



本课件仅用于教学使用。未经许可，任何单位、组织和个人不得将课件用于该课程教学之外的用途(包括但不限于盈利等)，也不得上传至可公开访问的网络环境

1

数据科学导论

Introduction to Data Science

第四章 数据挖掘基础

陈恩红，黄振亚

Email: huangzhy@ustc.edu.cn, cheneh@ustc.edu.cn

课程主页:

<http://staff.ustc.edu.cn/~huangzhy/Course/DS2024.html>



无监督学习

2

- 数据挖掘任务 —— 无监督学习
 - 无标签数据是现实中最常见的数据
 - 例如，拍摄的照片等



该图片有关联



这张照片是哪里?

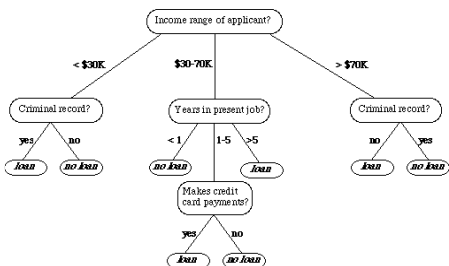




数据挖掘基础

数据挖掘——四个任务有哪些常用方法？

分类与预测



关联分析

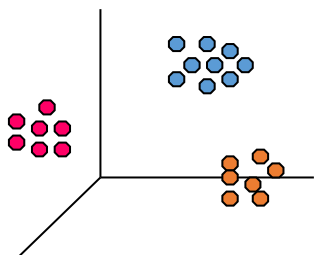


数据

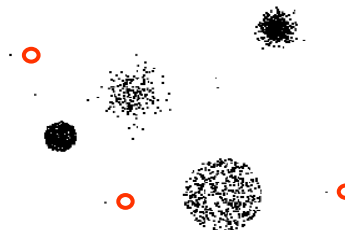
	T		H		P	
	L	H	L	H	L	H
J	-6.0	8.8	60	100	986	1044
F	-2.8	10.9	48	100	973	1025
M	-5.6	17.7	34	100	976	1037
A	-1.2	22.2	27	100	996	1036
M	-0.8	27.8	25	100	1003	1034
J	5.2	29.1	26	100	998	1030
J	9.8	30.6	23	99	997	1027
A	5.6	26.1	31	100	992	1029
S	5.2	24.8	35	100	998	1028
O	-0.4	21.3	42	100	990	1031
N	-7.6	17.3	55	100	963	1023
D	-10.4	9.2	53	100	987	1039

table 17a
2010 monthly weather variation, Cambridge (UK)

聚类



异常检测



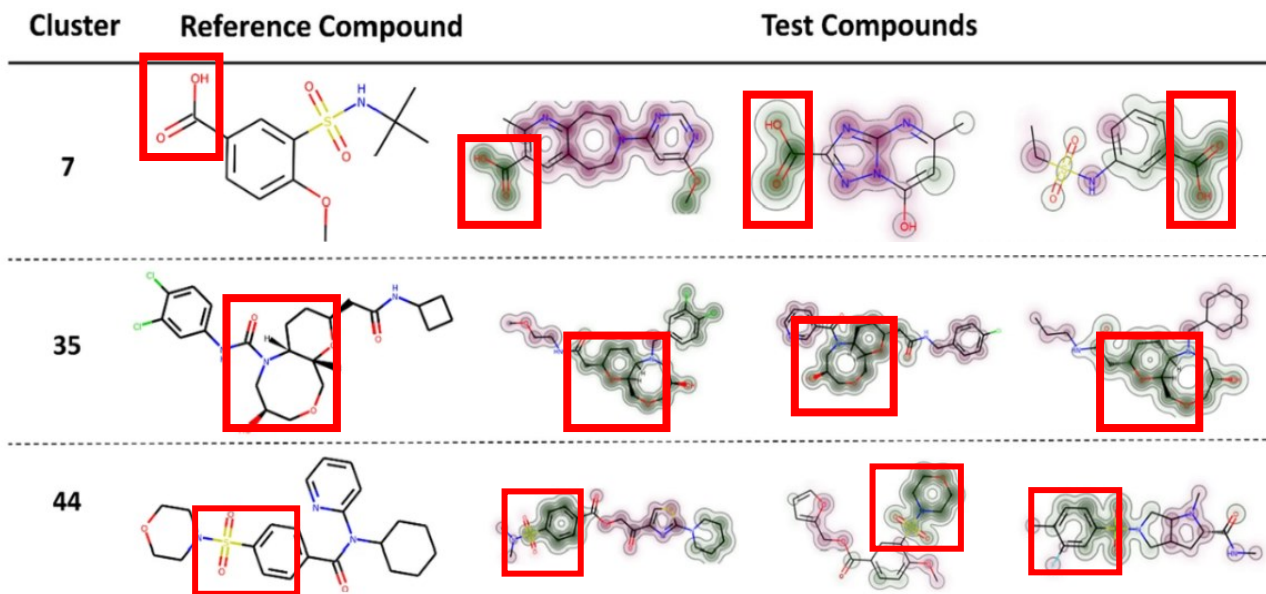


聚类分析: 应用实例

4

案例一：分子与药物分析

- 输入：生物医药分子
- 结构相似度更高的分子被分配到一个聚簇



- 第7簇中含有芳香族羧酸酯
- 第35簇中含有芳基卤化物
- 第44簇中含有磺胺

➤ Hadipour, H., Liu, C., Davis, R., Cardona, S.T., & Hu, P. (2022). Deep clustering of small molecules at large-scale via variational autoencoder embedding and K-means. *BMC Bioinformatics*, 23.

聚类分析: 应用实例

5

案例二：疫情溯源

- ▣ 输入：纽约市病例人群的信息：地理位置等
- ▣ 对纽约市的冠状病毒病(COVID-19)爆发场所聚类，定位的感染源

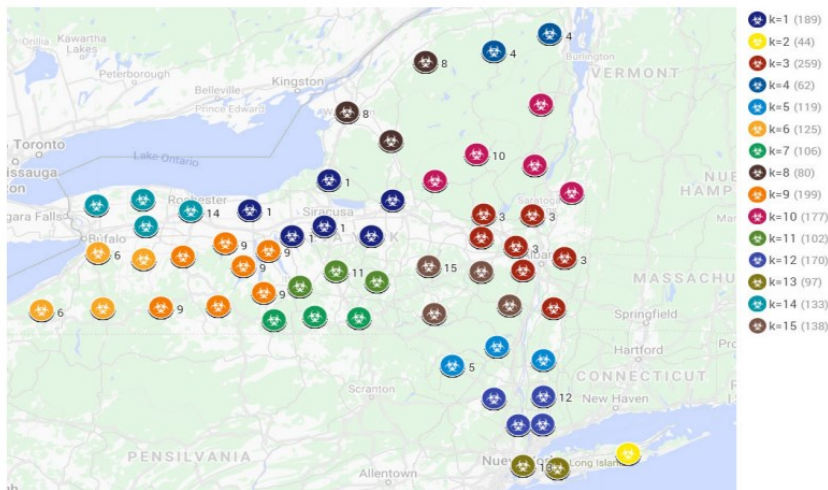


FIGURE 5. K-means clustering ($k = 15$) in New York state.

