# 2020-2021学年秋季学期

## 数据科学导论
### *The Introduction of Data Science*

授课团队：沙瀛　周川

助　　教：梁棋

数据科学导论

*The Introduction of Data Science*
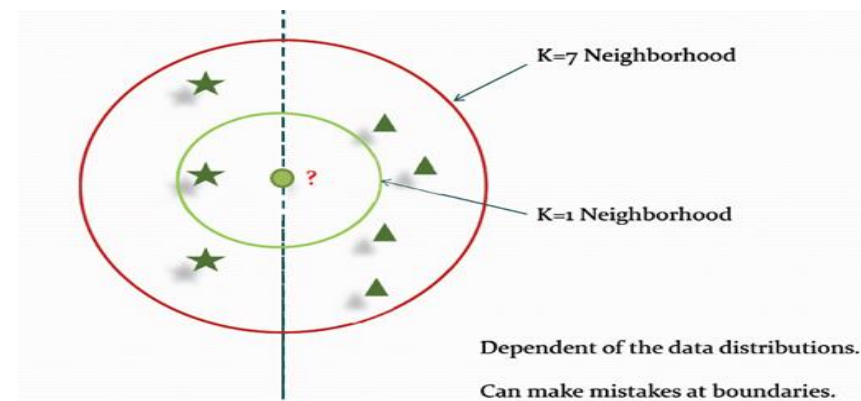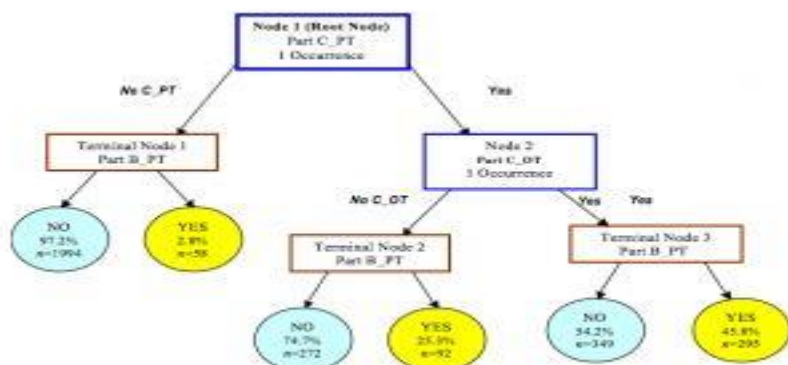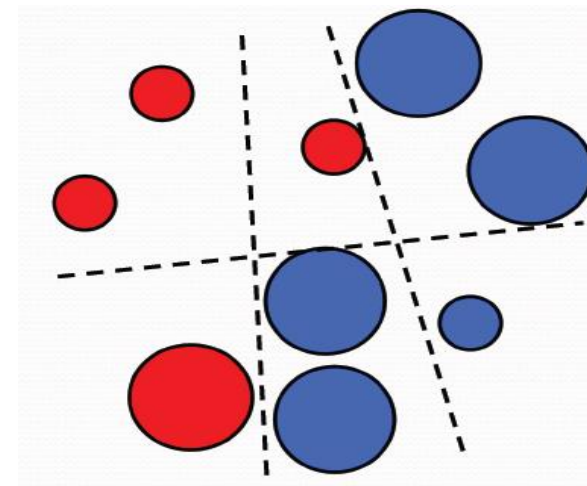
# [第 章] 总结与展望

