

# 数据挖掘与应用

## 数据挖掘简介

授课教师：周晟

浙江大学 软件学院

2023.09.19



# 课程内容

## 1 课程简介

- 课程信息
- 准备知识

## 2 数据挖掘简介

- 从数据到知识
- 数据为中心的人工智能
- 数据挖掘的定义
- 数据挖掘与学术
- 课程大纲

## 3 数据类型与度量

- 数据类型
- 数据的相关性度量

## 4 数据的统计与可视化

- 数据的统计描述
- 数据可视化

## 5 数据预处理与数据降维

- 数据质量
- 主成分分析 PCA



# 课程内容

## 1 课程简介

- 课程信息
- 准备知识

## 2 数据挖掘简介

- 从数据到知识
- 数据为中心的人工智能
- 数据挖掘的定义
- 数据挖掘与学术
- 课程大纲

## 3 数据类型与度量

- 数据类型
- 数据的相关性度量

## 4 数据的统计与可视化

- 数据的统计描述
- 数据可视化

## 5 数据预处理与数据降维

- 数据质量
- 主成分分析 PCA



# 课程简介

## 课程信息

本课程主要讲述数据挖掘相关的基本概念，经典任务，前沿技术以及在实际生产生活中的应用。学习本课程有望在掌握数据挖掘相关知识的同时也培养相关的实践动手能力。

- ① 课程主页: <https://zhoushengisnoob.github.io/courses/index.html?course=data-mining-2023>
- ② 授课时间: 秋学期周二上午 1-4 节
- ③ 授课教师: 周晟
- ④ 课程助教: 郑卓男, 徐鸿嘉
- ⑤ 考核方式: 随堂测试 (2\*10%) + 期末报告 (80%)



# 准备知识

## 数学基础

- ① 线性代数基础（常见的矩阵理论）
- ② 概率论基础（基础贝叶斯理论，了解常见的概率分布）
- ③ 统计学基础（常见的统计方法与指标）

## 代码基础

- ① Python
- ② Numpy
- ③ Pandas
- ④ Scikit-Learn
- ⑤ Pytorch



# 参考资料

- ① 书籍：《数据挖掘》《统计学习方法》



- ② Survey: CSUR, TPAMI, TKDE  
③ Paper  
④ Online Resources



# 课程内容

## 1 课程简介

- 课程信息
- 准备知识

## 2 数据挖掘简介

- 从数据到知识
- 数据为中心的人工智能
- 数据挖掘的定义
- 数据挖掘与学术
- 课程大纲

## 3 数据类型与度量

- 数据类型
- 数据的相关性度量

## 4 数据的统计与可视化

- 数据的统计描述
- 数据可视化

## 5 数据预处理与数据降维

- 数据质量
- 主成分分析 PCA



# 数据的来源

人类从诞生之日起就不断地产生数据，例如文字，图像以及音乐等



# 数据的来源

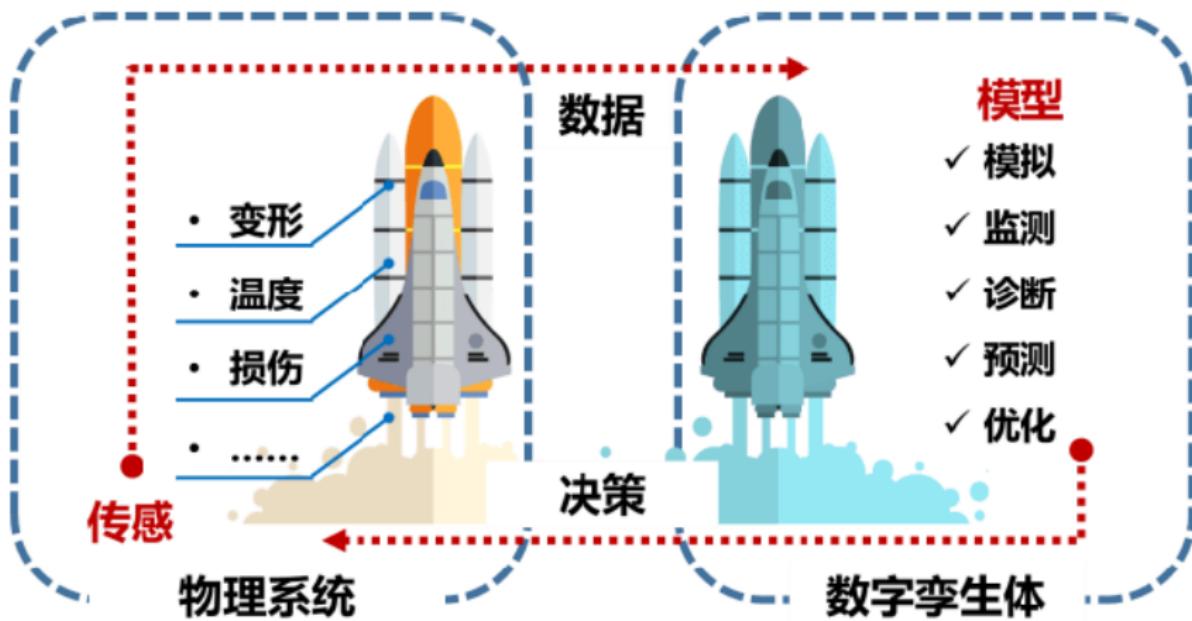
人类从诞生之日起就不断地产生数据，例如文字，图像以及音乐等



随着数据**采集、存储**技术的快速发展，人类采集和存储了海量的数据，包括商业，社会，科学，工程等领域



# 数据的来源



数字孪生等技术使得万物皆可变成数据

# 数据来源



现在，每年会产生 16.3 泽 ( $1,000,000,000,000,000,000,000,000$  ( $10^{21}$ )) 字节的数据

# 常见的数据种类

目前常用的数据类型主要包括：

- ① 数据库数据
- ② 图像数据
- ③ 音频数据
- ④ 视频数据
- ⑤ 网络数据
- ⑥ ...

# 常见的数据种类

目前常用的数据类型主要包括：

- ① 数据库数据
- ② 图像数据
- ③ 音频数据
- ④ 视频数据
- ⑤ 网络数据
- ⑥ ...



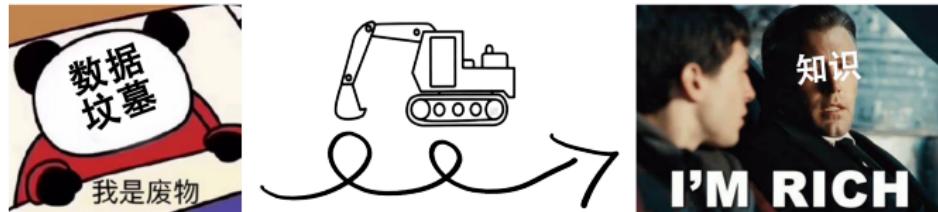
数据越多越好吗？

# 从数据到知识

## 直接使用数据面临的困难

- ① 数据量大 (大海捞针)
- ② 噪声和异常多
- ③ 数据特征高维 (盲人摸象)
- ④ 模式特征不显著 (雾里看花)

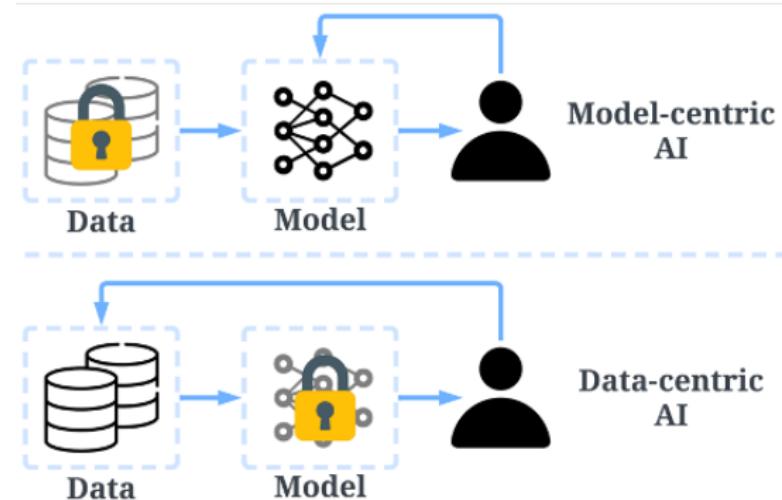
收集在大型数据库中的数据变成了“数据坟墓”，有必要系统地开发数据挖掘工具，将数据坟墓转换成知识“金块”。



We are drowning in data, but starving for knowledge!