按病例的位置聚类，得到K个聚簇区域

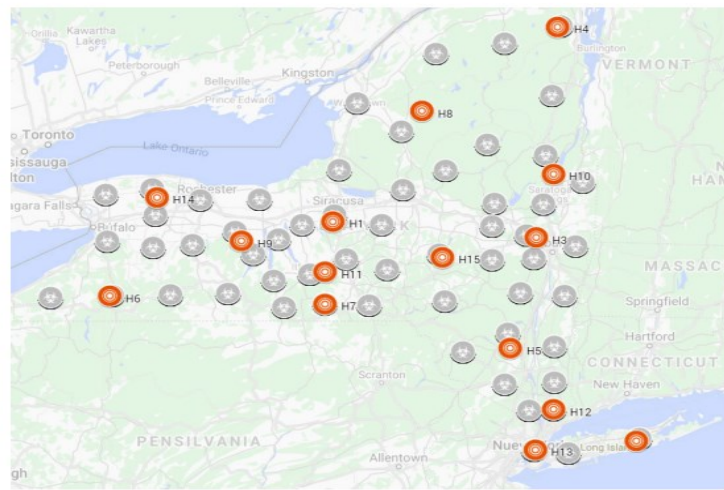


FIGURE 7. Hot spots H_k for each cluster in New York state (orange circles).

聚簇区域中心被视为感染源

- Guevara C, Peñas M S. Surveillance Routing of COVID-19 Infection Spread Using an Intelligent Infectious Diseases Algorithm[J]. Ieee Access, 2020, 8: 201925-201936.

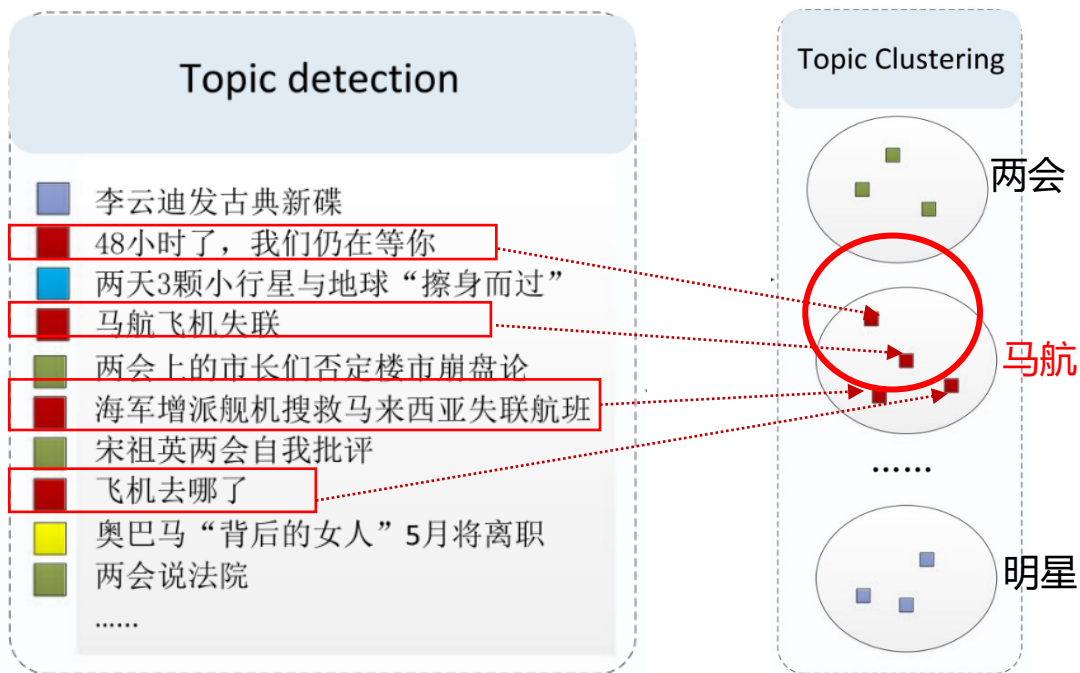


聚类分析: 应用实例

6

案例三：网络舆情分析-话题挖掘

- 输入：社交媒体中的评论与话题
- 话题聚类，同一话题簇中出现的**关键词相似**



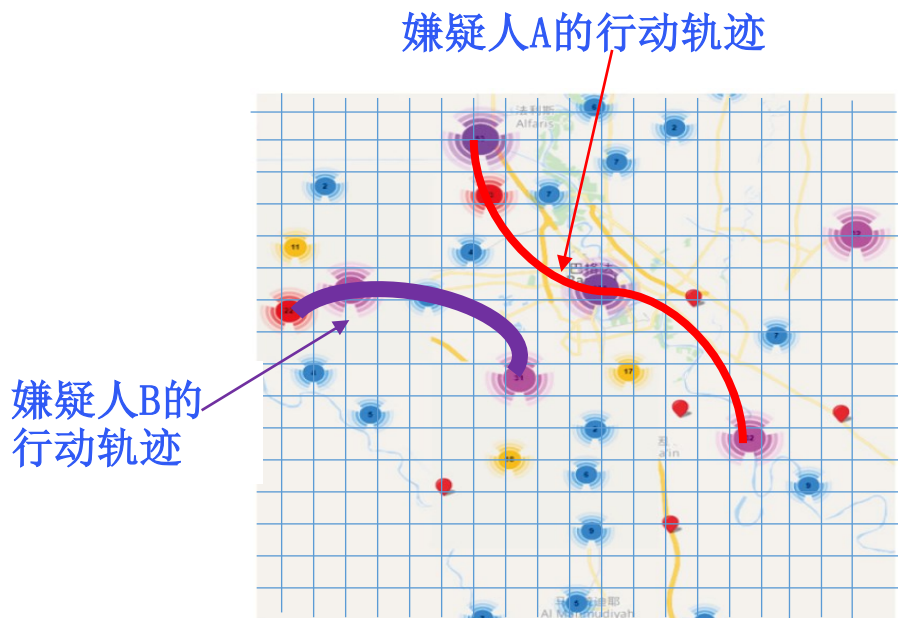
➤ Cai Y, Wu X, Xie X, et al. A topic mining method for multi-source network public opinion based on improved hierarchical clustering[C]//2019 IEEE DSC. IEEE, 2019: 439-444.