授课教师：周川

授课时间：2020年12月25日

# 主要章节

1. 引言
2. 数据科学生命周期
3. 相关与因果-批判性思维
4. 数据预处理
5. 数据分析与建模
6. 数据可视化
7. 项目实施与沟通

8. 分类
9. 聚类
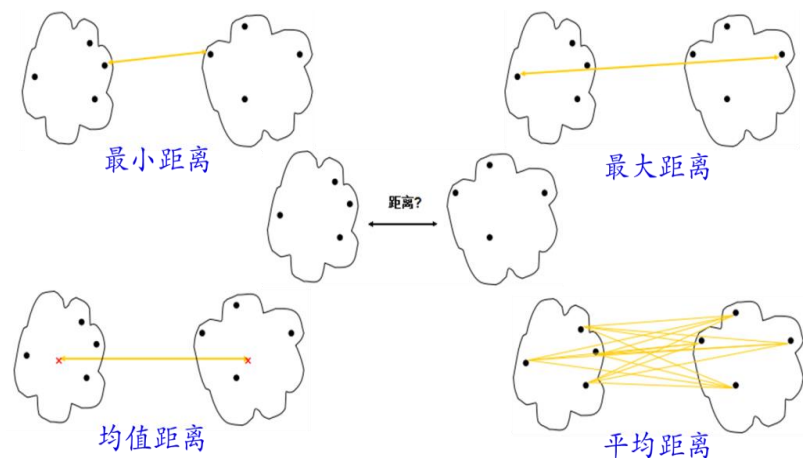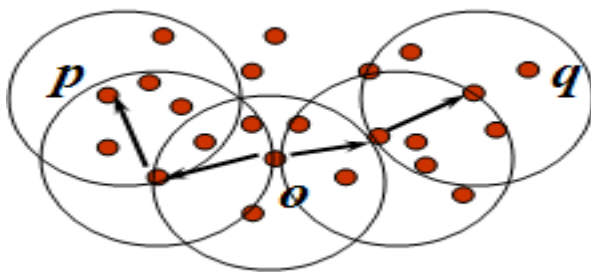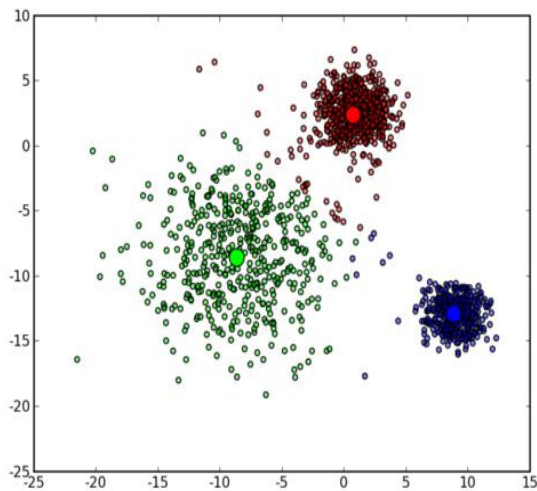10. 回归
11. 关联规则分析
12. 异常检测
13. 数据降维
14. 时间序列分析

# 第8章 分类

- 决策树算法
- 朴素贝叶斯分类器
- 最近邻分类器
- Logistics回归
- 提升方法（集成学习）

# 第9章 聚类

- 接近性度量/差异性度量
- 顺序聚类算法
- 划分聚类算法（K-means聚类算法）
- 层次聚类算法（ AGNES 、 DIANA ）
- 密度聚类算法（ DBSCAN ）

最小距离 最大距离 距离?

均值距离 平均距离

# 第10章 回归

- 一元线性回归

- 多元线性回归

- 非线性回归

- 附录：逐步回归、岭回归、Lasso回归

$$\hat{\beta} = (X^{\tau}X)^{-1}X^{\tau}Y$$

$$\begin{cases} y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \varepsilon \\ \varepsilon \sim N(0, \sigma^2) \end{cases}$$

