



《人工智能及应用》

第六讲：机器学习-无监督学习



授课人：张鑫 zhangxin@uestc.edu.cn

专 业：机器人工程

提纲



1、无监督学习的概念

2、K均值聚类

3、主成份分析

1、有监督学习与无监督学习

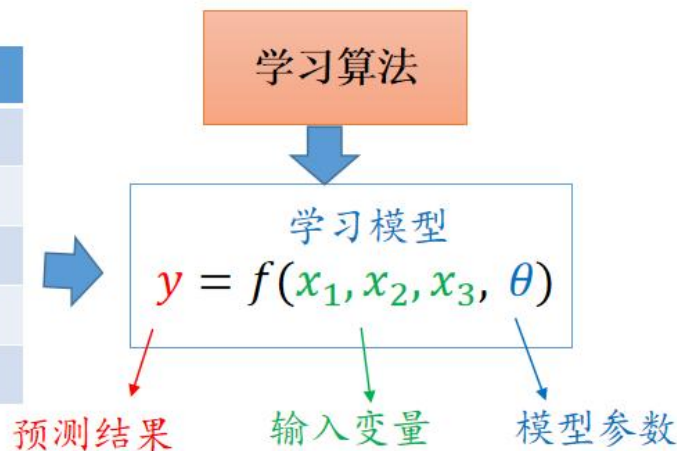


有监督学习

Supervised learning

训练集
(带标签)

	x_1	x_2	x_3	y
编号	就业	教育水平	居住时间	信用卡?
1	是	研究生	5	Yes
2	是	高中	2	No
3	否	大学	1	No
4	是	高中	10	Yes
...				



从带**标签**的训练数据中，学习从数据到其标签的映射关系

数据的标注作为**监督**信号

最优化目标为预测结果与标注尽可能一致

1、有监督学习与无监督学习



有监督学习

Supervised learning

训练集
(带标注)

	x_1	x_2	x_3	y
编号	就业	教育水平	居住时间	信用卡?
1	是	研究生	5	Yes
2	是	高中	2	No
3	否	大学	1	No
4	是	高中	10	Yes
...				

推理阶段

编号	就业	教育水平	居住时间	信用卡?
1	否	研究生	3	?
2	是	大学	2	?
...				

归纳



学习算法



学习模型

$$y = f(x_1, x_2, x_3, \theta)$$

应用模型

(演绎)

$$y = f(x_1, x_2, x_3, \theta)$$

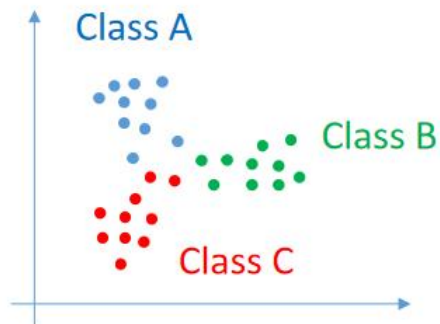
1、有监督学习与无监督学习



有监督学习

Supervised learning

训练数据带标签，在
标签的监督下，发掘
数据与标签之间的映
射关系



老虎



狮子



大象

1、有监督学习与无监督学习

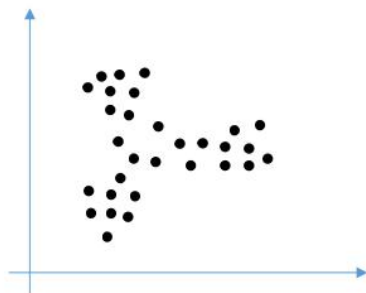


无监督学习

Unsupervised learning

训练数据**无标签**，如何仅利用数据本身的信息，发掘数据样本中蕴含的模式？

- **无监督学习**：是从无标注的数据样本中，学习数据中蕴含的模式，完成如**聚类**或**降维**等特定任务。



人类学习中，无监督学习占据主导地位：通过观察客观事务去洞悉其内在的规律。



1、无监督学习的概念

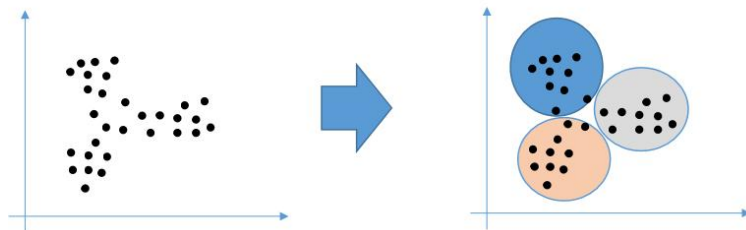
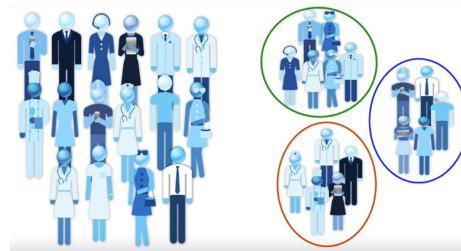
1.1 概念——聚类、降维



- 1. **聚类**：根据事物的相似程度，将他们分成不同类别。

物以类聚、人以群分。

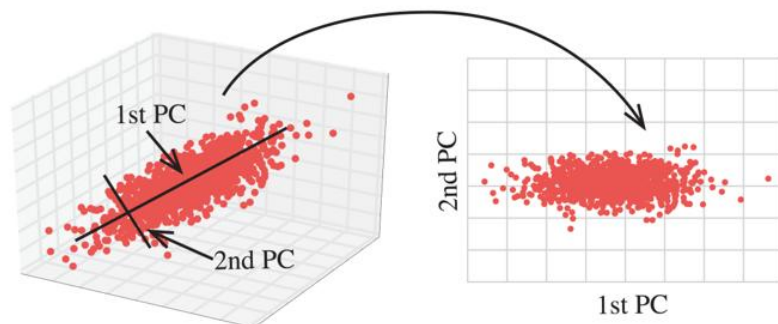
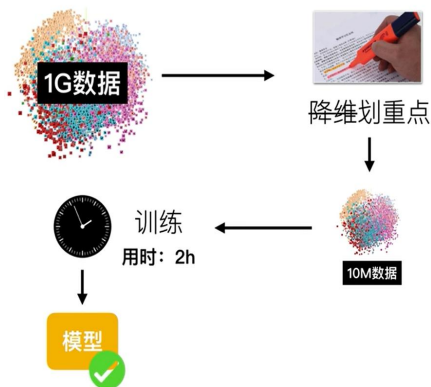
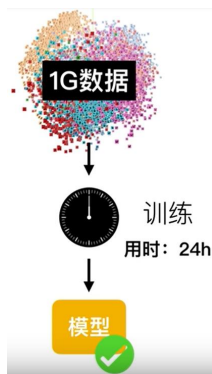
- 目标用户的群体分类
- 图像切割
- 基因聚类



