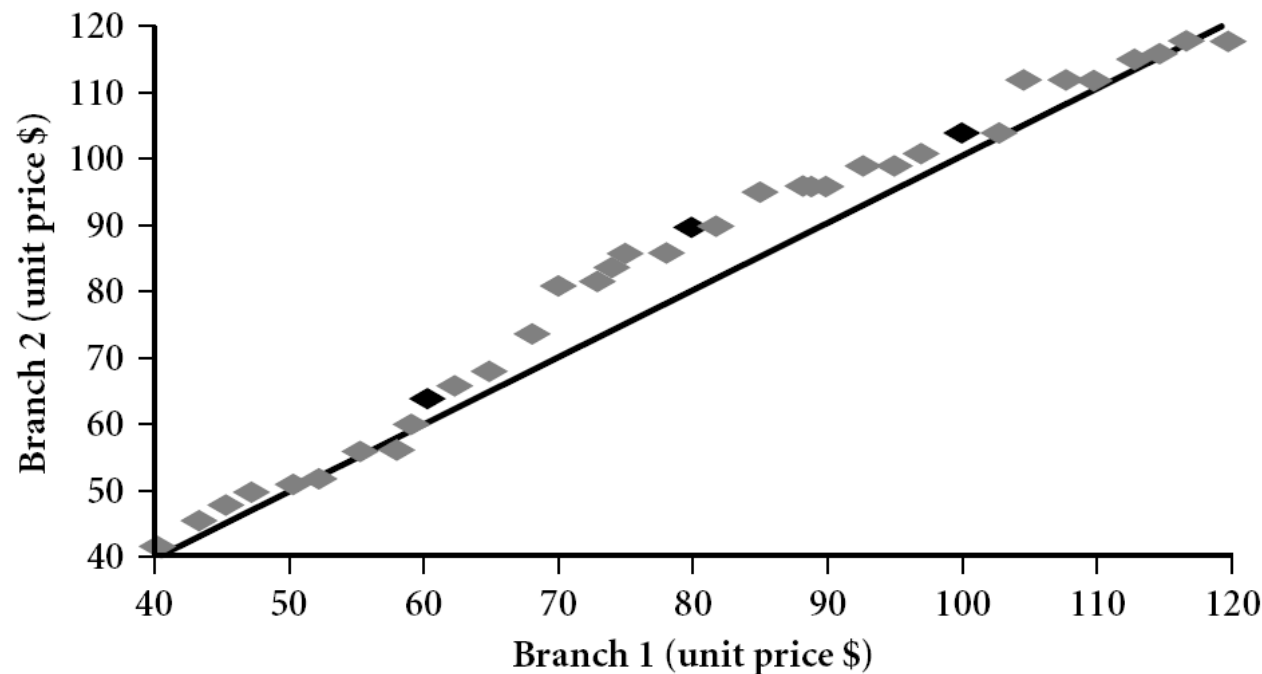
# Quantile Plot

- Plots **quantile** information
  - For a dataset consisting of points  $x_i$  sorted in increasing order,  $f_i$  indicates that approximately  $100 f_i\%$  of the data are below or equal to the value  $x_i$