# 第11章 关联规则分析

- 频繁相集/关联规则
- Apriori算法
- PCY算法
- FP-Growth算法
- 序列模式挖掘



图 8.5 PCY 算法中哈希表的构造和应用

# 第12章 异常检测

- 定义与主要挑战

- 基于图形的方法
  - 箱型图

- 基于统计的方法

- 分类和聚类

- 基于距离和基于密度的方法
  - 基于K-means的方法
  - 局部异常因子（LOF）方法

# 第13章 数据降维

- 主成分分析
  - 几何意义、总体主成分、样本主成分

- 因子分析
  - 因子载荷、因子旋转、因子得分

- SVD分解

- 低维嵌入

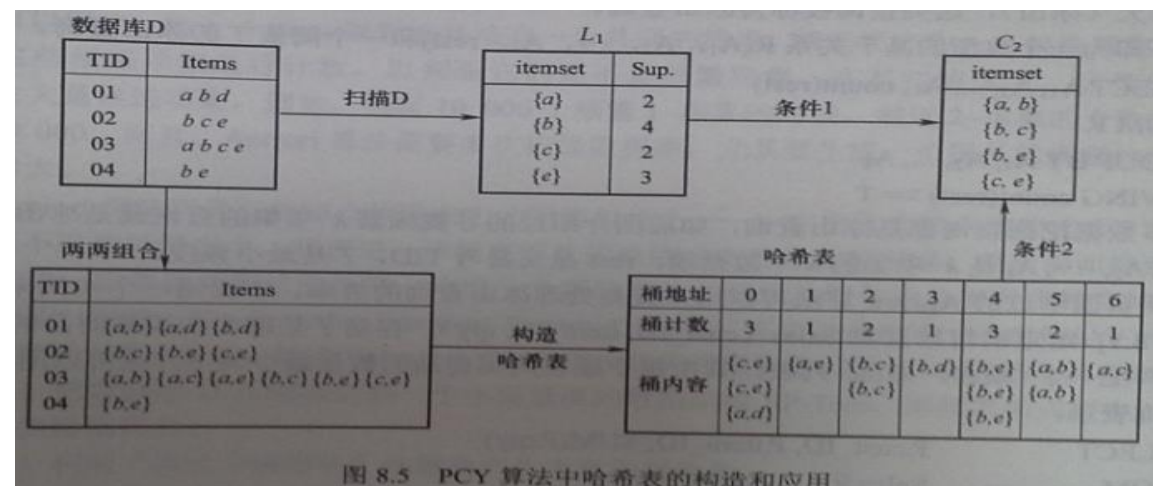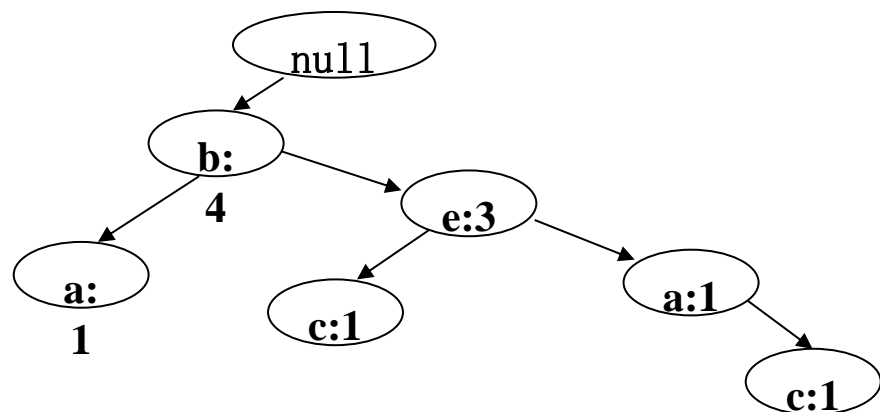$$\Sigma = \mathbf{A}\mathbf{A}' + \mathbf{D} = \mathbf{U} \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_p \end{bmatrix} \mathbf{U}'$$

$$\boldsymbol{A}_{m \times n} = \boldsymbol{U}_{m \times m} \boldsymbol{\Sigma}_{m \times n} \boldsymbol{V}_{n \times n}^{\mathbf{T}}$$

$$\varepsilon(\boldsymbol{W}) = \sum_{i=1}^{n} \left| \boldsymbol{x}_i - \sum_{j=1}^{k} w_{ij} \boldsymbol{x}_{i\_j} \right|^2 \qquad \Phi(\boldsymbol{W}) = \sum_{i=1}^{n} \left| \boldsymbol{y}_i - \sum_{j=1}^{k} w_{ij} \boldsymbol{y}_{i\_j} \right|^2$$