# 数据为中心的人工智能

Data-centric AI (DCAI) is the discipline of systematically engineering the data used to build an AI system. –Andrew Ng



## 数据为中心的人工智能<sup>1</sup>

<sup>1</sup>Zha, Daochen, et al. Data-centric AI: Perspectives and Challenges. SDM, 2023. ↗ ↘ ↙ ↘

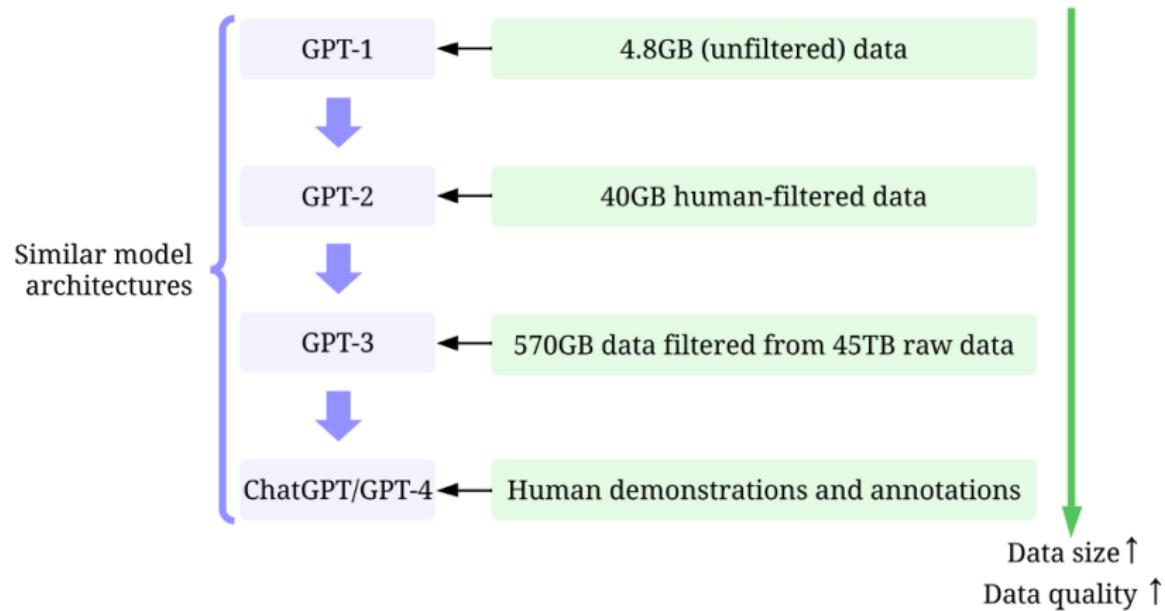
# 人工智能的发展离不开高质量的数据

时间	人工智能突破	数据
1994	Human-level spontaneous speech recognition	Spoken Wall Street Journal articles and other texts (1991)
1997	IBM Deep Blue defeated Garry Kasparov	700,000 Grandmaster chess games (1991)
2012	AlexNet, one of the first successful CNNs	ImageNet corpus of 1.5 million labeled images (2010)
2021	AlphaFold, AI for science	Annotated protein sequence (2017)
Now	Large language models	Large text data

## 数据驱动下的人工智能突破<sup>2</sup>

<sup>2</sup><http://www.spacemachine.net/views/2016/3/datasets-over-algorithms>

# 大模型与大数据



## GPT 模型用到的数据<sup>3</sup>

<sup>3</sup>Zha, Daochen, et al. Data-centric Artificial Intelligence: A Survey. *arXiv*, 2023.

# 大模型时代是否需要数据挖掘

一般认为，人工智能技术发挥作用的三要素



大数据，大模型，大算力三者缺一不可<sup>4</sup>

大模型可以开源，大算力可以购买，（有价值的）大数据从哪来？

<sup>4</sup>Zhihua Zhou, Abductive Learning, CCF-GAIR 2020

# 数据挖掘的定义

Data Mining(knowledge discovery from data)

- Data Mining is the process of automatically extracting **interesting and useful hidden** patterns from usually **massive**, incomplete and noisy data.
- 从**大量**数据中提取**有价值的模式或知识**的过程

## 别名

- 数据库中的知识发现 (Knowledge Discovery in Databases, **KDD**)
- 知识提取 (Knowledge Extraction)
- 数据/模式分析 (Data/Pattern Analysis)

不是所有数据分析都是数据挖掘

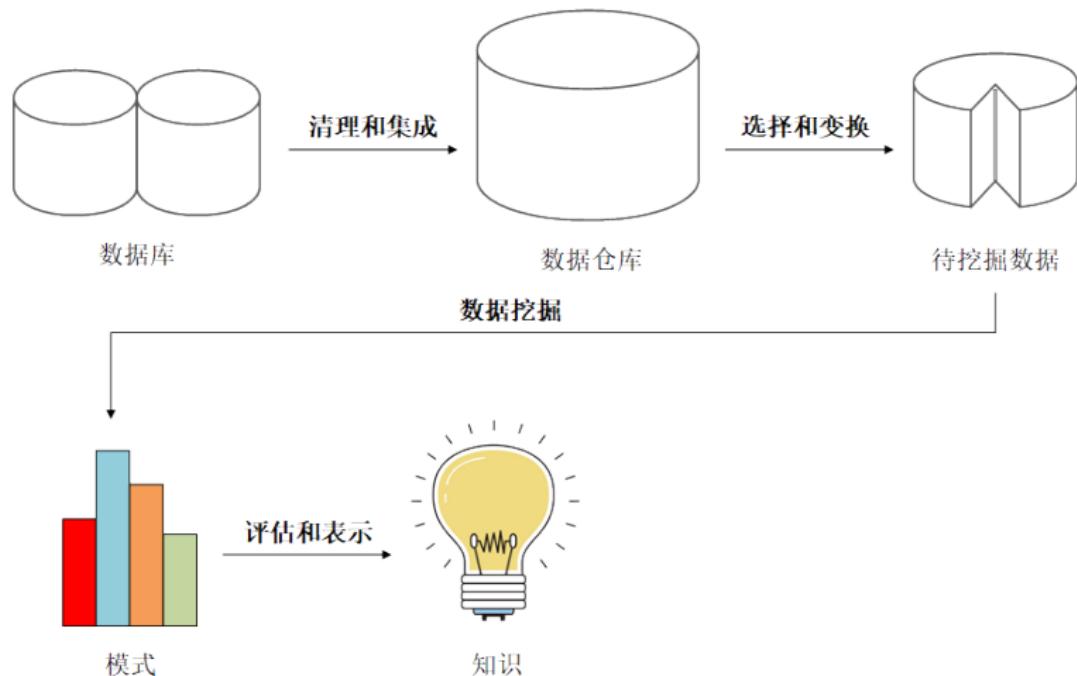
- 简单的检索和查询处理
- 演绎专家系统

# 数据挖掘的步骤

数据挖掘是将**数据**转化为**知识**的过程，其主要包含如下的步骤：

- ① 数据清理（消除噪声和删除不一致数据）
- ② 数据集成（多种数据源/模态可以组合在一起）
- ③ 数据选择（从数据库中提取与分析任务相关的数据）
- ④ 数据变换（通过汇总或聚集操作，把数据变换和统一成适合挖掘的形式）
- ⑤ **知识挖掘**（基本步骤，使用智能方法提取数据模式）
- ⑥ **模式评估**（根据兴趣度度量，识别代表知识的真正有趣的模式）
- ⑦ 知识表示（使用可视化和知识表示技术，识别代表知识的真正有趣的模式）

# 数据挖掘的步骤



## 数据挖掘的主要步骤

# 知识——有一定价值的模式

数据挖掘可以产生海量模式或规则，但并不是所有模式都是有用的。

有价值的模式：

- ① 易于被人理解
- ② 在某种确信度上，对于新的或检验数据是有效的
- ③ 是潜在有用的
- ④ 是新颖的

模式价值度量

- ① 客观度量：基于所发现模式的结构和相关的统计量
- ② 主观度量：基于用户对数据的观念，反映用户需要和兴趣

# 从数据中我们可以挖掘什么知识

描述性 (descriptive) 知识 【本门课程重点】:

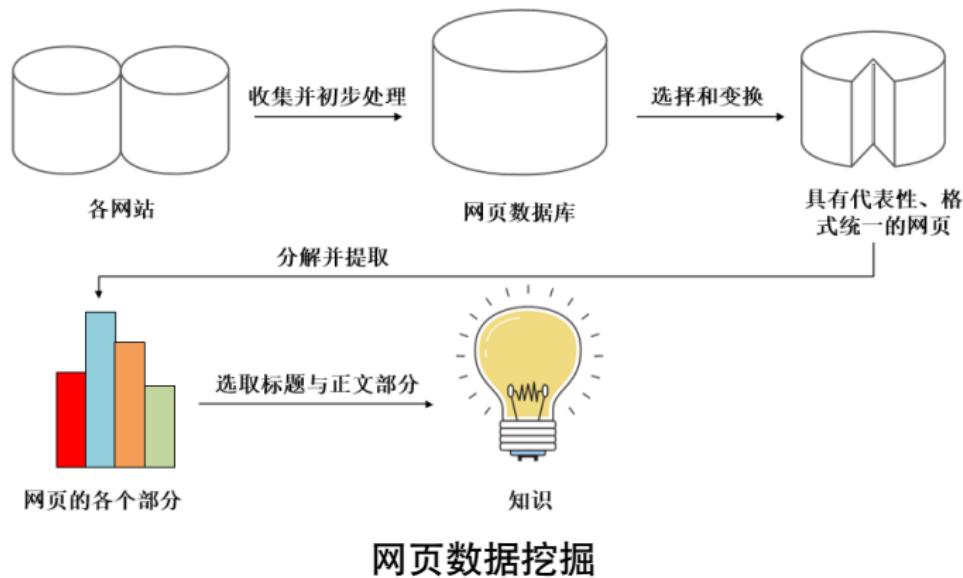
- ① 分类 (离散标号) 与回归 (连续数值)
- ② 聚类分析 (数据内在关系, 不考虑标号。)
- ③ 异常检测
- ④ 关联分析, 因果推断

预测性 (predictive) 知识:

