



基于空间聚类的 出租车轨迹数据分析

汇报人：钟代琪

2022.10.12



Part 1 上下客热点分析

- K-means
- DBSCAN

Part 2 异常轨迹分析

- 层次聚类



上下客热点分析



K-means核心思想：通过**迭代**把数据对象划分到不同的簇中，以求**目标函数最小化**，从而使生成的簇尽可能地紧凑和独立。

步骤：

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

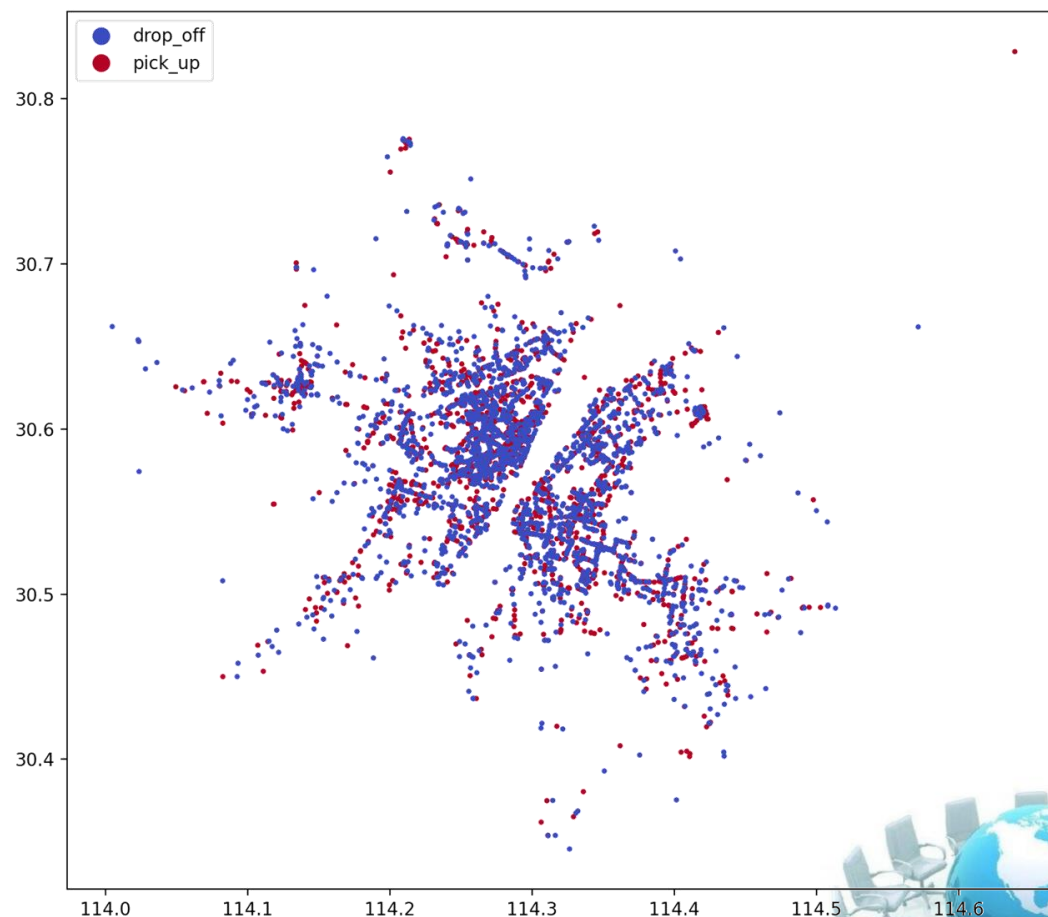
Step1. **随机**选取K个对象作为初始的**K个簇的质心**；

Step2. 将其余对象根据其与各个簇质心的距离**分配到最近的簇**；

再求新形成的簇的质心；