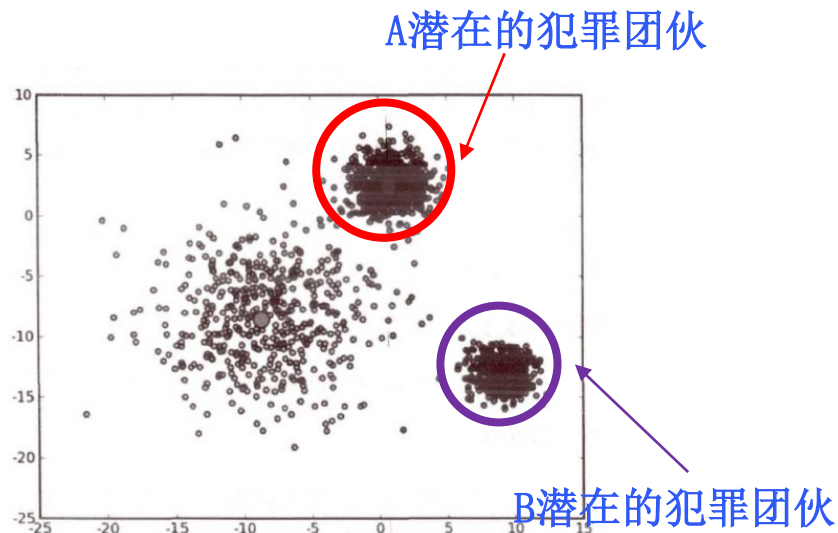
聚类分析: 应用实例

7

- 案例四：安防与维稳-犯罪团伙识别
 - 输入：人员轨迹时空数据：如网吧、酒店、车站等，
 - 对嫌疑人的轨迹信息进行聚类，找出犯罪团伙。



地理空间网格化



轨迹信息聚类结果



聚类分析: 应用实例

案例五：教育问题的聚类

- 输入：数学应用题
- 题目聚类，同类题目的解答模板一样

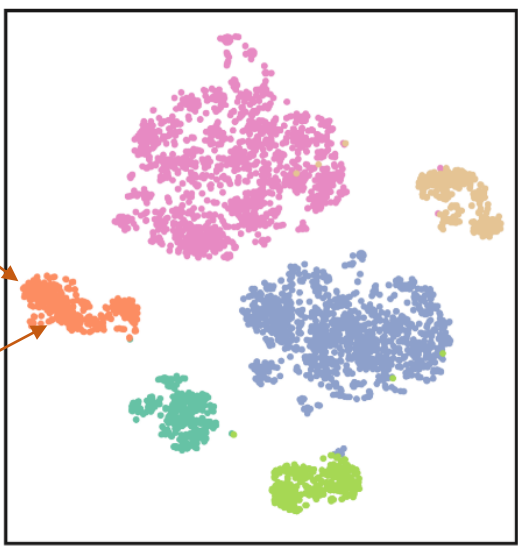
数学应用题

Prob. A: Norma has 88 cards. She loses 70. How many cards will Norma have?

Eq: $88 - 70$

Prob. B: Joyce starts with 75 apples. She gives 52 to Larry. How many apples does Joyce end with?

Eq: $75 - 52$



基础运算的模式不同

- $n_1 + n_2$
- n_1 / n_2
- $n_1 - n_2$
- $(n_1 + n_2) * n_3$
- $n_1 * n_2$
- $(n_1 + n_2) / n_3$

- Li, Z., Zhang, W., Yan, C., Zhou, Q., Li, C., Liu, H., & Cao, Y. (2022). Seeking Patterns, Not just Memorizing Procedures: Contrastive Learning for Solving Math Word Problems. ArXiv, abs/2110.08464.
- Huang, Z., Lin, X., Wang, H., Liu, Q., Chen, E., Ma, J., Su, Y., & Tong, W. (2021). DisenQNet: Disentangled Representation Learning for Educational Questions. Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining.