# 第14章 时间序列分析

- 时间序列分析的预处理
  - 平稳性检验
  - 随机性检验

- 平稳时间序列分析
  - ARMA模型
  - 平稳序列建模与学习
  - 序列预测

- 非平稳时间序列分析
  - 序列分解
  - ARIMA模型
  - 残差自回归模型
  - 指数平滑预测模型

# What is Data Science?

Extraction of **knowledge from large volumes of data** that are structured or unstructured.

It is a continuation of the fields **data mining** and **predictive analytics**

# My Definition for Data Science

The application of **data centric**, **computational**, and **inferential thinking** to

*understand the world* **&** *solve problems*

_____

**Science**

_____

**Engineering**

➢ *Data science is fundamentally <u>interdisciplinary</u>*

# Data Science

# Data Science Tools

# DATA SCIENCE LIFECYCLE: AN ALTERNATE VIEW



Figure 1.1 The lifecycle of a data science project: loops within loops

# DATA SCIENCE LIFECYCLE

# 采集的目标：快、准、全

```
Twitter          机器采集行为的限         ┐
                    制                 ├──→  用户行为模拟  ──→  基于Ajax模拟
Facebook         数据源自身的限制                               基于浏览器测试组件
                  （登陆、好友）      ┘

Google+          找到海量数据            ──→  智能采集策略  ──→  关键用户-关键词-普通用户
                 目标消息                                    联动的采集策略
                 与用户

QQ               单机采集能力的限         ──→  多机          ──→  分布式采集架构
                    制                     （分布式采集）
```

目标数据源          难点          关键技术突破          核心技术

# Data Cleanup

Real data is messay, often needs to cleaned up before useful.

- Bad forms

- Missing Data

- Useless Variables

- Wrong Data

# Data Cleanup (Contd.)

- Transform variables (date formats, String to int)
- Create derived variables
  - Derive county from IP
  - Age from ID card number
- Normalize strings
  - Different spelling and nicknames (William -> Bill)
- Feature value rescaling
- Enrich

# Data Exploration

Understand, and get a feel for what is expected (models => densities, constraints) and unexpected/ residuals (errors, outliers)

- think what this is data about? domain, background, how it is collected, what each fields mean and range of values.

- head, tail, count, all descriptives (Mean, Max, median, percentiles .. ) - Five number Summary. Min. 1st Qu. Median Mean 3rd Qu. Max.

- run a bunch of count/group-by statements to gauge if it's corrupt

# Data Exploration (Contd.)

- Plot - take random sample and explore ( scatter plot)
  - e.g. Draw scatter plot or Trellis Plot
- Find Dependencies between fields
  - Calculate Correlation
  - Dimensionality reduction
  - Cluster and look visualize clusters
- Look at frequency distribution of each field and try to find a known distribution if possible.

# Data Exploration (Contd.)

# Feature Engineering



Features

- Feature engineering is the art of finding feature that leads simplest decision algorithm. ( Good features allow a simple model to beat a complex model.)
- Best features may be a subset, or a combination, or transformed version of the features.