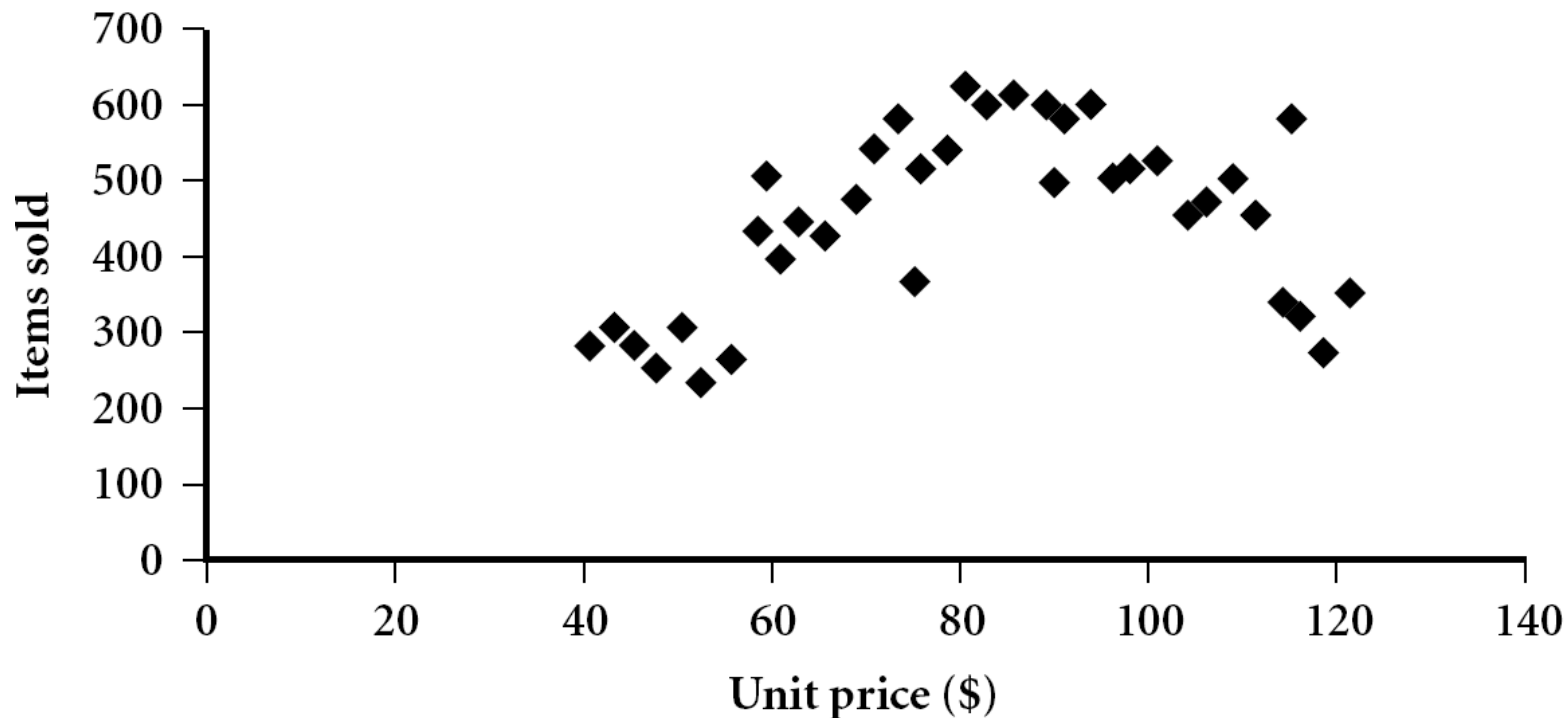
# Quantile-Quantile (Q-Q) Plot

- Graphs the quantiles of one univariate distribution against the corresponding quantiles of another
  - Example shows unit price of items sold at Branch 1 vs. Branch 2 for each quantile. Unit prices of items sold at Branch 1 tend to be lower than those at Branch 2.