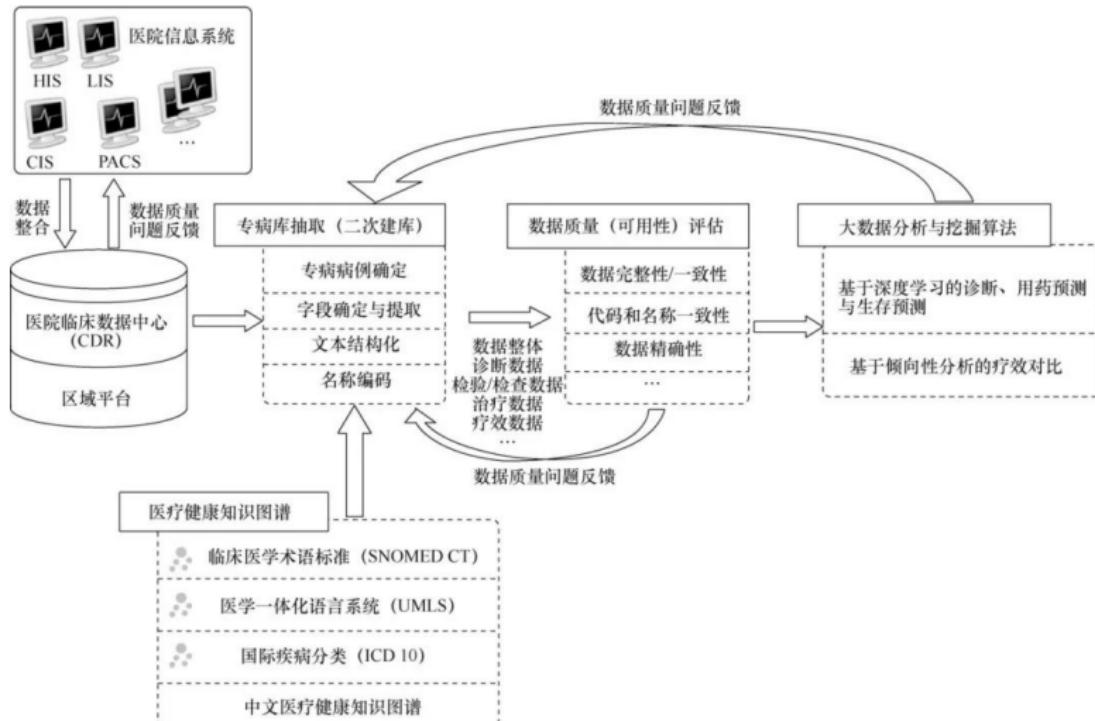
在当前数据上进行归纳, 对未来进行预测

# 数据挖掘的应用场景（一）

## 服务于视力障碍者的网页主体数据提取与展示

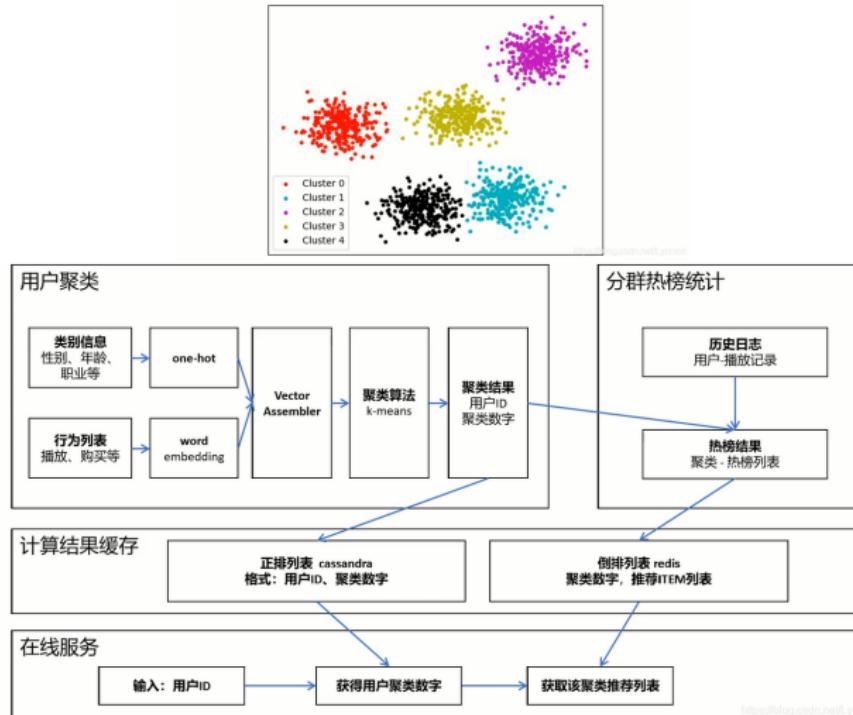


# 数据挖掘的应用场景（二）



## 基于电子病历的临床数据挖掘

# 数据挖掘的应用场景（三）



## 基于用户聚类的推荐系统

# 数据挖掘的应用场景（四）

工具中心为基础

- 存储工具
- 处理工具

规则中心为驱动

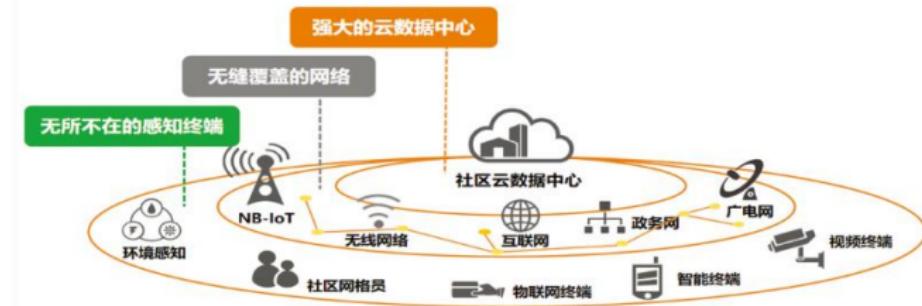
- 规则管理
- 调度管理

能力中心为核心

- 数据采集
- 数据安全
- 数据存储
- 数据质量
- 数据处理
- 数据资产
- 元数据

开放中心为索引

- 数据应用
- 数据运营



基于数据挖掘的智慧城市<sup>5</sup>

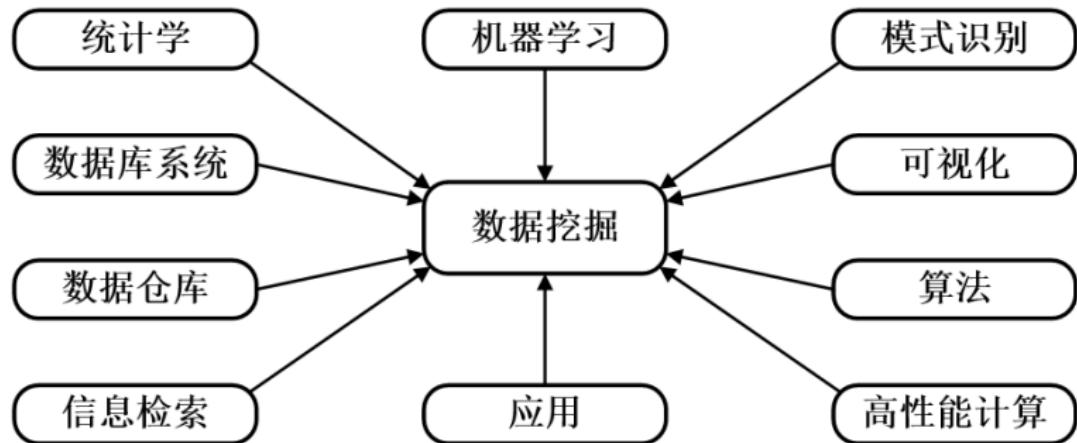
<sup>5</sup><http://www.asiainfodata.com/article/detail/30>

# 数据挖掘的主要挑战

虽然数据挖掘得到了广泛的研究，但是仍然面临如下的挑战：

- ① 新的知识类型 (需求驱动挖掘)  $\Rightarrow$  不要做无意义的数据挖掘
- ② 数据特征高维 (Curse of Dimensionality)  $\Rightarrow$  深度学习
- ③ 跨学科交叉 (AI4Science)  $\Rightarrow$  打 (Jiang) 破 (Wei) 壁 (Da) 垒 (Ji)
- ④ 实时性要求高 (Efficiency, Online)  $\Rightarrow$  复杂度分析
- ⑤ 数据不确定性强, 噪声大 (Noisy, Uncertainty)  $\Rightarrow$  落地 VS 假设
- ⑥ 数据隐私保护 (Federated Learning, Swarm Learning[1])  $\Rightarrow$  红线
- ⑦ 挖掘知识难以理解和解释 (Causal Inference)  $\Rightarrow$  可解释性研究

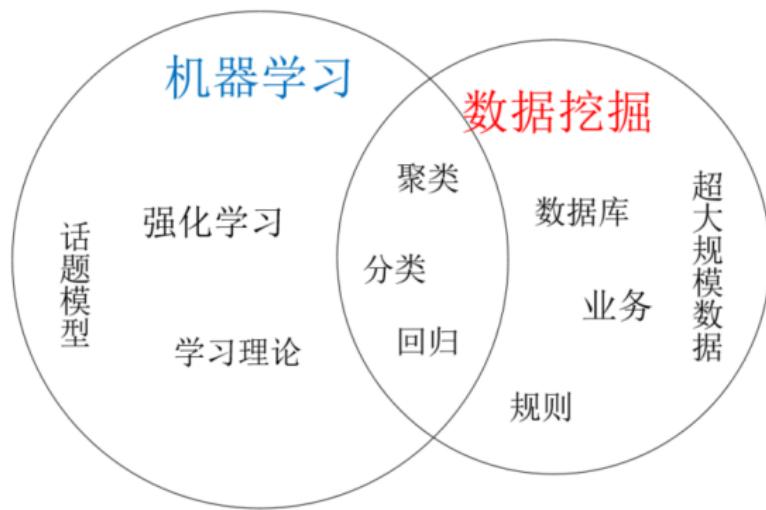
# 数据挖掘的常用技术



数据挖掘的常用技术

# 数据挖掘 VS 机器学习

- ① 机器学习研究通常关注模型的准确率（数据 -> 模型）
- ② 除准确率外，数据挖掘研究非常强调挖掘方法在大型数据集上的有效性和可伸缩性，以及处理复杂数据类型的方法，开发新的，非传统的方法。（模型 -> 知识）