1、无监督学习的概念

1.1 概念——聚类、降维

- 2. 降维：实际上就是划重点、提精华。



Principal Components Analysis (PCA): 主成份分析

提纲



1、无监督学习的概念

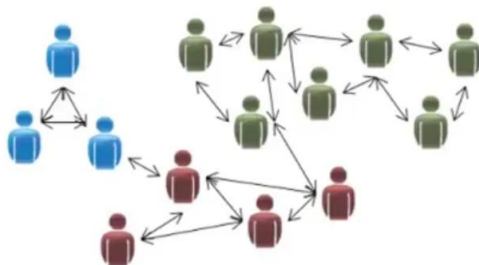
2、K均值聚类

3、主成份分析

2、K均值聚类

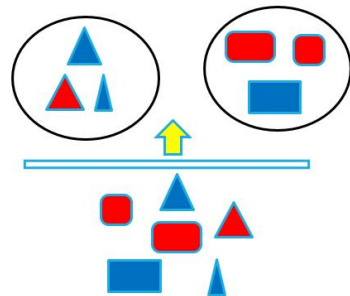
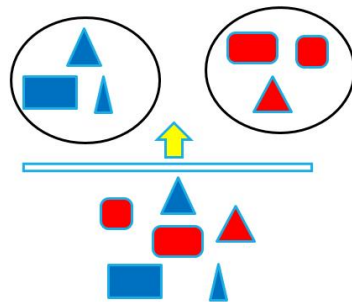
2.1 概念

- 物以类聚，人以群分（《战国策·齐策三》）
 - 目的：将 n 个数据聚类到 k 个集合（也称为类簇）
 - 输入： n 个数据（无任何标注信息），聚类数 k
 - 输出： k 个聚类结果
-
- 换个说法：如何将 n 个数据依据**相似程度**将它们分别聚类到 k 个集合，使得每个数据仅属于一个聚类集合。



对人群进行社交关系分类，
自动推荐可能认识的人

社交关系分组



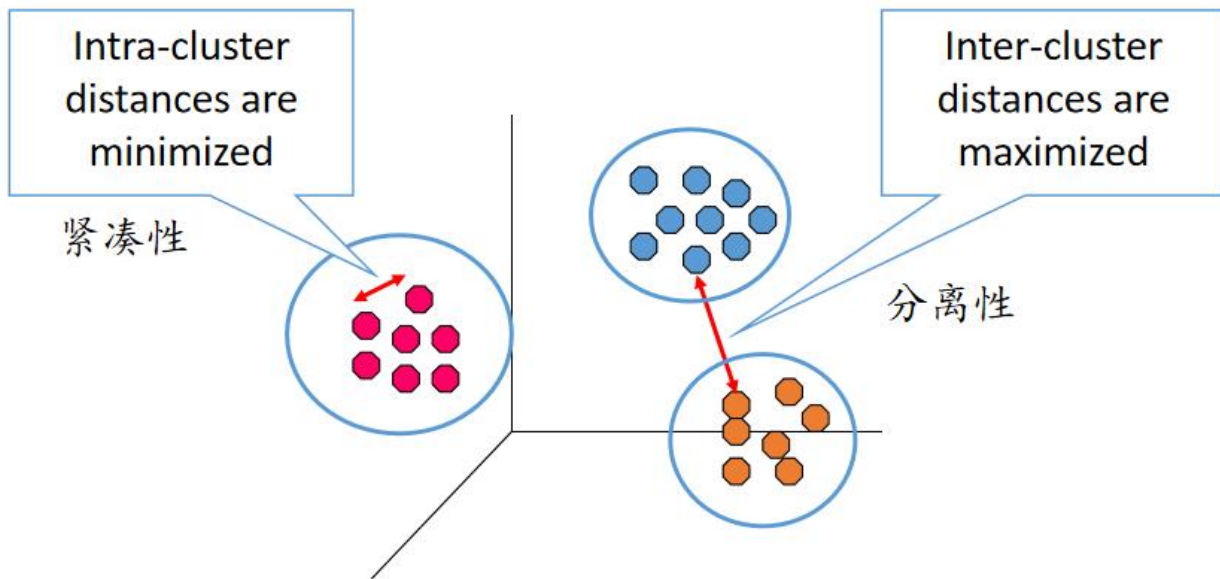
2、K均值聚类

2.1 概念



聚类的特点

- 同一个聚类中的样本相似（或距离近）
- 不同聚类间的样本不相似（或距离远）



2、K均值聚类

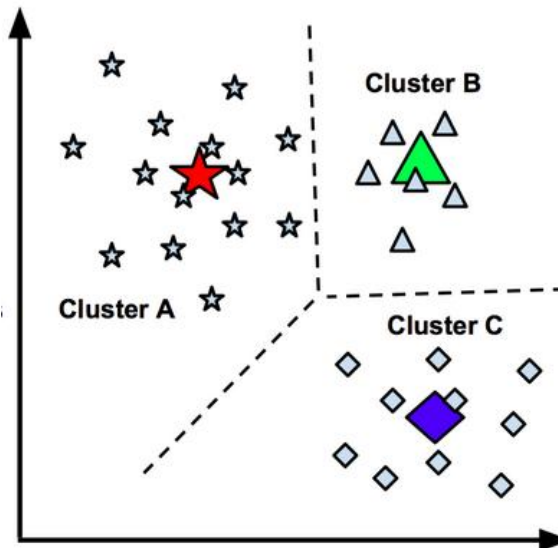
2.1 概念

是一种基于中心的分割型聚类方法

- 聚类个数 k 需指定
- 每个聚类有一个对应的中心点 (Centroid)
- 每个样本被划分到聚类最近的中心点

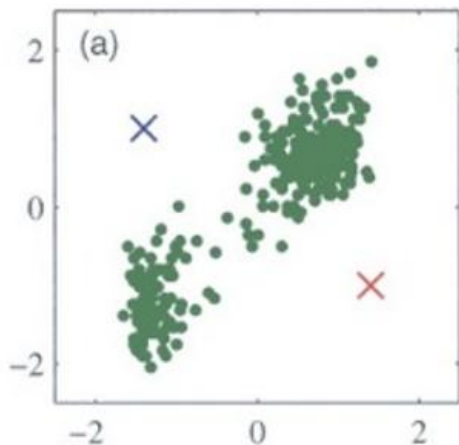
算法基本过程

- 1: Select K points as the initial centroids.
- 2: **repeat**
- 3: Form K clusters by assigning all points to the closest centroid.
- 4: Recompute the centroid of each cluster.
- 5: **until** The centroids don't change