聚类分析: 应用实例

9

- 案例六：学业数据分析——优化教师教学
 - 输入：试验学校的学生考试数据
 - 聚类发现教师教学模式的规律

根据考试数据对班级进行简单聚类，根据聚类结果，发现70%的类里，两个班级是同位授课教师



聚类

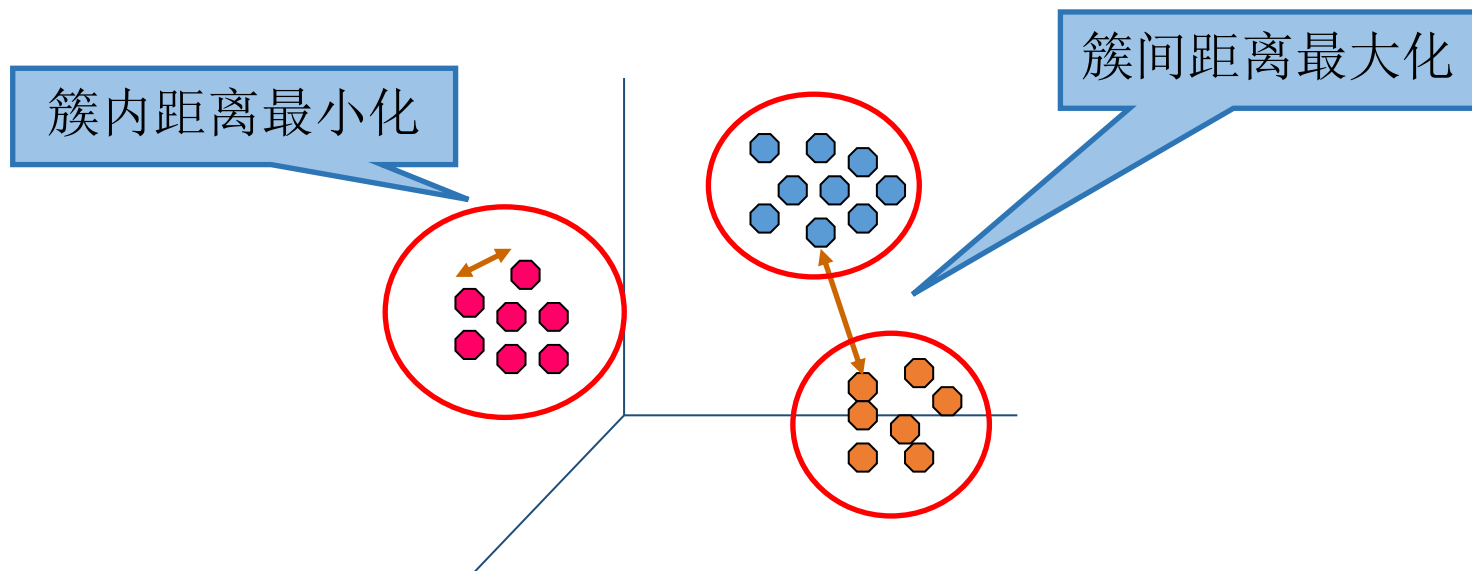




聚类分析

10

- 数据挖掘任务 —— 聚类(Clustering)
 - 目标：对数据进行“**群体性**”分析，将样本分为若干个**簇** (Clusters)
 - 其中，每个簇都由相似的样本所组成
 - 簇的特点：簇内相似（距离近），簇间相异（距离远）





聚类分析

11

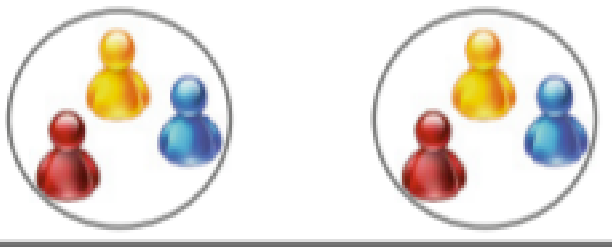
聚类分析要解决三个问题

1. 如何定义簇？：即，思考我们的目标（但具有主观性）

- “群体性”的依据：不同的“群体性”立场，可以得到不同的簇
- 例如，学生分组应该考虑 **技能互补**？ 还是 **能力相近**？

2. 如何定义相关性？ 即，度量数据之间的相似性

- 相似性度量往往存在一定局限性，未必反映聚类的真实意图
- 例如，常用向量表征数据(人的爱好)， 但是否绝对相似？



技能互补？ 能力相近？



图片本身相似，但代表的类别完全不同？



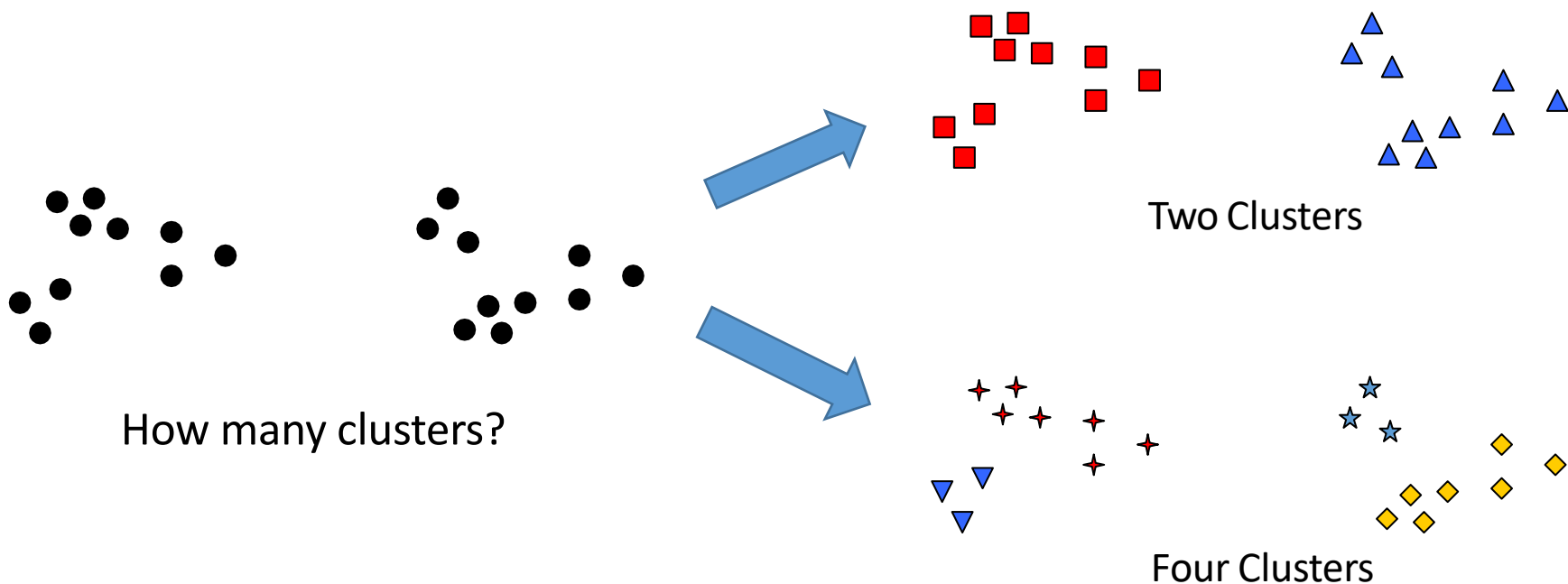
聚类分析

12

聚类分析要解决三个问题

3. 如何决定簇的数量？即，选择合适数量的簇

- 数据没有天然标签，簇的数量往往是个开放性问题
- 避免过大或过小的簇，会导致失去代表性，但这未必可通过簇数调节



哪一种方式更合适？



聚类分析

13

- 聚类方法：最常见的无监督学习算法
- 常用方法
 - K均值聚类(K-means)
 - 层次聚类(Hierarchical Clustering)
 - 密度聚类(Density-based Clustering)
 - 聚类效果验证
 - 前沿聚类方法



聚类分析：K-means

14

□ K-means的基本概念

- **数据**：视作高维空间中的一个点，表示为向量
- **中心点**：簇的中心，反应簇的共有属性
- **簇的数量**：人为设定
- **数据的关系度量**：用“平均”的方式簇中心与簇中数据





聚类分析: K-means

15

- K-means算法: 设定K个中心, 形成K个数据簇
 - 根据问题目标, 预先指定簇的个数K
 - 每一个簇存在一个中心点
 - 每一个数据属于最近中心点对应的簇
 - 簇中心的更新: 依赖簇中数据的算术平均
 - 簇中心更新后: 根据度量, 数据将重新分配至不同的簇
 - 算法是迭代过程: 簇中心与簇数据迭代更新, 直至稳定