# 数据挖掘的会议与期刊

## ● 数据挖掘会议

- ACM SIGKDD Int. Conf. on Knowledge Discovery in Databases and Data Mining (**KDD**)
- ACM International Conference on Information and Knowledge Management (**CIKM**)
- SIAM Data Mining Conf. (**SDM**)
- (IEEE) Int. Conf. on Data Mining (**ICDM**)
- European Conf. on Machine Learning and Principles and practices of Knowledge Discovery and Data Mining (**ECML-PKDD**)
- Pacific-Asia Conf. on Knowledge Discovery and Data Mining (**PAKDD**)
- Int. Conf. on Web Search and Data Mining (**WSDM**)

## ● 其他相关会议

- DB 会议: ACM SIGMOD, VLDB, ICDE, EDBT, ICDT, ...
- Web and IR 会议: WWW, SIGIR, WSDM
- ML 会议: ICML, NIPS      PR 会议: CVPR

## ● 期刊

- IEEE Trans. On Knowledge and Data Eng. (TKDE)
- ACM Trans. on KDD

# 课程结构

## 数据挖掘与应用课程结构

理论	应用
认识数据, 数据预处理	
支持向量机 SVM	
	决策树与 Boosting 算法
经典聚类算法	
	深度聚类算法
经典异常检测方法	
	深度异常检测方法
	数据挖掘实战

以数据挖掘主要任务为基础, 以实用为主要目标。

# Q&A

# 课程内容

## 1 课程简介

- 课程信息
- 准备知识

## 2 数据挖掘简介

- 从数据到知识
- 数据为中心的人工智能
- 数据挖掘的定义
- 数据挖掘与学术
- 课程大纲

## 3 数据类型与度量

- 数据类型
- 数据的相关性度量

## 4 数据的统计与可视化

- 数据的统计描述
- 数据可视化

## 5 数据预处理与数据降维

- 数据质量
- 主成分分析 PCA

# 数据挖掘中常用的数据类型

虽然自然界存在海量的数据，但是仅能够被计算机存储和识别的数据才能作为数据挖掘的对象。而在当前数据挖掘系统中常用的数据类型主要包括：

- ① 记录型数据（数据之间相互独立，共享特征空间）
- ② 序列化数据（数据之间通过时间维度进行排列）
- ③ 网络数据（数据之间通过关系显式链接）
- ④ 结构化数据（数据的结构固定，只在对应的特征不同）

记录型数据是基本类型，不同类型的数据之间可以相互转化

# 数据的表示

数据由样本构成, 每个样本由特征 (attribute) 进行描述, 数据通常以  $N \times K$  的矩阵/张量形式进行表示

pandas.DataFrame:  $N=4$ ,  $K=5$

	姓名	年龄	性别	平均绩点	是否选课
0	张三	22	男	3.83	True
1	李四	21	女	3.87	False
2	李雷	20	男	2.80	True
3	韩梅梅	22	女	4.00	True

# 数据的特征

用来描述一个给定对象的一组属性称作属性向量（或特征向量）

常见的属性类型包括：

- ① 标称属性 (nominal/categorical attribute)：每个值代表某种类别，编码或状态。不一定要用文本表示，也可以用数字表示，除非想要理解具体的语义。如果只是用作分类，只需要数字即可。常见的包括用户名，昵称，身份证号
- ② 离散属性 (binary/bool attribute)：一共只有两种状态
- ③ 序数属性值之间具有意义的序 (ranking) (优秀, 良好, 中等, 差)
- ④ 连续特征 numeric/quantitative
  - ① 区间标度 Interval-scaled
  - ② 比率标度 Ratio-scaled

姓名(string) 年龄(int) 性别(string) 平均绩点(float) 是否选课(bool)

# 数据的相关性度量

数据间的**相关性**是海量数据挖掘的基础

如何基于属性对数据之间的关系进行描述？

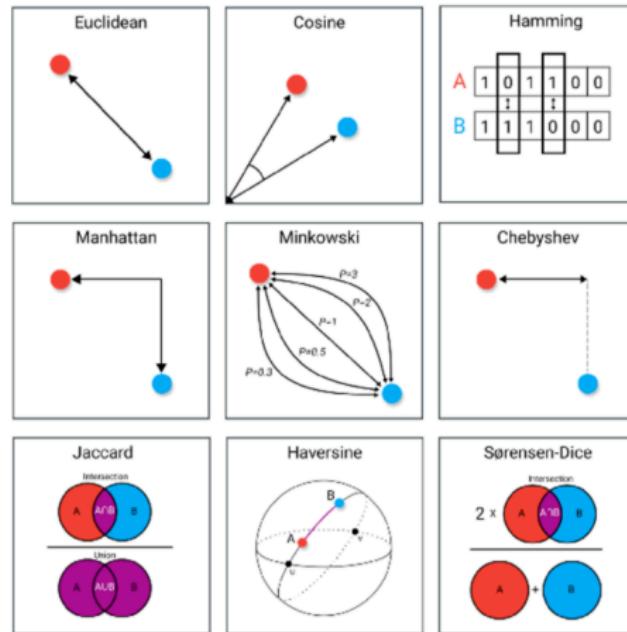
- ① 标称属性的距离度量
- ② 二元属性的距离度量 (Jaccard 以及 TP/FN 等等)
- ③ **数值属性的距离度量**
- ④ 序数属性的距离度量

# 数值属性的距离度量

在许多数据挖掘的任务中，数据（样本）之间的相似性或距离是最基础的信息单元。带有数值属性的数据之间的距离度量得到了广泛的研究。距离度量 (distance measure) 满足的基本性质

- ① 非负性:  $dist(x_i, x_j) \geq 0$
- ② 同一性:  $dist(x_i, x_j) = 0 \iff x_i = x_j$
- ③ 对称性:  $dist(x_i, x_j) = dist(x_j, x_i)$
- ④ 传递性:  $dist(x_i, x_j) \leq dist(x_i, x_k) + dist(x_k, x_j)$

# 数值属性的距离度量



## 常见的距离度量<sup>6</sup>

<sup>6</sup><https://toutiao.io/posts/8r7zk9b/preview>

# 闵可夫斯基距离

## 定义

给定样本  $x_i = (x_{i1}, x_{i2}, \dots, x_{id}) \in \mathcal{R}^d$  和  $x_j = (x_{j1}, x_{j2}, \dots, x_{jd}) \in \mathcal{R}^d$ , 闵可夫斯基距离 (Minkowski Distance) 定义为:

$$\left( \sum_{k=1}^d |x_{ik} - x_{jk}|^p \right)^{1/p}$$

上述定义也称为  $L_p$  范数 (norm)。

闵可夫斯基距离是数据挖掘/机器学习中最常用的一组距离度量。

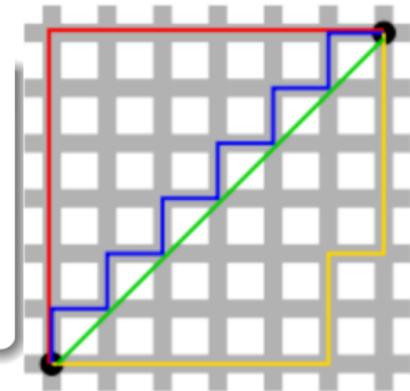
# 闵可夫斯基距离的一般形式

曼哈顿距离 Manhattan Distance

当  $p = 1$  时，闵可夫斯基距离为曼哈顿距离：

$$\text{dist} (x_i, x_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_1 = \sum_{k=1}^d |x_{ik} - x_{jk}|$$

曼哈顿距离也称为街区距离 (City Block Distance)



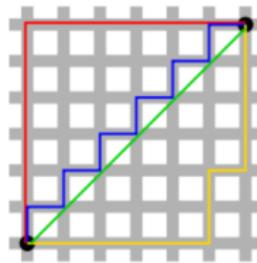
曼哈顿距离与欧式距离

# 闵可夫斯基距离的一般形式

欧氏距离 Euclidean Distance

当  $p = 2$  时，闵可夫斯基距离为欧氏距离：

$$dist(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \cdots + (x_n - y_n)^2}$$



曼哈顿距离与欧式距离

随着**数据维度**的增加，欧几里得距离的作用就越小。

# 闵可夫斯基距离的一般形式

切比雪夫距离 Chebyshev Distance

当  $p = \infty$  时，闵可夫斯基距离为切比雪夫距离：

$$dist(x, y) = \max_d(|x_i - y_i|) = \lim_{k \rightarrow \infty} \left( \sum_{i=1}^d |p_i - q_i|^k \right)^{1/k}$$

理解：两个向量各个维度特征的数值差的最大值

# 闵可夫斯基距离的一般形式

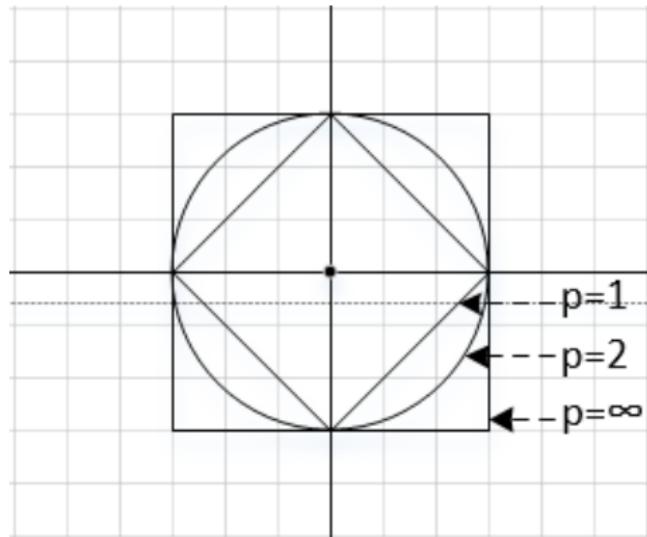
国际象棋中，国王可以直行、横行、斜行。国王走一步，可以移动到相邻的 8 个方格的任意一个。国王从格子  $(X_1, Y_1)$  到格子  $(X_2, Y_2)$  最少需要多少步？这个距离就是切比雪夫距离。

	a	b	c	d	e	f	g	h	
8	5	4	3	2	2	2	2	2	8
7	5	4	3	2	1	1	1	2	7
6	5	4	3	2	1	1	1	2	6
5	5	4	3	2	1	1	1	2	5
4	5	4	3	2	2	2	2	2	4
3	5	4	3	3	3	3	3	3	3
2	5	4	4	4	4	4	4	4	2
1	5	5	5	5	5	5	5	5	1
	a	b	c	d	e	f	g	h	