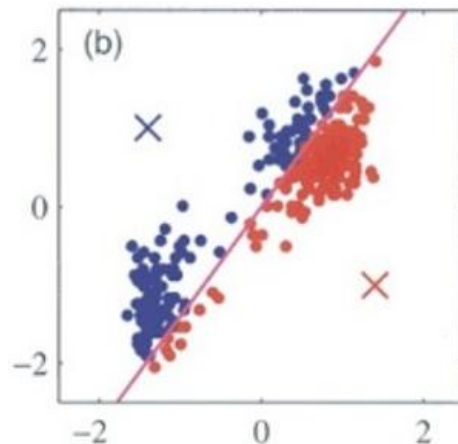
2、K均值聚类

2.2 算法



Initial centroids

初始中心

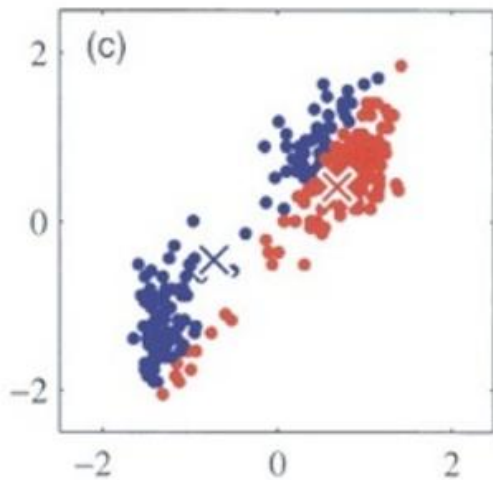


Initial clusters

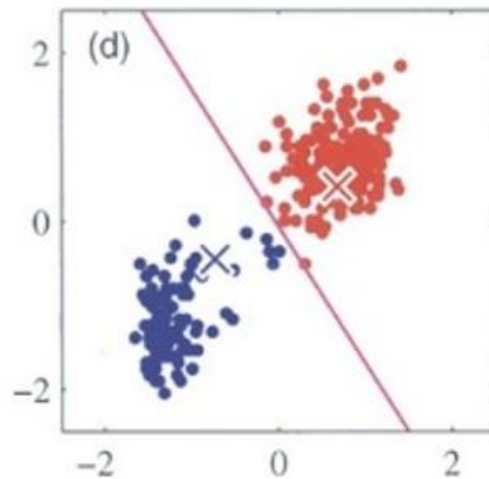
初始聚类

2、K均值聚类

2.2算法



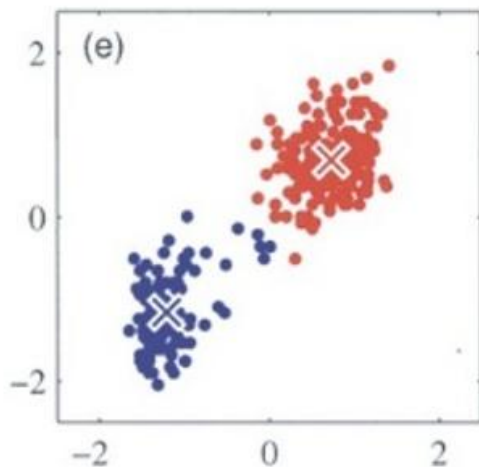
更新中心



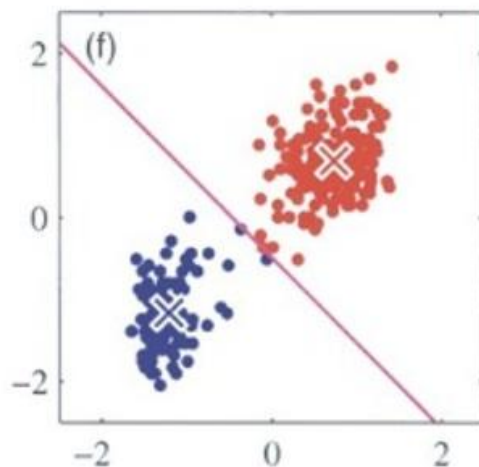
更新聚类

2、K均值聚类

2.2 算法



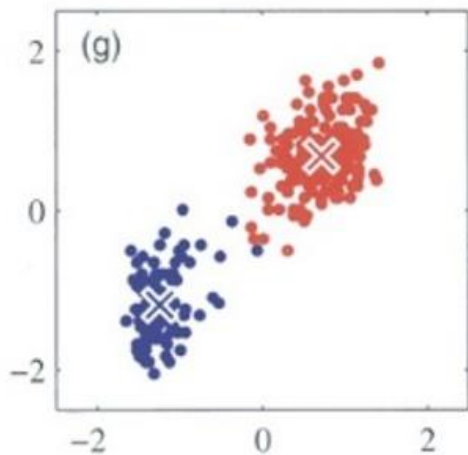