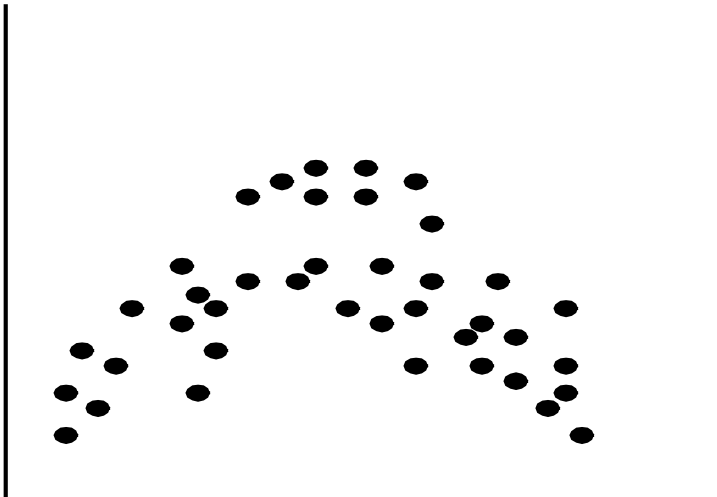
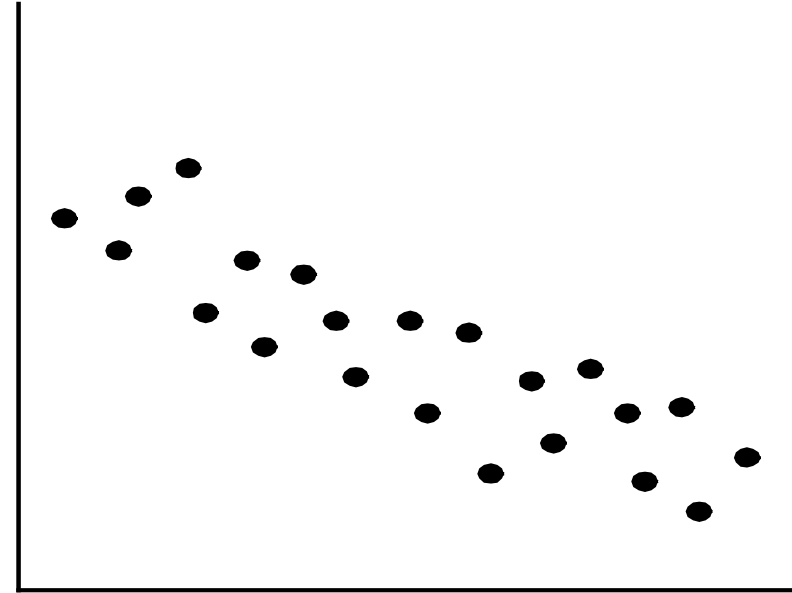
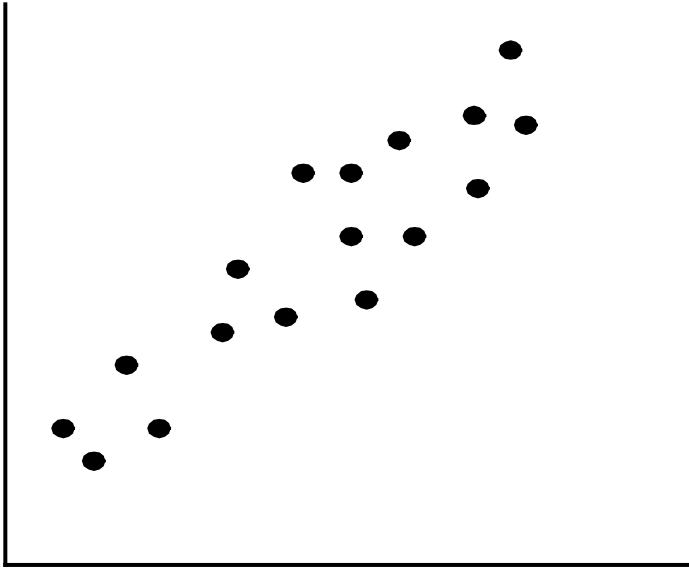


# Scatter plot

- Provides a first look at bivariate data to see clusters of points, outliers, etc
- Each pair of values is treated as a pair of coordinates and plotted as points in the plane



# Positively and Negatively Correlated Data



- The left half fragment is positively correlated
- The right half is negative correlated

# Visualizations

- **Boxplot:** graphic display of five-number summary
- **Histogram:** x-axis are values, y-axis repres. frequencies
- **Quantile plot:** each value  $x_i$  is paired with  $f_i$  indicating that approximately  $100 f_i \%$  of data are  $\leq x_i$
- **Quantile-quantile (q-q) plot:** graphs the quantiles of one univariant distribution against the corresponding quantiles of another
- **Scatter plot:** each pair of values is a pair of coordinates and plotted as points in the plane

# Summary

- Data engineering involves knowledge discovery and data mining
- Business knowledge is important for adoption of technologies supporting data mining.
- Properties of qualitative and quantitative data
- Numerical summaries of data such as mean, median, mode
- Visualization of data including histogram and scatter-plot