切比雪夫距离

# 闵可夫斯基距离的可视化

闵可夫斯基距离的可视化如图所示：



常见的闵可夫斯基距离的可视化结果

# 从向量内积到余弦相似度

向量内积是计算向量距离的基本方式

## 定义

$$\text{Inner}(x, y) = \langle x, y \rangle = \sum_i x_i y_i$$

向量内积是没有界限的，一种解决方法是除以长度后再求内积，即余弦相似度（Cosine Similarity）

## 定义

$$\text{Cos Sim}(x, y) = \frac{\sum_i x_i y_i}{\sqrt{\sum_i x_i^2} \sqrt{\sum_i y_i^2}} = \frac{\langle x, y \rangle}{\|x\| \|y\|}$$

余弦相似度与向量的大小无关，只与向量的方向有关。余弦相似度不满足度量测度性质，因此被称为非度量测度（nonmetric measure）

# 从余弦相似度到皮尔逊相关系数

余弦相似度面临的问题是不满足平移不变性:

$$\cos(\vec{a} + \vec{c}, \vec{b} + \vec{c}) \neq \cos(\vec{a}, \vec{b})$$

而皮尔逊相关系数则可以满足平移不变性

## 定义

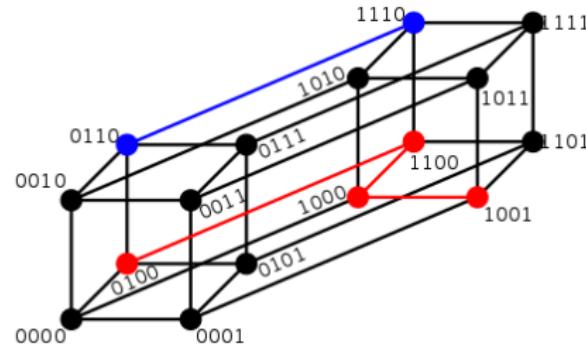
皮尔逊相关系数 (Pearson Correlation) 用于度量两组数据的变量 X 和 Y 之间的线性相关的程度, 定义为:

$$\begin{aligned}\text{Corr}(x, y) &= \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} = \frac{\langle x - \bar{x}, y - \bar{y} \rangle}{\|x - \bar{x}\| \|y - \bar{y}\|} \\ &= \text{Cos Sim}(x - \bar{x}, y - \bar{y})\end{aligned}$$

皮尔逊相关系数具有平移不变性和尺度不变性, 计算出了两个向量 (维度) 的相关性。

# 汉明距离 (Hamming Distance)

在信息理论中，Hamming Distance 表示两个等长字符串在对应位置上不同字符的数目。对于两个数字来说，汉明距离就是转成二进制后，对应的位置值不相同的个数。对于二进制串  $a$  和  $b$  来说，汉明距离等于  $a \text{ XOR } b$  中 1 的数目



汉明距离示例

# 非数值特征数据的关系度量

除了数据特征数据之间的关系度量，数据挖掘任务中还有很多特殊形式的数据：

- ① 集合数据  $\Rightarrow$  Jaccard 距离
- ② 概率分布  $\Rightarrow$  Divergence (散度)
  - ① KL-Divergence
  - ② JS-Divergence
- ③ Wasserstein Distance

# 常见距离度量的代码实现

目前大部分经典距离度量都已有开源实现，在 Sklearn<sup>7</sup>，Pytorch 等框架中已经原生支持，可以直接调用。

<code>metrics.pairwise.additive_chi2_kernel(X[, Y])</code>	Computes the additive chi-squared kernel between observations in X and Y.
<code>metrics.pairwise.chi2_kernel(X[, Y, gamma])</code>	Computes the exponential chi-squared kernel X and Y.
<code>metrics.pairwise.cosine_similarity(X[, Y, ...])</code>	Compute cosine similarity between samples in X and Y.
<code>metrics.pairwise.cosine_distances(X[, Y])</code>	Compute cosine distance between samples in X and Y.
<code>metrics.pairwise.distance_metrics()</code>	Valid metrics for pairwise_distances.
<code>metrics.pairwise.euclidean_distances(X[, Y, ...])</code>	Considering the rows of X (and Y=X) as vectors, compute the distance matrix between each pair of vectors.
<code>metrics.pairwise.haversine_distances(X[, Y])</code>	Compute the Haversine distance between samples in X and Y.
<code>metrics.pairwise.kernel_metrics()</code>	Valid metrics for pairwise_kernels.
<code>metrics.pairwise.laplacian_kernel(X[, Y, gamma])</code>	Compute the laplacian kernel between X and Y.
<code>metrics.pairwise.linear_kernel(X[, Y, ...])</code>	Compute the linear kernel between X and Y.
<code>metrics.pairwise.manhattan_distances(X[, Y, ...])</code>	Compute the L1 distances between the vectors in X and Y.
<code>metrics.pairwise.nan_euclidean_distances(X)</code>	Calculate the euclidean distances in the presence of missing values.
<code>metrics.pairwise.pairwise_kernels(X[, Y, ...])</code>	Compute the kernel between arrays X and optional array Y.
<code>metrics.pairwise.polynomial_kernel(X[, Y, ...])</code>	Compute the polynomial kernel between X and Y.
<code>metrics.pairwise.rbf_kernel(X[, Y, gamma])</code>	Compute the rbf (gaussian) kernel between X and Y.
<code>metrics.pairwise.sigmoid_kernel(X[, Y, ...])</code>	Compute the sigmoid kernel between X and Y.
<code>metrics.pairwise.paired_euclidean_distances(X, Y)</code>	Computes the paired euclidean distances between X and Y.
<code>metrics.pairwise.paired_manhattan_distances(X, Y)</code>	Compute the L1 distances between the vectors in X and Y.
<code>metrics.pairwise.paired_cosine_distances(X, Y)</code>	Computes the paired cosine distances between X and Y.
<code>metrics.pairwise.paired_distances(X, Y, *[...])</code>	Computes the paired distances between X and Y.
<code>metrics.pairwise_distances(X[, Y, metric, ...])</code>	Compute the distance matrix from a vector array X and optional Y.
<code>metrics.pairwise_distances_argmin(X, Y, *[...])</code>	Compute minimum distances between one point and a set of points.
<code>metrics.pairwise_distances_argmin_min(X, Y, *)</code>	Compute minimum distances between one point and a set of points.
<code>metrics.pairwise_distances_chunked(X[, Y, ...])</code>	Generate a distance matrix chunk by chunk with optional reduction.

<sup>7</sup><https://scikit-learn.org/stable/modules/classes.html#module-sklearn.metrics>

//scikit-learn.org/stable/modules/classes.html#module-sklearn.metrics ↻ ↻ ↻

# 课程内容

## 1 课程简介

- 课程信息
- 准备知识

## 2 数据挖掘简介

- 从数据到知识
- 数据为中心的人工智能
- 数据挖掘的定义
- 数据挖掘与学术
- 课程大纲

## 3 数据类型与度量

- 数据类型
- 数据的相关性度量

## 4 数据的统计与可视化

- 数据的统计描述
- 数据可视化

## 5 数据预处理与数据降维

- 数据质量
- 主成分分析 PCA

# 数据的统计描述

用于描述数据整体特征的表征通常可以分为两大类：

- ① 中心趋势性指标
- ② 数据分散性指标

## 中心趋势性指标

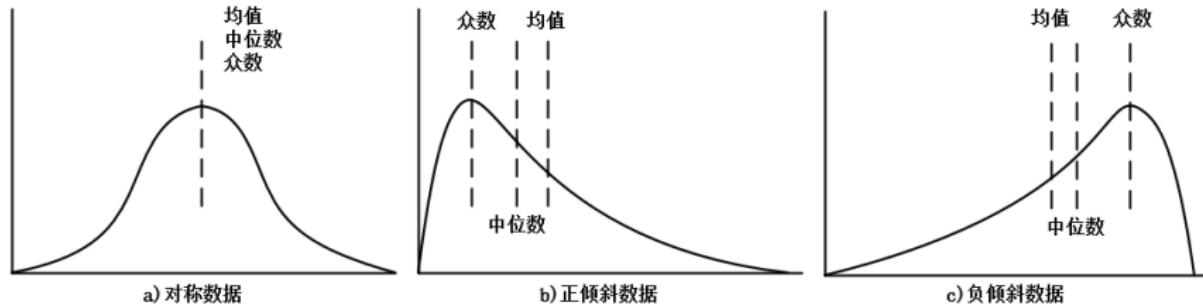
- ① 均值（加权均值，截尾均值）
- ② 中位数
- ③ 众数
- ④ 中列数：数据集的最大值和最小值的平均值

## 数据分散性指标

- ① 极差
- ② 分位数
  - ① 四分位数
  - ② 百分位数
  - ③ 四分位数极差
- ③ 方差/标准差

# 数据的统计描述

虽然每个数据集都可以得到对应的统计描述，但是不同的数据分布应由不同的统计指标进行度量。



不同数据分布对统计信息的影响

 新浪科技 新浪科技>互联网>腾讯2021年Q4及全年财报专题 > 正文

新闻

请输入关键词



腾讯11万员工去年人均年薪84.7万元，同比上涨3.59万元

# 数据统计的可视化

常见的数据可视化工具：

- ① Excel
- ② MATLAB
- ③ Matplotlib (Python)
- ④ Seaborn (Python)
- ⑤ Tikz (Latex)