更新中心



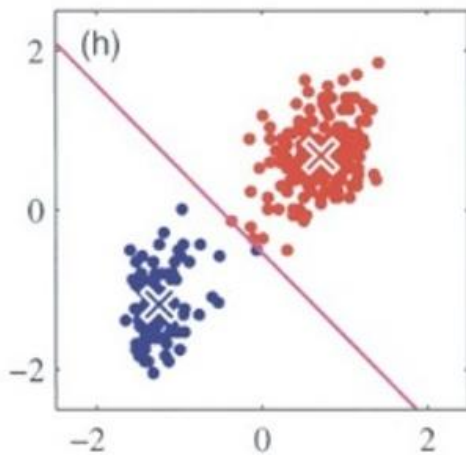
更新聚类

2、K均值聚类

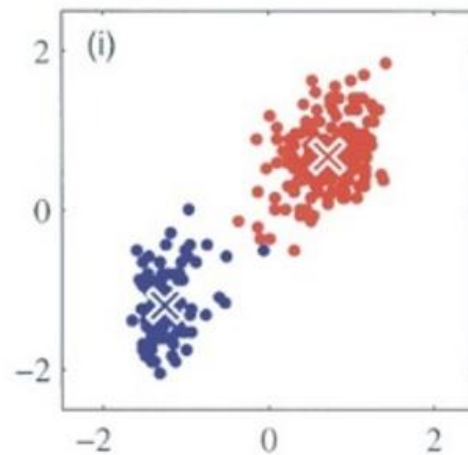
2.2 算法



更新中心



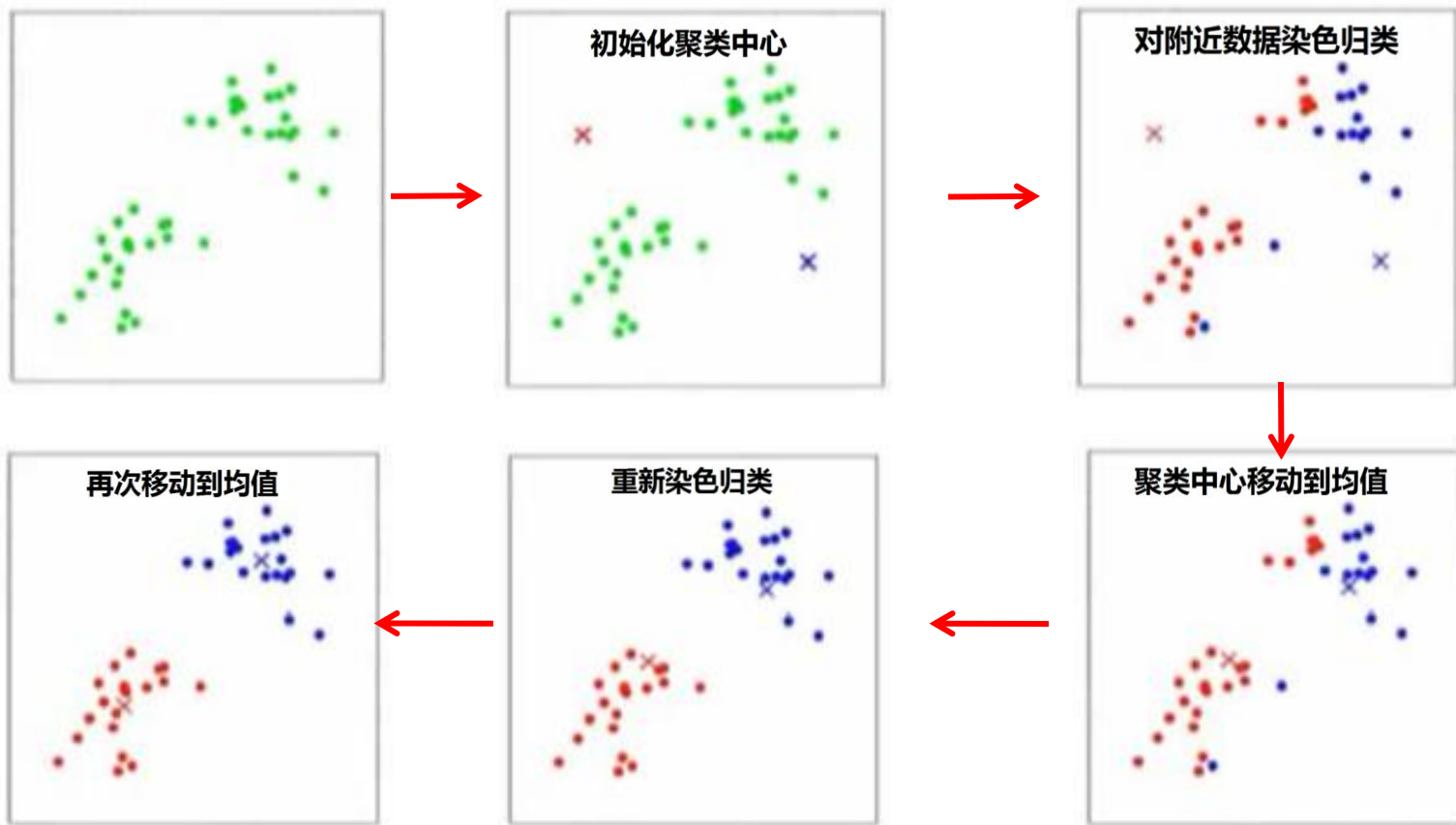
更新聚类



最终结果

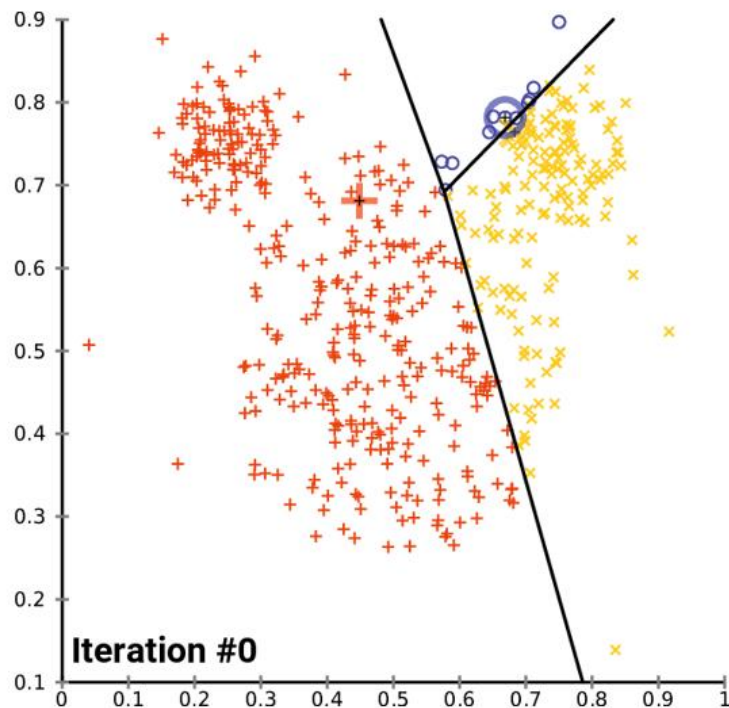
2、K均值聚类

2.2 算法



2、K均值聚类

2.2算法



输入：聚类个数 K

Step1: 设有 K 个聚类 $C = \{C_1, C_2, \dots, C_K\}$, 随机选择 K 个样本分别作为 K 个聚类的中心点 $\{m_1, m_2, \dots, m_K\}$