# How to do Feature Engineering?

- Manually pick by domain experts and trial and error.
- Search the possible combinations by training and combining subsets (e.g. Random Forest)
- Use statistical concepts like correlation and information criteria
- Reduce the features to a low dimension space using techniques like PCA.
- Automatic Feature Learning though Deep Learning

# Analysis

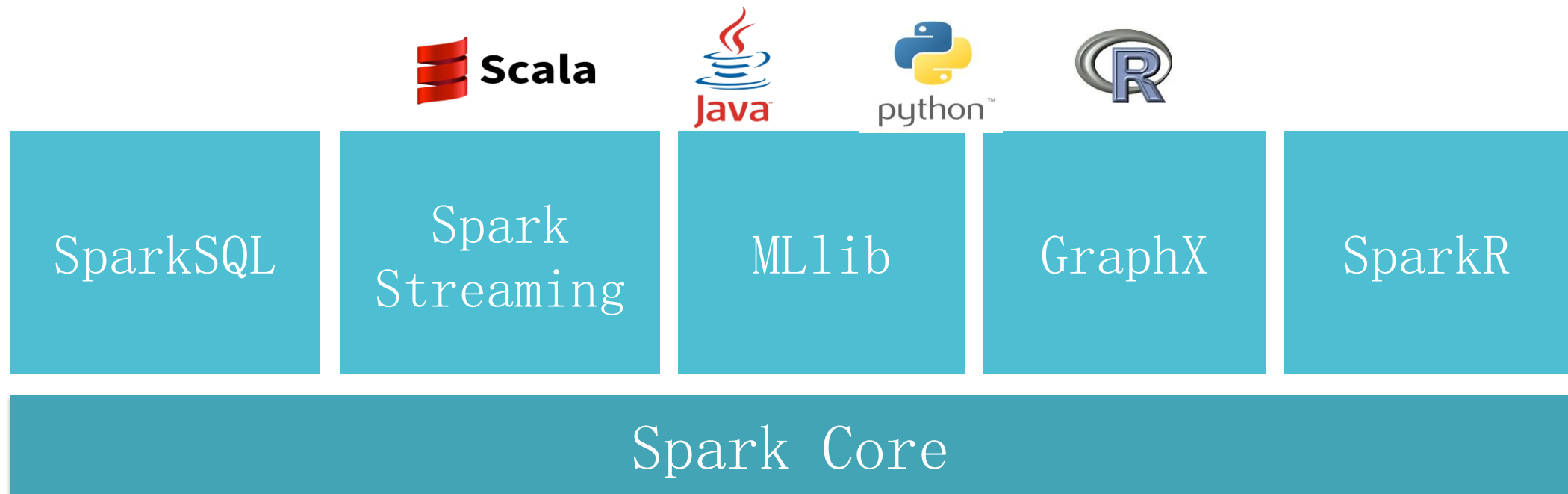- Goal of analysis is to extract knowledge

- This knowledge usually come in one of the two forms

    - KPI (Key Performance Indicators)

        - Describe key measurement for what is being measured. (e.g. revenue per year, profit margin, revenue for sqft in retail, revenue per employer)