数据可视化也是科研的重要部分

# 数据统计的可视化

常见的数据可视化工具：

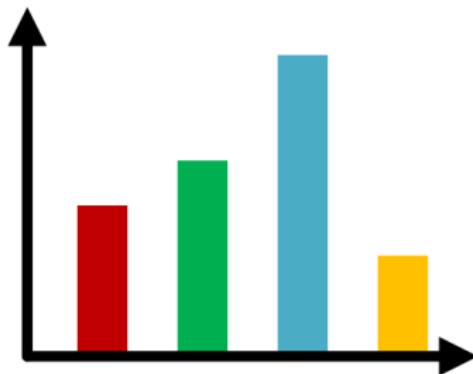
- ① Excel
- ② MATLAB
- ③ Matplotlib (Python)
- ④ Seaborn (Python)
- ⑤ Tikz (Latex)

数据可视化工具很多，核心是可视化的目标是否明确，呈现是否清晰简洁。盲目的可视化只会引起误解。

数据可视化也是科研的重要部分

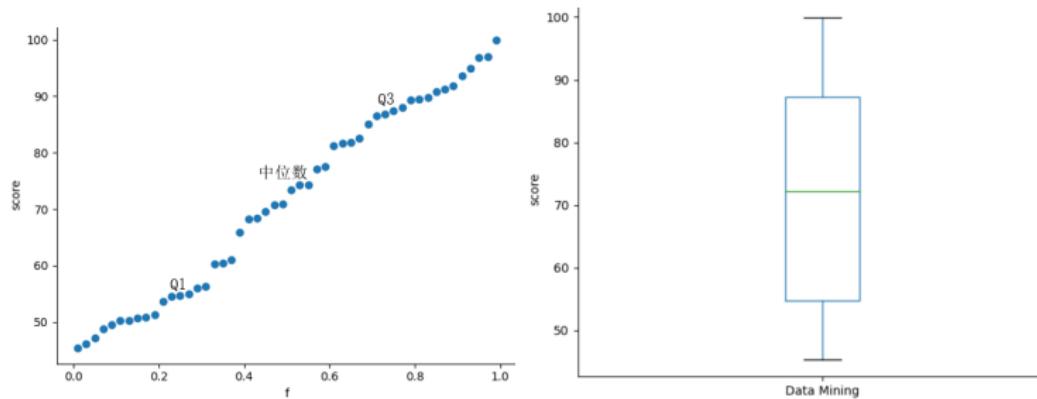
# 数据统计的可视化

为了展示数据集中数据的频率变化以及主要高频分量，通常使用柱状图和词云图等进行可视化。



# 数据统计的可视化

为了展示数据的分布以及分布的关键分段指标，通常使用分位数和盒图等进行可视化。



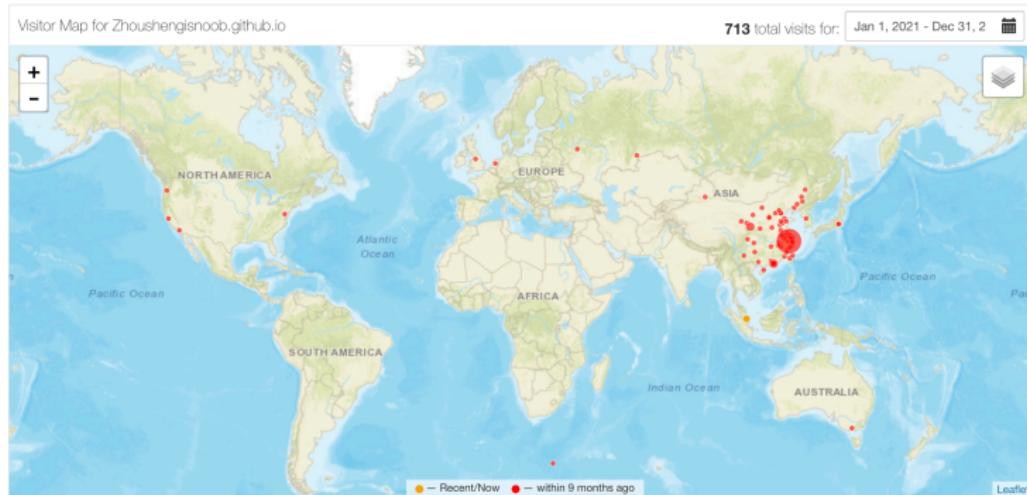
# 数据统计的可视化

为了展示数据的全局数值分布，通常使用热力图进行可视化。



# 数据统计的可视化

为了展示数据的地理分布，可以使用地理数据可视化工具。



# 课程内容

## 1 课程简介

- 课程信息
- 准备知识

## 2 数据挖掘简介

- 从数据到知识
- 数据为中心的人工智能
- 数据挖掘的定义
- 数据挖掘与学术
- 课程大纲

## 3 数据类型与度量

- 数据类型
- 数据的相关性度量

## 4 数据的统计与可视化

- 数据的统计描述
- 数据可视化

## 5 数据预处理与数据降维

- 数据质量
- 主成分分析 PCA

# 数据的质量

高质量的数据是高效数据挖掘的前提和保障，而真实场景中数据质量却往往难以满足实用需求：**GIGO: garbage in garbage out**

## 数据质量的三大问题

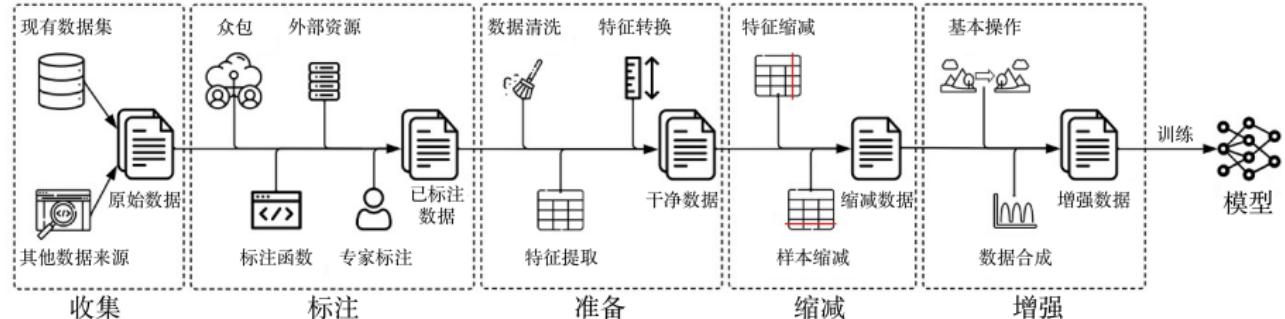
- ① 准确性（数据采集软硬件故障，人为故障，用户主观不愿意被采集数据，传输过程的信息损失）
- ② 完整性（数据的特征采集由人工制定）
- ③ 一致性（不同数据源的数据形式不统一）

数据预处理是数据挖掘的必要操作，也是数据分析、算法工程师重要的日常工作之一。

# 数据预处理的主要步骤

数据预处理主要包含如下的步骤：

- ① 数据清理（Data Cleaning）⇒ 去除缺失，噪声数据
- ② 数据集成（Data Integration）⇒ 去除冗余的数据和特征
- ③ 数据归约（Data Reduction）⇒ 降低数据集的规模
- ④ 数据变换（Data Transformation）⇒ 提升数据挖掘效率



# 数据清理

数据清理通常包含如下的操作：填补缺失值，识别并去除噪声数据。

填补缺失值的方式：

- ① 手工填写
- ② 全局常量
- ③ 统计变量（均值，中位数）
- ④ 丢弃属性

清理噪声数据的方式：

- ① 统计数据识别
- ② 规则识别噪声
- ③ 异常检测算法识别噪声
- ④ 回归方式识别噪声

Q: 离群点是不是就是异常点？

# 数据集成

数据集成的主要目的是将不同数据源的数据整合到统一的数据中，其中核心的问题是避免数据的冗余，判断是否冗余则主要依赖相关分析。

## 数值数据的相关分析

### ① 相关系数

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n\sigma_A\sigma_B}$$

### ② 协方差矩阵

$$\begin{aligned} Cov(A, B) &= E((A - \bar{A})(B - \bar{B})) \\ &= \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n} \end{aligned}$$

# 数据归约

数据归约的主要目的是减少数据量，但仍然接近于保持原始数据的完整性。

数据归约的主要方式包括：

① 维度归约 (dimension reduction)

- ① 小波变换
- ② 主成分分析
- ③ 表征学习

② 数量归约 ( numerosity reduction )

- ① 聚类
- ② 采样

③ 数据压缩 (data compression)

# 数据变换

## 为什么要做数据变换

原始数据（Raw Data）无法直接被计算机存储和识别，且直接处理原始数据往往不利于数据挖掘的开展

### 常见的数据变换方法：

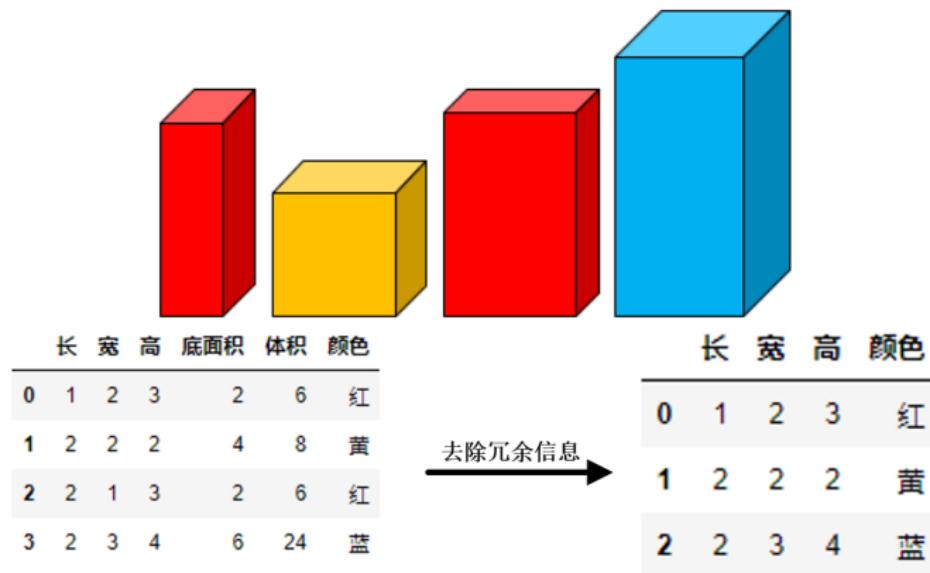
- ① 标准化（Normalization）
- ② 数据采样
- ③ 数据的编码