Step2: 对每个样本 x_j , 将其分配到距离最近的聚类中心所对应的聚类

Step3: 对第 i 个聚类, 重新计算其中心

$$m'_i = \frac{1}{n_i} \sum_{x_j \in C_i} x_j$$

其中 n_i 为聚类 C_i 中的样本个数

Step4: 重复step2和step3, 直到满足收敛条件 (达到指定的迭代次数, 聚类中心不再变化等)

输出：聚类结果

2、K均值聚类

2.2算法



- 相似度定义：用**欧式距离**描述数据之间的相似性。
- 有n个数据 $\{x_1, x_2, \dots, x_n\}$, $(1 \leq i \leq n)$
- 每个数据是m维的, $x_i = [x_{i1} \quad x_{i2} \quad \dots \quad x_{im}]^T$
- 两个m维数据之间的**欧氏距离**定义为:

$$d(x_i, x_j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{im} - x_{jm})^2}$$

2、K均值聚类

2.2 算法



- 需要事先确定**聚类数目K**，很多时候我们并不知道数据应被聚类的数目
- 需要**初始化聚类质心**，初始化聚类中心对聚类结果有较大的影响
- 算法是迭代执行，时间开销非常大
- 欧氏距离假设数据每个维度之间的重要性是一样的

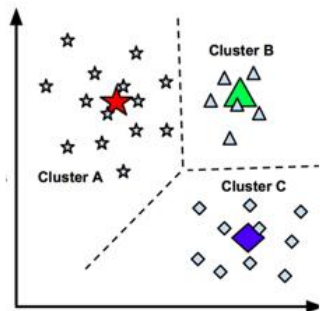
2、K均值聚类

2.2算法

需要事先确定聚类数目K



- 从目标优化的角度看待：误差平方和（Sum of Squared Error, SSE）

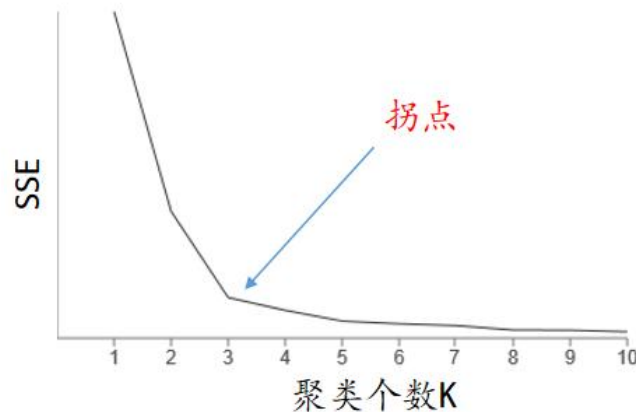


$$SSE = \sum_{i=1}^K \sum_{x \in C_i} dist^2(m_i, x)$$

x 是聚类 C_i 中的一个样本

m_i 是聚类 C_i 的中心

- 聚类的目标是实现SSE的最小化
 - 前提是K需要事先给定
 - 如何自动确定K的值？
 - 根据SSE的拐点



2、K均值聚类

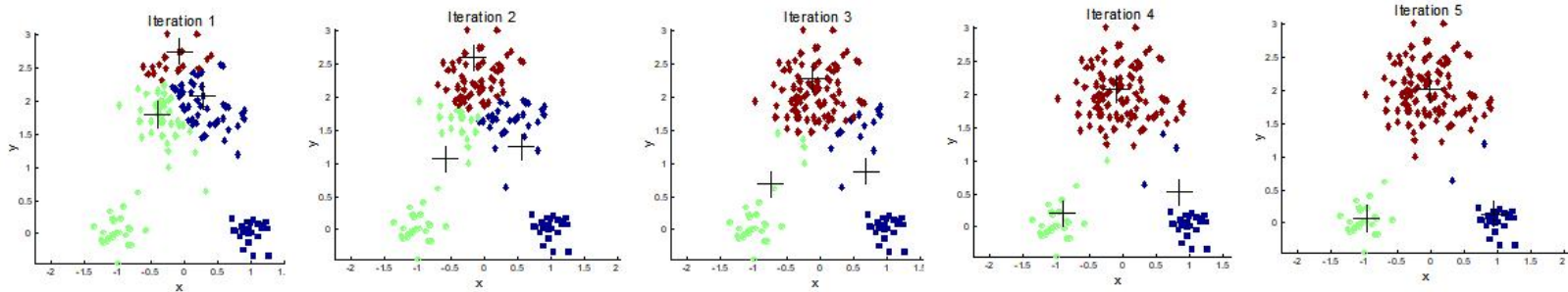
2.2 算法

需要初始化聚类质心

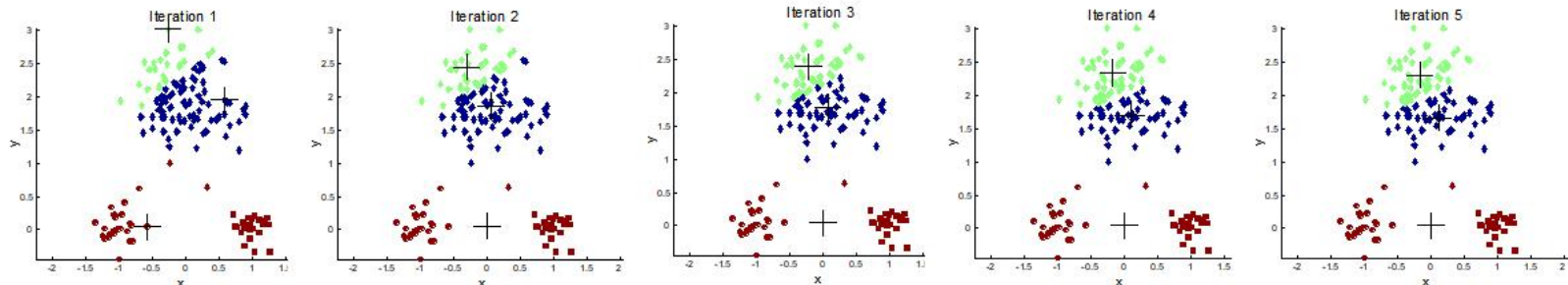


- 不是全局最优（对初始化敏感）
- 聚类的边界是线性的
- K需要事先指定
- 计算开销较大

好的初始化



不好的初始化



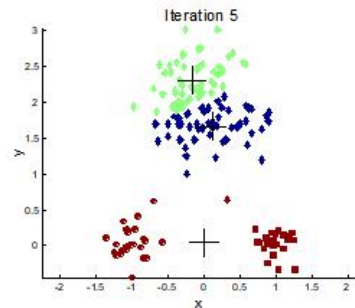
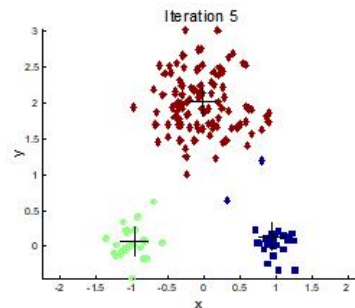
2、K均值聚类

2.2算法



如何解决对初始化敏感的问题？

- 1) 多次运行, 选择SSE最小的一次结果
- 2) 通过采样或者分层次聚类得到更稳定的初始化
- 3) 后处理
 - 去除小的聚类 (有可能是离群点Outlier)
 - 将松散 (较高的SSE值) 的聚类分解
 - 将两个邻近的小聚类融合



2、K均值聚类

2.3举例和应用