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



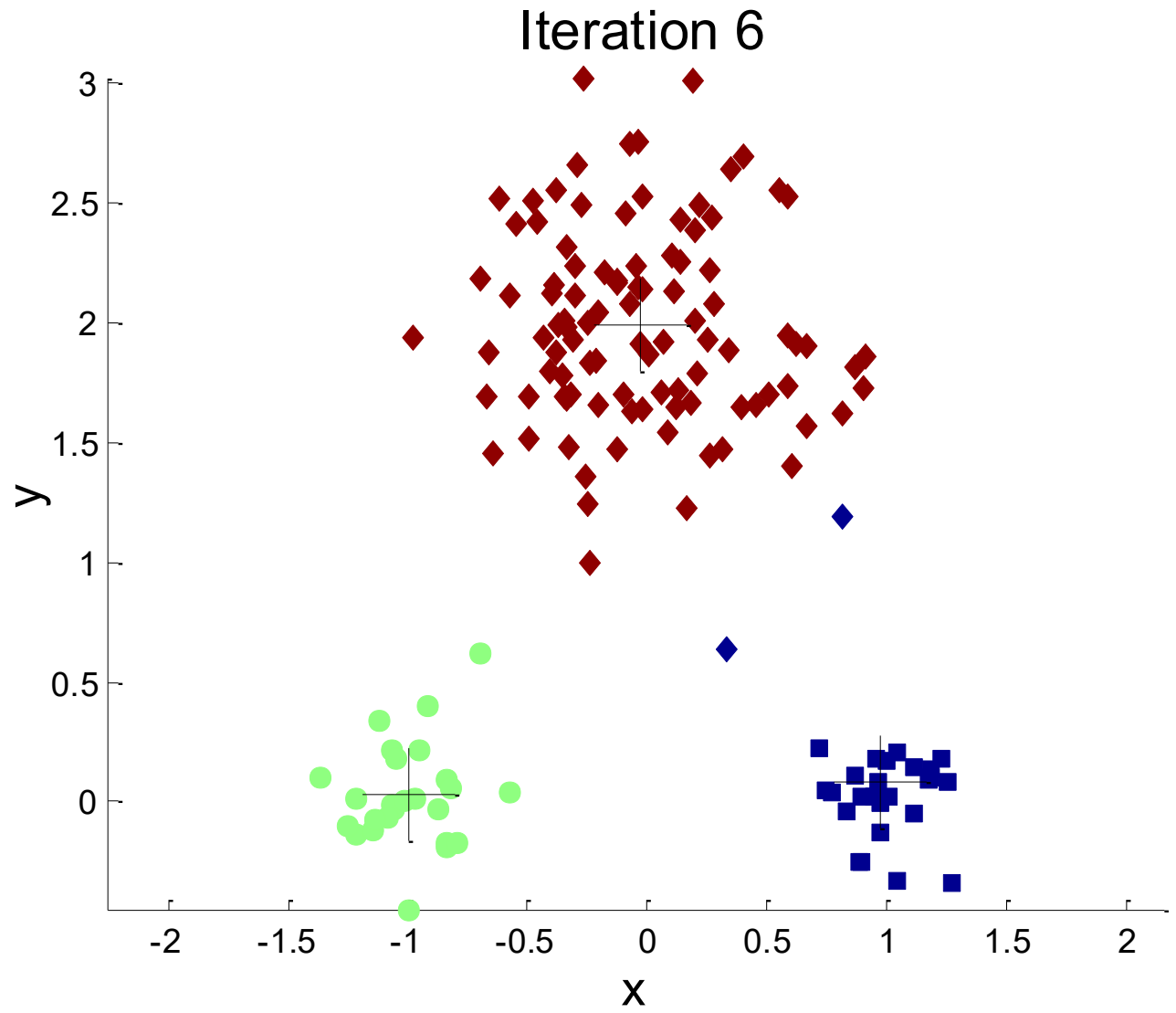
聚类分析：K-means

16

- 一个例子展示：K-means过程
- 3个簇中心
- 数据：二维空间的点
- “+”：表示簇中心
- 颜色：不同的簇



聚类分析: K-means

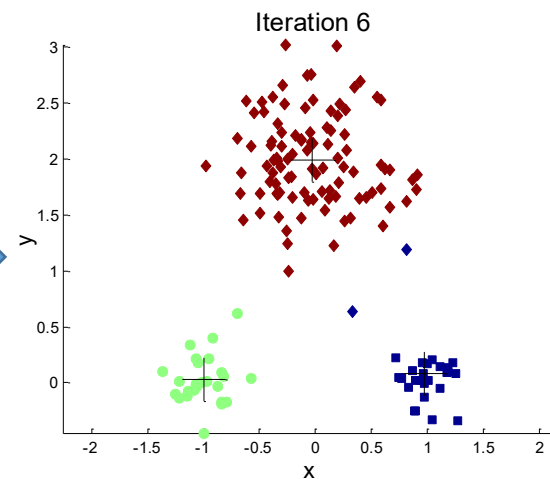
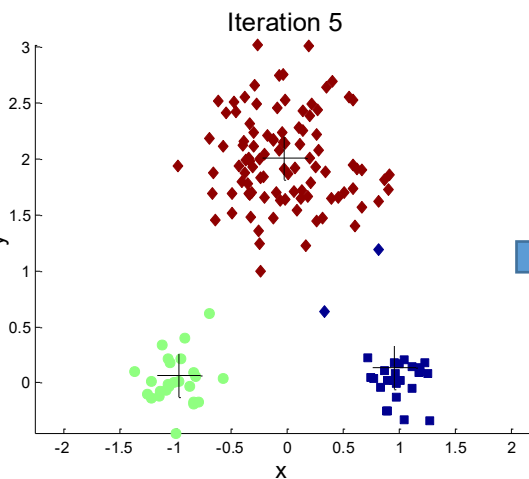
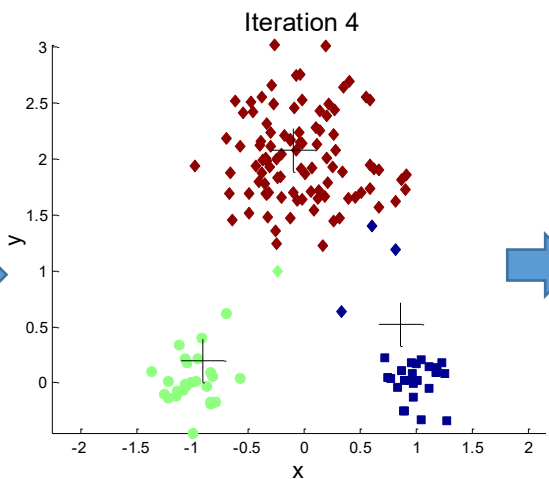
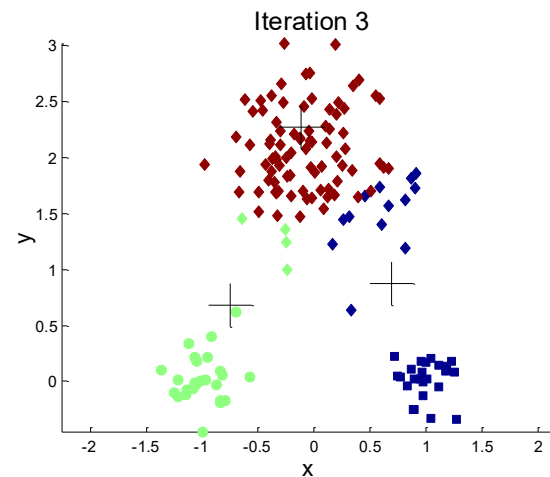
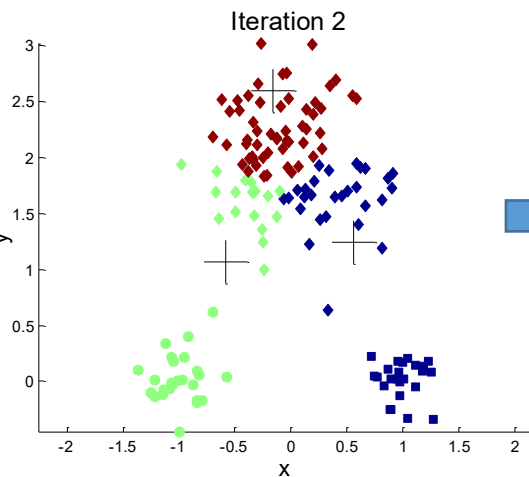
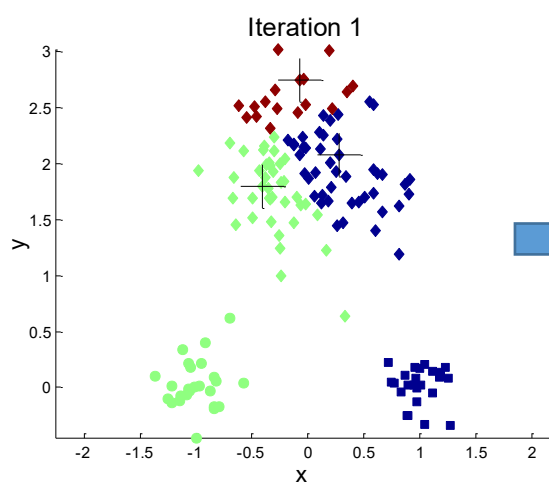




聚类分析: K-means

18

□ K-means示例





聚类分析：K-means

19

□ 如何评估K-means的效果

- 指标：平方误差和 (Sum of Squared Error , SSE)
- 算法目标：优化数据与簇中心的距离
 - 定义每个数据的聚类误差：样本数据与最近簇的距离