### 常见的数据变换案例：

- ① 离散类别的编码
- ② 图像、音视频、文本、图数据的编码

# 为何要做主成分分析？

特征冗余：在多维特征中，有些维度没有意义——听君一席话，如听一席话



使用较少的维度描述数据，利于存储和数据挖掘。

# 内积与正交基

向量的内积操作讲两个向量映射为实数：

$$(a_1, a_2, \dots, a_n)^\top \cdot (b_1, b_2, \dots, b_n) = a_1b_1 + a_2b_2 + \dots + a_nb_n$$

以二维空间为例，向量的内积可以理解为两个向量的模长乘积乘以两个向量的夹角余弦：

$$A \cdot B = |A||B| \cos(a) = |A| \cos(a)|B|$$

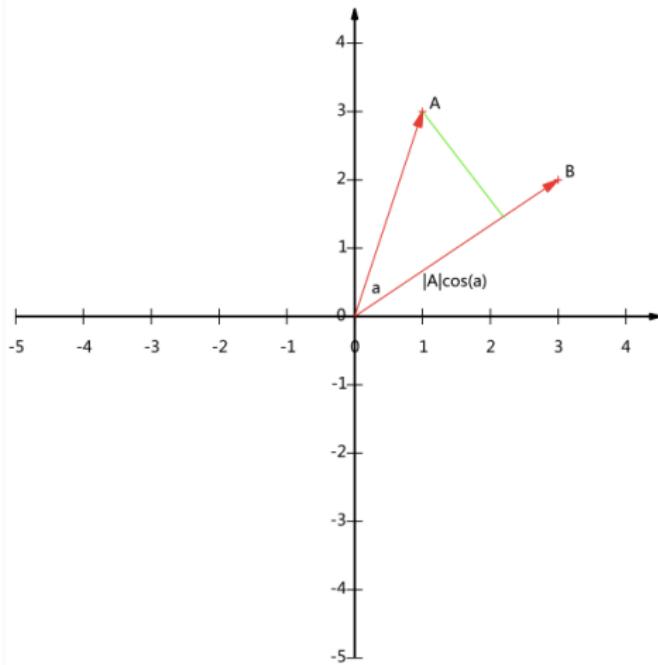
其中  $|A|$  是向量  $A$  的模（标量长度）， $a$  是两个向量之间的夹角。 $|A|\cos(a)$  称为向量  $A$  投影的矢量长度。

## 标量长度与矢量长度

标量长度总是大于等于 0，值就是线段的长度；而矢量长度可能为负，其绝对值是线段长度，而符号取决于其方向与标准方向相同或相反。

# 内积与正交基

二维空间中的内积操作：



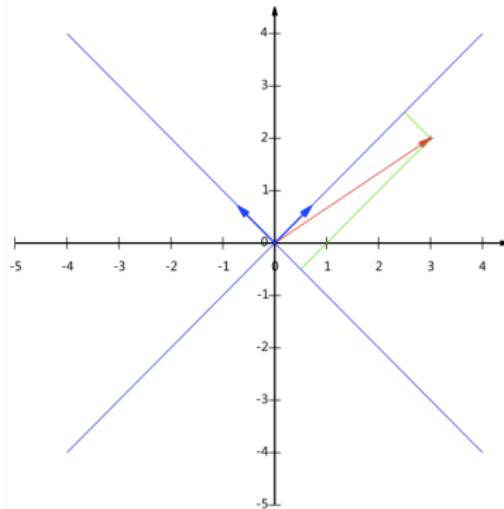
## 基向量

如果向量  $B$  的模为 1，则向量  $A$  和向量  $B$  的内积等于  $A$  向  $B$  所在直线投影的矢量长度。此时向量  $B$  也称为基向量。

## 向量的内积操作

# 内积与正交基

要准确描述向量，首先要确定一组基，然后给出在基所在的各个直线上的投影值。任何两个线性无关（不在同一直线上）的二维向量都可以成为一组基，但是我们通常选取正交的一组向量作为基，也称为正交基。



二维空间中的两组正交基<sup>8</sup>

<sup>8</sup><https://www.cnblogs.com/wi-1314/p/8032780.html>

# 基变换

对于原空间的一个向量，如果用在使用新基表示的空间中重新描述该向量，只需将基向量对应的矩阵成一原向量，就可以作为基变换之后新的坐标。单向量的基变换：

$$\begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} 5/\sqrt{2} \\ -1/\sqrt{2} \end{pmatrix}$$

多向量的基变换：

$$\begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix} \begin{pmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \end{pmatrix} = \begin{pmatrix} 2/\sqrt{2} & 4/\sqrt{2} & 6/\sqrt{2} \\ 0 & 0 & 0 \end{pmatrix}$$

# 基变换与数据降维

基变换可以将原空间中的数据，变换到新的空间中，并用新空间的基进行表示。（向量的维度有新空间中基的个数决定）

如果新空间的基的个数少于原空间中基的个数，那么基变换就实现了数据降维

## 什么是最优的数据降维？

高维空间中可以存在无数组基，对应无数种数据降维方式，但是哪一种数据降维方式是最优的呢？

# 主成分分析的两种目标

## 最大可分性

在正交属性空间中，存在这样的一个超平面，使得样本点在这个超平面上的投影尽可能分开

从信息熵的角度，数据足够分开，整个数据集包含的信息量越大。

## 最近重构性

在正交属性空间中，存在这样的一个超平面，使得样本点到超平面的距离足够接近

PCA 的目标转化为寻找最优的正交基，使得上述的两种性质可以被满足。

# 数据在属性维度的方差

给定某一属性维度，所有数据在这个维度上的投影分量组成了属性的向量表示。数据投影的离散程度可用这个该属性维度的方差表示：

## 属性维度的方差

给定一个  $n$  维的属性向量，其方差可以定义为：

$$\text{Var}(a) = \frac{1}{n} \sum_{i=1}^n (a_i - \bar{a})^2$$

如果该属性维度进行了标准化，则方差可以简化为：

$$\text{Var}(a) = \frac{1}{n} \sum_{i=1}^n a_i^2$$

# 属性的协方差

两个属性维度之间的相关性，通常用协方差来表示。

## 属性间的协方差

给定两个  $n$  维的向量，向量之间的协方差可以定义为

$$Cov(a, b) = \frac{\sum_{i=1}^n (a_i - \bar{a})(b_i - \bar{b})}{n}$$

如果这两个向量进行了标准化，则协方差可以简化为：

$$Cov(a, b) = \frac{\sum_{i=1}^n a_i b_i}{n}$$

# 主成分分析优化目标

## 目标

给定一组  $m$  维向量降维  $k$  维 ( $0 < k < m$ )，主成分分析目标是选择  $k$  个单位（模为 1）正交基，使得原始数据变换到这组基上后，各属性两两间协方差为 0，而字段的方差则尽可能大

一句话总结：在正交的约束下，取方差最大的  $K$  个属性维度

# 矩阵形式目标

上述优化目标，能否使用统一的形式进行表示？

# 矩阵形式目标

上述优化目标，能否使用统一的形式进行表示？  
给定  $n$  个二维数据（向量）

$$X = \begin{pmatrix} a_1 & a_2 & \cdots & a_n \\ b_1 & b_2 & \cdots & b_n \end{pmatrix}$$

该矩阵乘以自身的转置为：

$$\frac{1}{n} XX^\top = \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n a_i^2 & \frac{1}{n} \sum_{i=1}^n a_i b_i \\ \frac{1}{n} \sum_{i=1}^n a_i b_i & \frac{1}{n} \sum_{i=1}^n b_i^2 \end{pmatrix}$$

一个矩阵同时包含了单个属性维度的方差（对角线元素）和不同属性维度之间的协方差（非对角线元素）

# 基变换与协方差矩阵推导

给定特征矩阵  $X$  以及对应的协方差矩阵  $C$ ,  $P$  是一组基向量,  $Y$  是  $X$  经过  $P$  对应的基变换得到的特征矩阵, 其对应的协方差矩阵为  $D$ , 则  $D$  与  $C$  的关系如下:

$$\begin{aligned} D &= \frac{1}{n}YY^\top \\ &= \frac{1}{n}(PX)(PX)^\top \\ &= \frac{1}{n}PXX^\top P^\top \\ &= P\left(\frac{1}{n}XX^\top\right)P^\top \\ &= PCP^\top \end{aligned}$$

# 主成分分析优化目标

## 矩阵形式优化目标

寻找一个矩阵  $P$ , 满足  $PCP^\top$  是一个对角矩阵, 并且对角元素按照从大到小依次排列, 那么  $P$  的前  $k$  行就是要寻找的基。用  $P$  的前  $k$  行组成的矩阵乘以  $X$  就使得  $X$  从  $m$  维降到了  $k$  维并满足上述优化条件。

如何快速找到这样的矩阵呢?