将1、2、3、4、5、6、7、8、9九个整数，用K均值分为3类。

第0步：初始化三个中心点： $E(c_1^0) = 1$, $E(c_2^0) = 2$, $E(c_3^0) = 3$

第1a步：分别计算各样本距离中心点的距离，将样本进行分类；

$$c_1^0 = \{1\} \quad c_2^0 = \{2\} \quad c_3^0 = \{3,4,5,6,7,8,9\}$$

第1b步：计算分别分类的均值：

$$E(c_1^1) = 1 \quad E(c_2^1) = 2 \quad E(c_3^1) = 6$$

2、K均值聚类

2.3举例和应用



将1、2、3、4、5、6、7、8、9九个整数，用K均值分为3类。

第2a步：分别计算各样本距离中心点的距离，将样本进行分类；

$$E(c_1^1) = 1 \quad E(c_2^1) = 2 \quad E(c_3^1) = 6$$

$$c_1^1 = \{1\} \quad c_2^1 = \{2, 3, 4\} \quad c_3^1 = \{5, 6, 7, 8, 9\}$$

第2b步：计算分别分类的均值：

$$E(c_1^2) = 1 \quad E(c_2^2) = 3 \quad E(c_3^2) = 7$$

2、K均值聚类

2.3举例和应用



将1、2、3、4、5、6、7、8、9九个整数，用K均值分为3类。

第3a步：分别计算各样本距离中心点的距离，将样本进行分类；

$$E(c_1^2) = 1 \quad E(c_2^2) = 3 \quad E(c_3^2) = 7$$

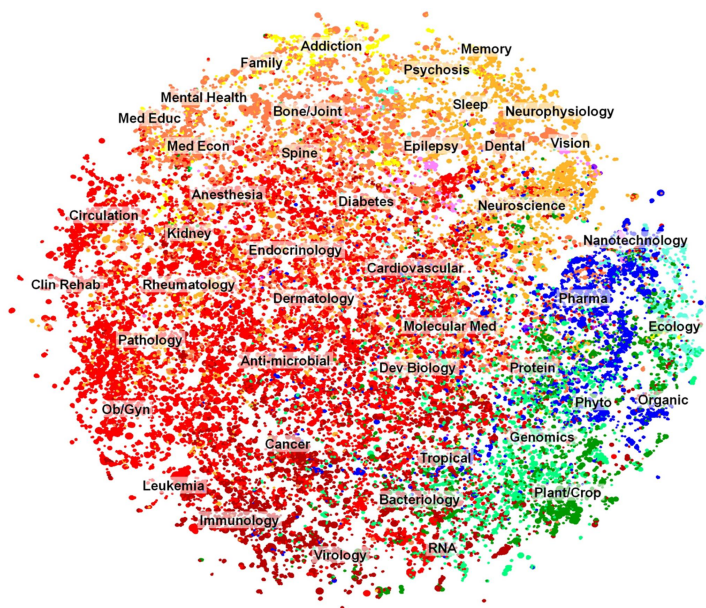
$$c_1^2 = \{1, 2\} \quad c_2^2 = \{3, 4, 5\} \quad c_3^2 = \{6, 7, 8, 9\}$$

第3b步：计算分别分类的均值：

$$E(c_1^3) = 2 \quad E(c_2^3) = 4 \quad E(c_3^3) = 8$$

2、K均值聚类

2.3应用



文本分类：将200多万篇论文聚类到29,000个类别，包括化学、工程、生物、传染疾病、生物信息、脑科学、社会科学、计算机科学等及给出了每个类别中的代表单词

提纲



- 1、无监督学习的概念
- 2、K均值聚类
- 3、主成份分析

3、主成份分析

3.1 概念



- 主成份分析是一种特征降维方法。人类在认知过程中会主动“化繁为简”
- Principle Component Analysis (PCA)
- 奥卡姆剃刀定律 (Occam's Razor) : “如无必要, 勿增实体”, 即“简单有效原理”



3、主成份分析

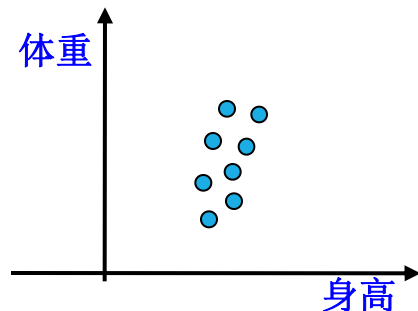
3.1 概念



- 主成份分析主要用于发现数据中的基本结构(数据之间的关联)，是数据分析的主要工具；
- 主成份分析，可把数据用主要成份来表示，这可理解为对数据降维。
- 主成份分析属于多元统计分析的经典方法，最早由Pearson于1901年提出，当时只针对非随机变量，1933年由Hotelling推广到多维随机变量；