□ SSE定义为

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

- x 是簇 C_i 的样本, m_i 是簇 C_i 的质心, 可证明 m_i 是簇 C_i 平均 (mean)
- 注：面对多个聚类结果时, 倾向于SSE更小的方式
 - 当簇数量K增加时, SSE一般趋于下降, 因此尽量在相同K下比较SSE
 - K和SSE较小的聚类, 优于 K和SSE交大的聚类



聚类分析：K-means

20

□ K-means的特点

□ 关于中心点

- 中心点一般采用随机初始化，因此重复K-means得到的结果可能不同
- 中心点一般设置为簇内数据的平均向量（算术平均）

□ 关于相关性度量

- 可用欧几里得距离、余弦相似度、相关系数等度量

□ 关于算法运行

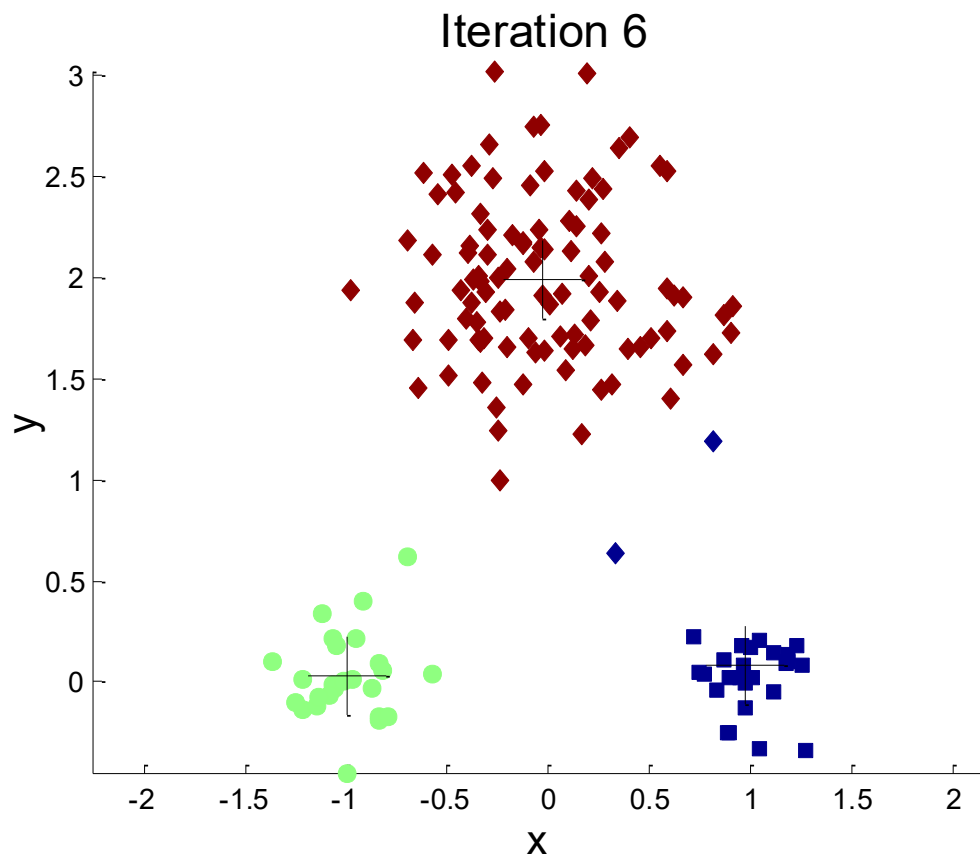
- K均值算法常常几轮就收敛
- 算法停止条件为：Until relatively few points change clusters
 - 低于一定数量的数据 更新簇的归属，即，每个簇中 只有 很少的数据 发生变化
- 算法复杂度 $O(n \times K \times I \times d)$
 - n = 样本总数, K = 簇的个数数, I = 迭代轮数, d = 特征维度