# 实对称矩阵的数学性质

一个  $n$  行  $n$  列的实对称矩阵  $C$  一定可以找到  $n$  个单位正交特征向量，设这  $n$  个特征向量为  $e_1, e_2, \dots, e_n$ ，将其按照列组成矩阵：

$$E = (e_1 \quad e_2 \quad \cdots \quad e_n)$$

该矩阵具有如下的性质：

$$E^\top C E = \Lambda = \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix}$$

其对角元素为各特征向量对应的特征值。

# 主成分分析的矩阵形式

根据上述推导，我们可以得到主成分分析实现数据降维过程为：

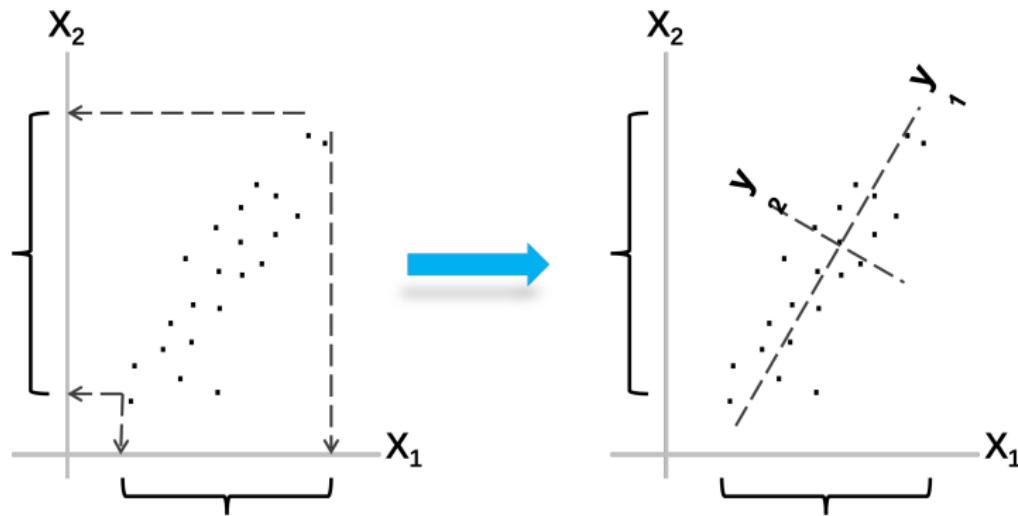
- ① 对所有样本进行中心化：

$$\mathbf{x}_i \leftarrow \mathbf{x}_i - \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$$

- ② 计算样本的协方差矩阵  $\mathbf{X}\mathbf{X}^\top$
- ③ 对协方差矩阵进行特征分解  $\mathbf{X}\mathbf{X}^\top$
- ④ 取  $k$  个最大的特征值对应的特征向量并组成矩阵

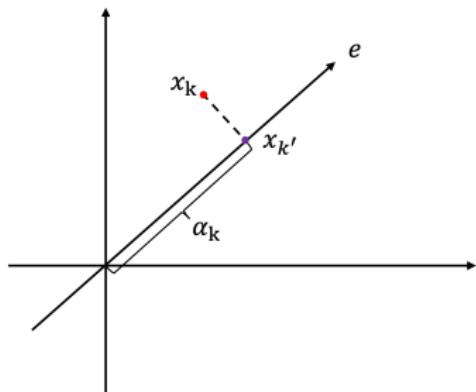
$$\mathbf{E} = (e_1 \quad e_2 \cdots e_n)$$

# 基于最近重构的主成分分析



# 基于最近重构的主成分分析

$$\begin{aligned}
 J(e) &= \sum_{i=1}^n \|x'_k - x_k\|^2 = \sum_{i=1}^n \|\alpha_k e - x_k\|^2 = \sum_{i=1}^n (\alpha_k e - x_k)^T (\alpha_k e - x_k) \\
 &= \sum_{i=1}^n (\alpha_k e^T - x_k^T)(\alpha_k e - x_k) \\
 &= \sum_{i=1}^n \alpha_k^2 \|e\|^2 + \sum_{i=1}^n \|x_k\|^2 \\
 &\quad - \sum_{i=1}^n \alpha_k e^T x_k - \sum_{i=1}^n \alpha_k x_k^T e \\
 &= \sum_{i=1}^n \alpha_k^2 \|e\|^2 + \sum_{i=1}^n \|x_k\|^2 - 2 \sum_{i=1}^n \alpha_k e^T x_k \\
 &= \sum_{i=1}^n \alpha_k (\alpha_k - 2e^T x_k) + \sum_{i=1}^n \|x_k\|^2
 \end{aligned}$$



# 基于最近重构的主成分分析

$$= - \sum_{i=1}^n \alpha_k^2 + \sum_{i=1}^n \|x_k\|^2 = - \sum_{i=1}^n (e^T x_k)^2 + \sum_{i=1}^n \|x_k\|^2$$

$$= - \sum_{i=1}^n e^T x_k x_k^T e + \sum_{i=1}^n \|x_k\|^2$$

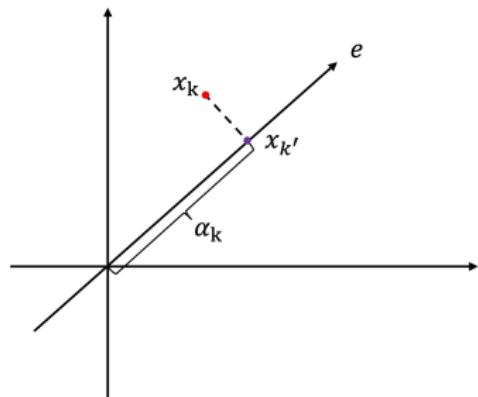
令  $S = \sum_{i=1}^n x_k x_k^t$

$$\max_e e^t S e \quad \text{s.t. } \|e\| = 1$$

使用拉格朗日乘子进行优化：

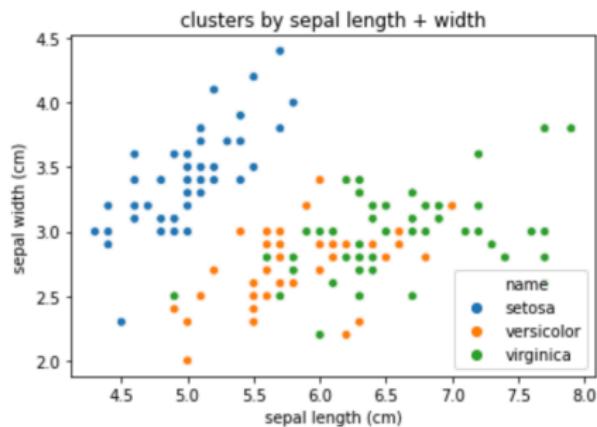
$$u = e^t S e - \lambda (e^t e - 1)$$

$$\frac{\partial u}{\partial e} = 2Se - 2\lambda e$$

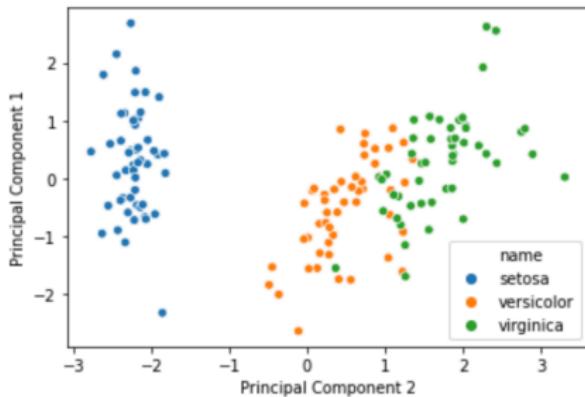


# PCA 实现

一行代码: `sklearn.decomposition.PCA`



原始特征

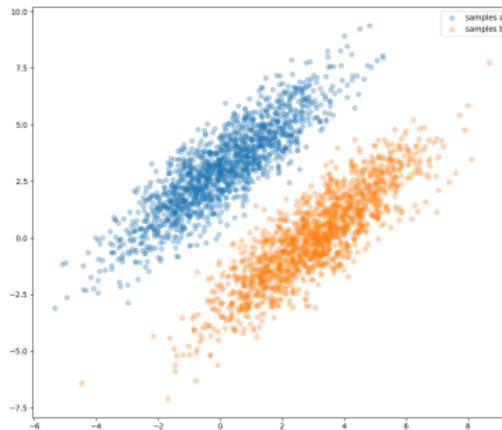


PCA 特征

# PCA 的主要缺陷

PCA 是经典的数据降维方法，但是仍然面临如下缺陷：

- ① 计算复杂度高
- ② 维度诅咒
- ③ 只能处理特征向量数据



PCA 并不能总是满足需求

# 经典的数据降维技术

经典的数据降维技术包括：

- ① 主成分分析 (Principle component analysis)
- ② 等距特征映射 (Isometric Mapping)
- ③ 线性判别分析 (Linear Discriminant Analysis)
- ④ T-SNE (t-distributed stochastic neighbor embedding)

经典的深度学习降维技术包括：

- ① Auto-Encoder
- ② Deep Generative Model
- ③ Unsupervised Learning
- ④ Self-supervised Learning

# 要点总结

本次课程的核心知识点如下：

- ① 数据挖掘的基本概念
- ② 常见的数据形式
- ③ 常见的数据特征关系度量
- ④ 主成分分析

# 参考文献



WARNAT-HERRESTHAL, S., SCHULTZE, H., SHAstry, K. L.,  
MANAMOHAN, S., MUKHERJEE, S., GARG, V., SARVESWARA, R.,  
HÄNDLER, K., PICKKERS, P., AZIZ, N. A., ET AL.

Swarm learning for decentralized and confidential clinical machine  
learning.

*Nature* 594, 7862 (2021), 265–270.