- 1) 有一群女神参加选美，身高、体重二维数据如右图；
- 2) 降维：现在需用一个指标(数据)来对女神排序；
选身高还是体重作为排序的依据？

选方差大的数据维度
因为对应的信息量大



3、主成份分析

3.1 概念



- 降维后的结果要保持原始数据固有结构
 - 原始数据中的结构
 - 图像数据中结构：视觉对象区域构成的空间分布
 - 文本数据中结构：单词之间的(共现)相似或不相似



200万像素点



60万像素点

3、主成份分析

3.2数据基础：均值、方差、协方差



- **均值**：假设有n个数据，记为 $X = \{x_i, i=1 \dots n\}$ ，均值 u 为：

$$u = E(X) = \frac{1}{n} \sum_{i=1}^n x_i$$

- **方差**：描述了样本数据的波动程度，等于各个数据与样本均值之差的平方和之平均数：

$$\text{var}(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - u)^2$$

- **协方差**：衡量两个变量之间的相关程度。假设样本由两个变量 X 、 Y 描述，观察到了n组值，记为 $(X, Y) = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

$$\text{cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - E(X))(y_i - E(Y))$$

$\text{cov}(X, Y) > 0$ ，正相关； $\text{cov}(X, Y) < 0$ 负相关； $\text{cov}(X, Y) = 0$ ，不相关

3、主成份分析

3.2数据基础：均值、方差、协方差



编号	x_i	y_i	$x_i - E(X)$	$y_i - E(Y)$	$[x_i - E(X)][y_i - E(Y)]$
1	1	7	-8.33	-16.67	138.89
2	3	11	-6.33	-12.67	80.22
3	6	17	-3.33	-6.67	22.22
4	10	25	0.67	1.33	0.89
5	15	35	5.67	11.33	64.22
6	21	47	11.67	23.33	272.22
	$E(X) = 9.33$	$E(Y) = 23.67$	$Var(X) = 48.22$	$Var(Y) = 192.89$	$E([x_i - E(X)][y_i - E(Y)]) = 96.44$

$$u = E(X) = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{var}(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - u)^2 \quad \text{cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - E(X))(y_i - E(Y))$$

3、主成份分析

3.2数据基础：协方差矩阵



- m维随机变量：** 假设m维变量，记为 $\mathbf{x}=[x_1, x_2, \dots, x_m]^T$ ，观测到了n组值， X 可写为矩阵形式：

$$X = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_n] = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}$$

- 协方差矩阵：** 定义如下：

$$\Sigma = \text{cov}(X, X) = \begin{bmatrix} \text{cov}(\mathbf{x}_1, \mathbf{x}_1) & \text{cov}(\mathbf{x}_1, \mathbf{x}_2) & \dots & \text{cov}(\mathbf{x}_1, \mathbf{x}_m) \\ \text{cov}(\mathbf{x}_2, \mathbf{x}_1) & \text{cov}(\mathbf{x}_2, \mathbf{x}_2) & \dots & \text{cov}(\mathbf{x}_2, \mathbf{x}_m) \\ \dots & \dots & \dots & \dots \\ \text{cov}(\mathbf{x}_m, \mathbf{x}_1) & \text{cov}(\mathbf{x}_m, \mathbf{x}_2) & \dots & \text{cov}(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix}$$

3、主成份分析

3.3主成份定义



- **线性变换：** 考虑 m 维变量 $\mathbf{x} = [x_1, x_2, \dots, x_m]^T$, 到 m 维随机变量 $\mathbf{y} = [y_1, y_2, \dots, y_m]^T$ 的线性变换

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_m \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \dots \\ a_m \end{bmatrix} \cdot \mathbf{x} = \mathbf{A}\mathbf{x}$$

- 如上式的线性变换，**主成份**定义如下：
 - 1)未知变换的系数向量 \mathbf{a}_i 为**单位向量**， $\|\mathbf{a}_i\|=1 \quad i=1,2,\dots,m$
 - 2) y_i 与 y_j 互不相关，即 $\text{cov}(y_i, y_j)=0 \quad i \neq j$
 - 3) y_1 是所有 \mathbf{x} 的线性变换中**方差最大**的， y_2 是 \mathbf{x} 的所有线性变换中与 y_1 不相关中方差最大的.....；称 y_1, y_2, \dots 为 \mathbf{x} 的**第一主成份、第二主成份.....**；

3、主成份分析

3.4主成份分析



- 主成份分析：** 实际上求解一组标准正交基， a_1, a_2, \dots, a_m ，使变换后的随机向量 y_1, y_2, \dots, y_m 的方差取极值，即：

$$\max_{a_1} \text{var}(y_1) = \max_{a_1} \text{var}(a_1 x) \quad s.t. \quad a_1 a_1^T = 1$$

推导如下： $\text{var}(y_1) = \text{var}(a_1 x) = \frac{1}{n-1} \sum_{i=1}^n (a_1 x^{(i)} - E(a_1 x^{(i)}))^T (a_1 x^{(i)} - E(a_1 x^{(i)}))^T$

$$\text{var}(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - u)^2$$

$$= \frac{1}{n-1} a_1 \sum_{i=1}^n (x^{(i)} - E(x^{(i)}))^T (x^{(i)} - E(x^{(i)}))^T a_1^T$$

$$\text{cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - E(X))(y_i - E(Y))$$

$$= a_1 \text{cov}(X, X) a_1^T = a_1 \Sigma a_1^T$$

原问题等价为： $\max_{a_1} \text{var}(y_1) = \max_{a_1} a_1 \Sigma a_1^T \quad s.t. \quad a_1 a_1^T = 1$

3、主成份分析

3.4主成份分析



- 等式约束的拉格朗日极值问题：

$$\max_{a_1} \text{var}(y_1) = \max_{a_1} a_1 \Sigma a_1^T \quad s.t. \quad a_1 a_1^T = 1$$