聚类分析: K-means

21

- K-means的特点: 初始中心如何选择?
 - 中心点初始化较好时 (回顾19页的例子)

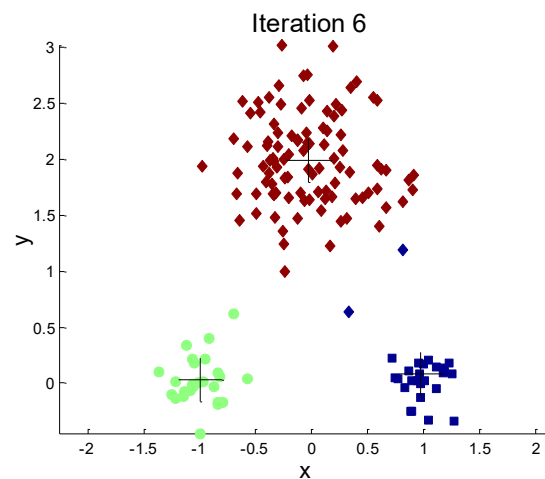
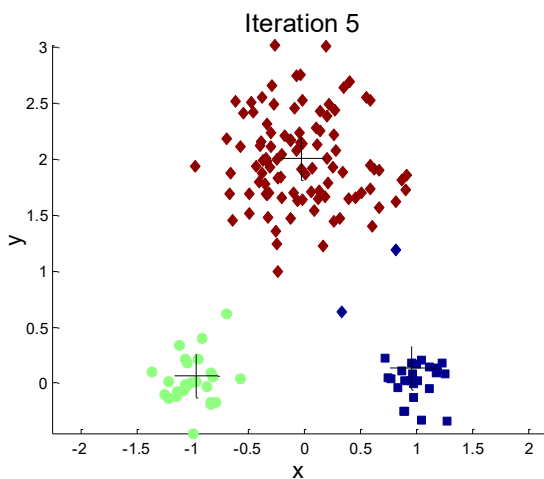
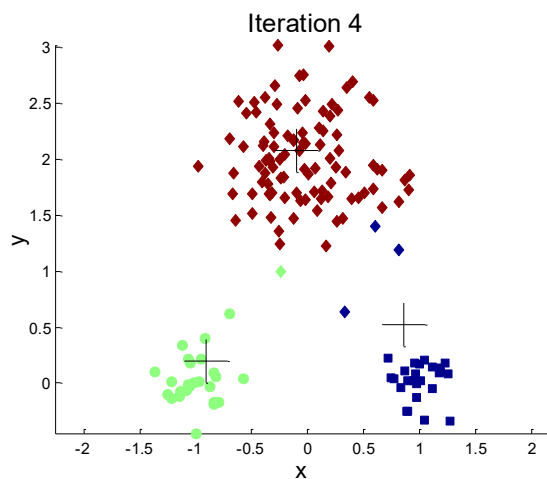
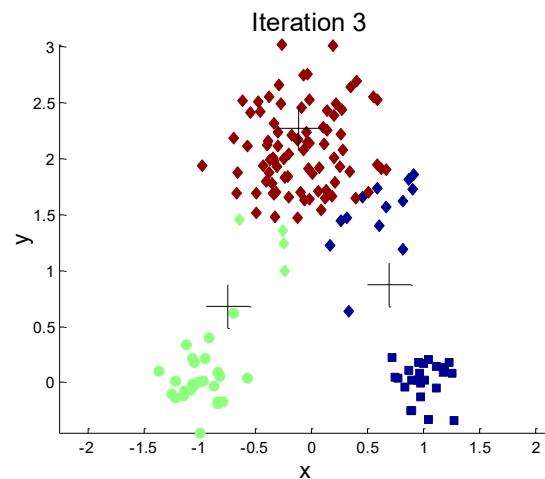
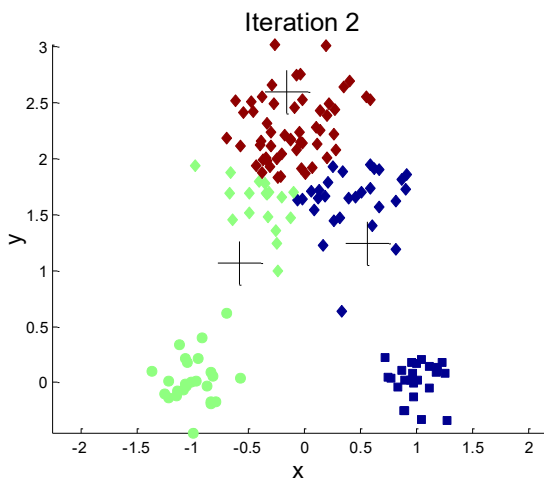
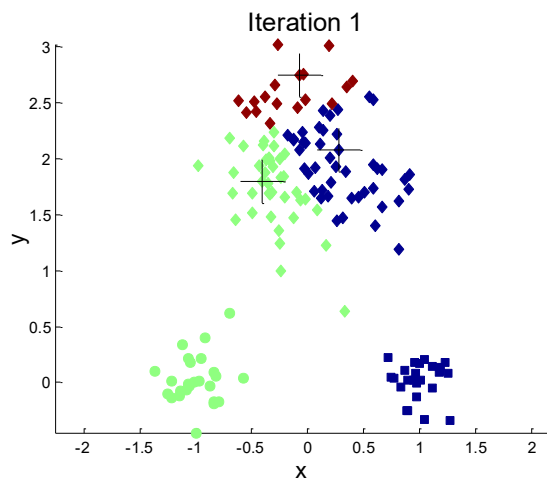


聚类结果好



数据挖掘基础

□ K-means的特点: 初始中心如何选择?