    - Models to describe or predict the data

        - e.g. Machine Learning models or Statistical models

# 4 Analysis types by time to decision
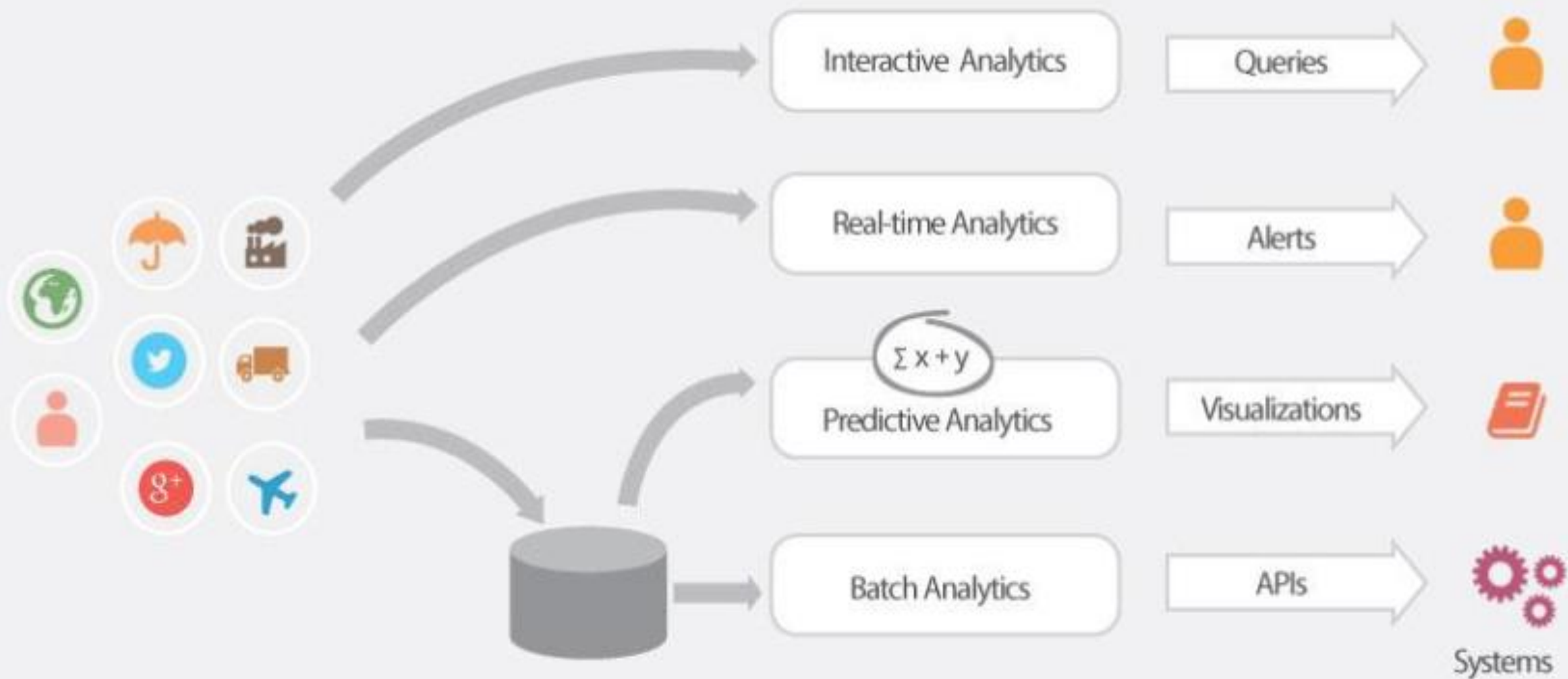
- Hindsight ( what happened?)
  - Done using Batch Analytics like MapReduce
- Oversight ( what is happening?)
  - Done using Realtime Analytics technologies like CEP
- Insight ( why things happening?)
  - Done with Data Mining and Unsupervised learning algorithms like Clustering
- Foresight ( what will happen?)
  - Done by building models using Machine learning or one of other techniques

# Data Analytics Tools : Apache Spark

➢ Unifies **batch, interactive, streaming** workloads

➢ Easy to build sophisticated applications
  ➢ Support iterative, graph-parallel algorithms
  ➢ Powerful APIs in Scala, Python, Java, R

Interactive Analytics → Queries

Real-time Analytics → Alerts

$\Sigma x + y$

Predictive Analytics → Visualizations

Batch Analytics → APIs

Systems

Collect Data → Analyze & Make Decisions → Communicate

# Inconvenient Truth About Data Science

- Data is **never clean**

- You will spend **most of your time** cleaning and preparing data.

- **95%** of tasks do not require deep learning

- In **90%** of cases generalized linear regression will do the trick

- Big Data is just a **tool**.

- You should embrace the **Bayesian** approach.

- No one cares **how you did it**.

- Academia and business are **two different** worlds.

- **Presentation is key** – be a master of Power Point.

- All models are false, but **some** are **useful**.

- There is no fully automated Data Science. You need to **get your hands dirty**.

# Data Science is extracting knowledge by analyzing data

**Hindsight**

**Oversight**

**Insight**

**Foresight**

# 谢谢大家！