- 用拉格朗日乘子法，构建拉格朗日函数：

$$L(a_1, \lambda) = a_1 \Sigma a_1^T - \lambda(a_1 a_1^T - 1)$$

对拉格朗日函数取偏导，并令其为0，

$$\frac{\partial L(a_1, \lambda)}{\partial a_1} = \Sigma a_1^T - \lambda a_1^T = 0$$

$$\Rightarrow \Sigma a_1^T = \lambda a_1^T$$

$$\Rightarrow \text{var}(y_1) = a_1 \Sigma a_1^T = \lambda a_1 a_1^T = \lambda$$

$$\Sigma = \text{cov}(X, X)$$

a_1 为协方差矩阵的单位特征向量；
 λ 为特征值。

3、主成份分析

3.4主成份分析



- 主成份的性质：

设 \mathbf{x} 是 m 维随机变量， Σ 是 \mathbf{x} 的协方差矩阵， Σ 的特征值分别为 $\lambda_1 \geq \lambda_2 \geq \dots \lambda_k$ 对应的特征向量分别为 a_1, a_2, \dots, a_k , 那么 \mathbf{x} 对应的第 k 主成份是：

$$y_k^{(j)} = a_k x^{(j)} \quad j = 1, 2, \dots, n$$

那么 \mathbf{x} 对应的第 k 主成份 y_k 的协方差是协方差矩阵的第 k 个特征值：

$$\text{var}(y_k) = \lambda_k$$

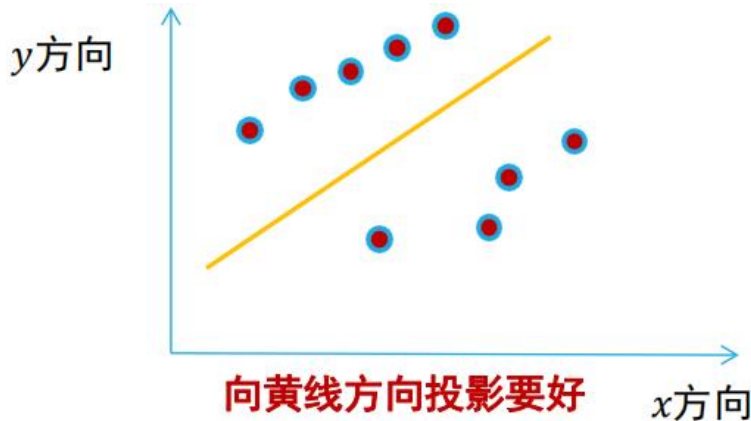
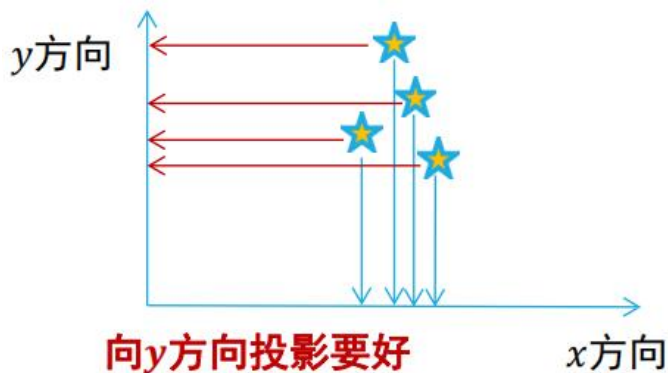
3、主成份分析

3.4主成份分析：算法动机

保证样本
投影后方差最大



- 在数理统计中，方差被经常用来度量数据和其数学期望（即均值）之间偏离程度，这个偏离程度反映了数据分布结构。
- 在许多实际问题中，研究数据和其均值之间的偏离程度有着很重要的意义。
- 在降维之中，需要尽可能将数据向方差最大方向进行投影，使得数据所蕴含信息没有丢失，彰显个性。
- 如左下图所示，向 y 方向投影（使得二维数据映射为一维）就好比向 y 方向投影结果在降维这个意义上而言要好；右下图则是黄线方向投影要好。



3、主成份分析

3.4主成份分析：算法动机



- 主成分分析思想是将 m 维特征数据映射到 k 维空间（ $k \gg m$ ），去除原始数据之间的冗余性（通过去除相关性手段达到这一目的）。
- 将原始数据向这些数据方差最大的方向进行投影。一旦发现了方差最大的投影方向，则继续寻找保持方差第二的方向且进行投影。
- 将每个数据从 m 维高维空间映射到 k 维低维空间，每个数据所得到最好的特征就是使得每一维上样本方差都尽可能大。
- 主成份维数 k 的选择, 可由线性变换后期望保留的信息量比例来决定：

$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_m} \geq \sigma \quad (\sigma = 0.95, \dots)$$

- 不足及改进：当数据量比较大时，协方差矩阵计算费时，特征值和特征向量求解困难。可采用矩阵的**奇异值分解**来求解特征值和特征向量(参考《矩阵分析》)

3、主成份分析

3.4主成份分析：算法动机



- 假设有 n 个 d 维样本数据所构成的集合 $D = \{x_1, x_2, \dots, x_n\}$ ，其中 $x_i (1 \leq i \leq n) \in R^d$ 。
- 集合 D 可以表示成一个 $n \times d$ 的矩阵 \mathbf{X} 。
- 假定每一维度的特征均值均为零（已经标准化）。
- 主成分分析的目的在于求取一个且使用一个 $d \times l$ 的映射矩阵 \mathbf{W} 。
- 给定一个样本 x_i ，可将 x_i 从 d 维空间如下映射到 l 维空间： $(x_i)_{1 \times d}(\mathbf{W})_{d \times l}$

? 如何求取
映射矩阵 \mathbf{W}

- 将所有降维后数据用 \mathbf{Y} 表示，有 $\mathbf{Y} = \mathbf{X} \mathbf{W}$

降维 原始 映射
结果 数据 矩阵

- $\mathbf{Y} = n \times l$
- $\mathbf{X} = n \times d$
- $\mathbf{W} = d \times l$

3、主成份分析

3.4主成份分析：算法动机



- 输入： n 个 d 维样本数据所构成的矩阵 \mathbf{X} ，降维后的维数 l
- 输出：映射矩阵 $\mathbf{W} = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_l\}$
- 算法步骤：

1: 对于每个样本数据 \mathbf{x}_i 进行中心化处理： $\mathbf{x}_i = \mathbf{x}_i - \mu, \mu = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$

2: 计算原始样本数据的协方差矩阵： $\Sigma = \frac{1}{n-1} \mathbf{X}^T \mathbf{X}$

3: 对协方差矩阵 Σ 进行特征值分解，对所得特征根按其值大到小排序 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_l$

4: 取前 l 个最大特征根所对应特征向量 $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_l$ 组成映射矩阵 \mathbf{W}

5: 将每个样本数据 \mathbf{x}_i 按照如下方法降维： $(\mathbf{x}_i)_{1 \times d} (\mathbf{W})_{d \times l} = 1 \times l$

谢谢