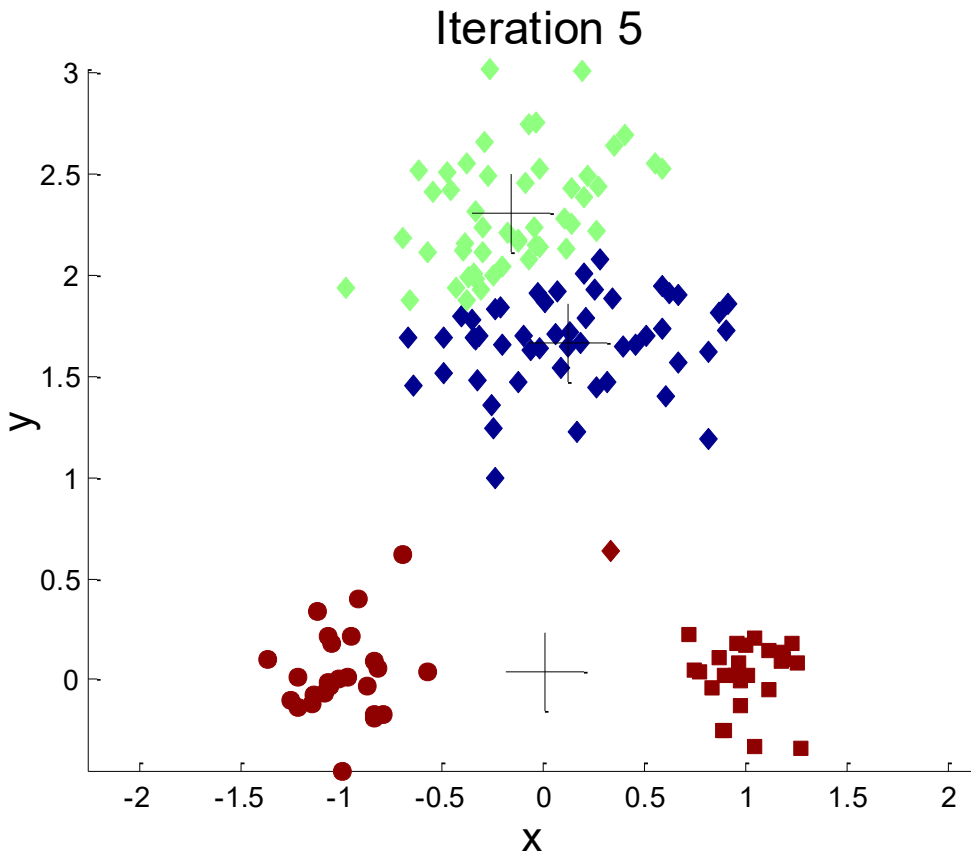
聚类结果好



聚类分析: K-means

23

- K-means的特点: 初始中心如何选择?
 - 中心点初始化较差时



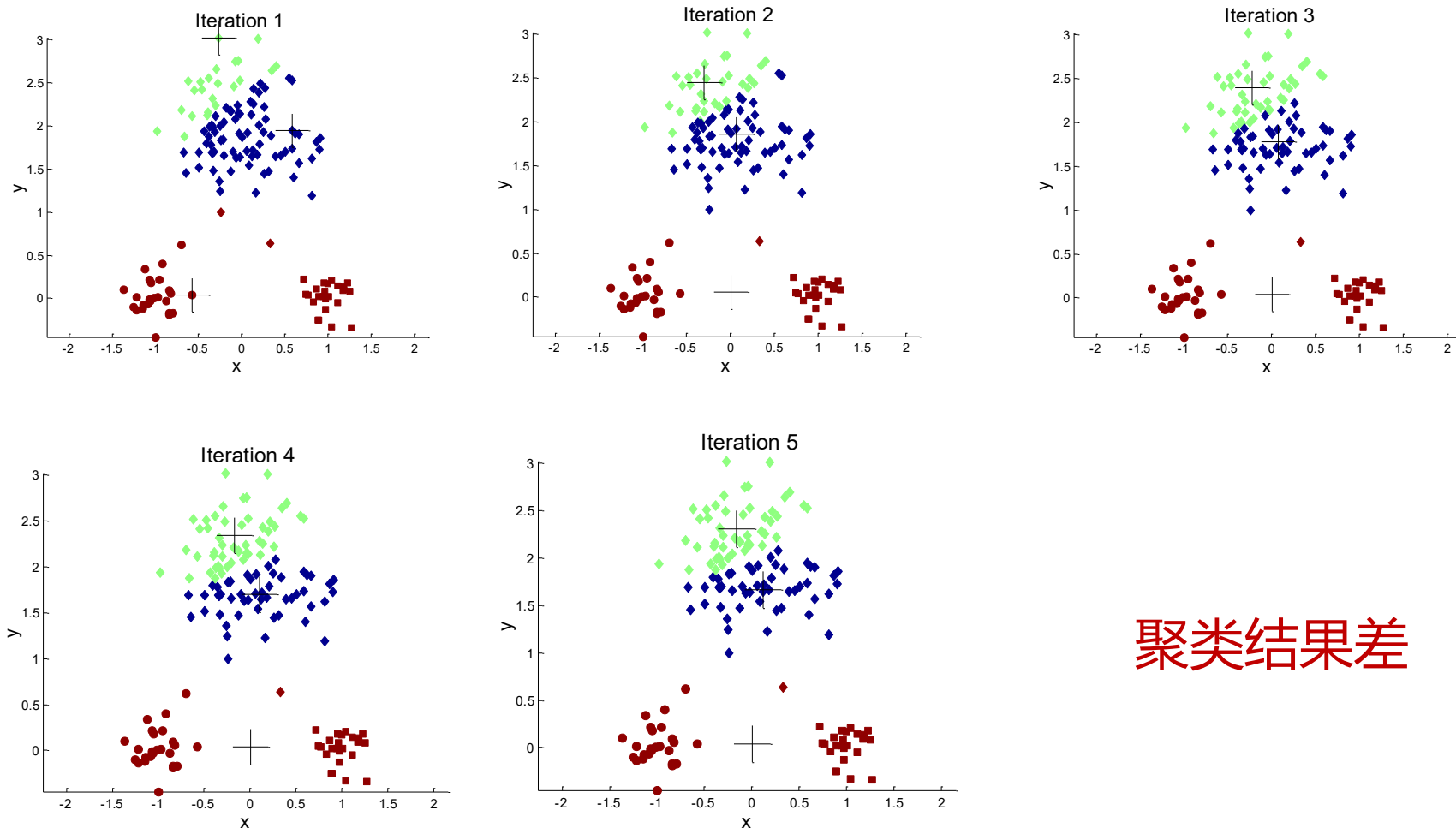
聚类结果差



聚类分析: K-means

24

□ K-means的特点: 初始中心如何选择?



聚类结果差