Step3. 不断重复**迭代**重定位过程，直到满足终止条件为止。

武汉市2018年11月5日出租车上下客记录



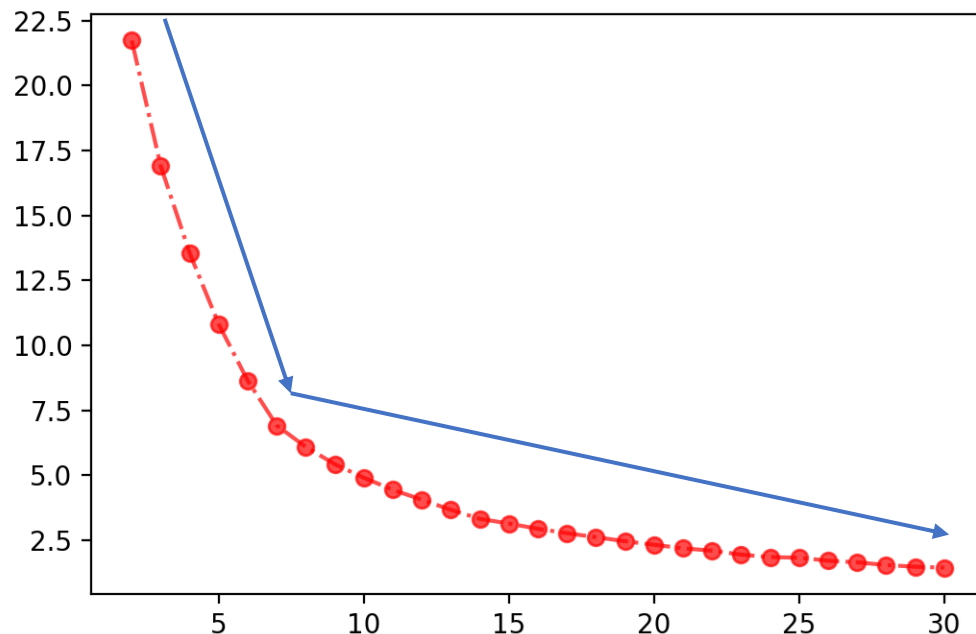
聚类数确定 — 手肘法



手肘法思想：

k-means以最小化样本与质心平方误差作为目标函数，将每个簇的质点与簇内样本点的平方距离误差和称为**畸变程度**。

畸变程度会**随着K的增加而降低**，在达到某个**临界点**时畸变程度会得到极大改善，之后缓慢下降，其函数形状类似一个手肘。该临界点可以考虑为**聚类性能较好**的参考点。



聚类数确定 — 轮廓系数



测绘与地理信息学院
COLLEGE OF SURVEYING AND GEO-INFORMATICS

轮廓系数 $S = \sum_i S(i)$

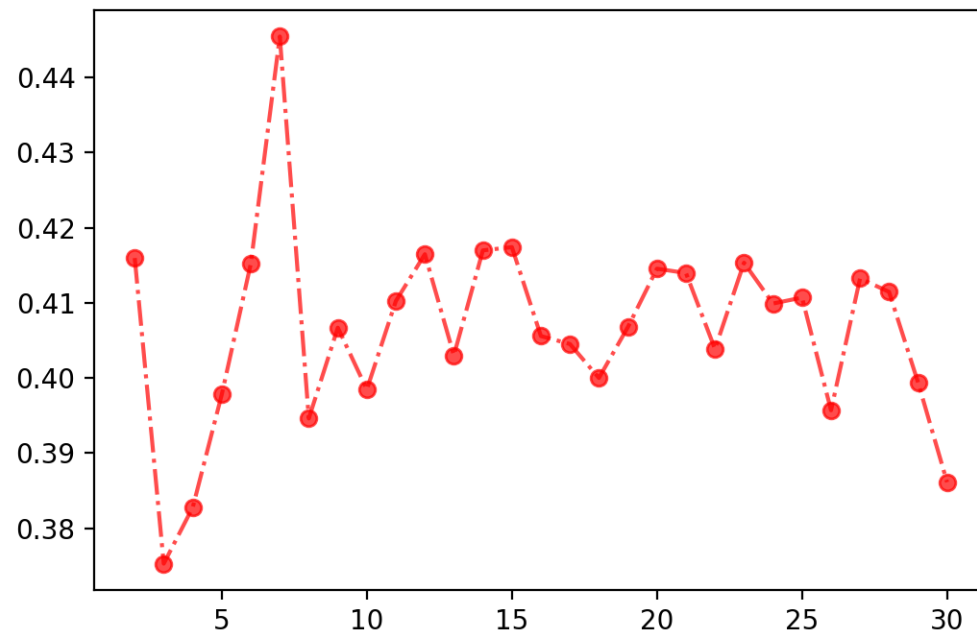
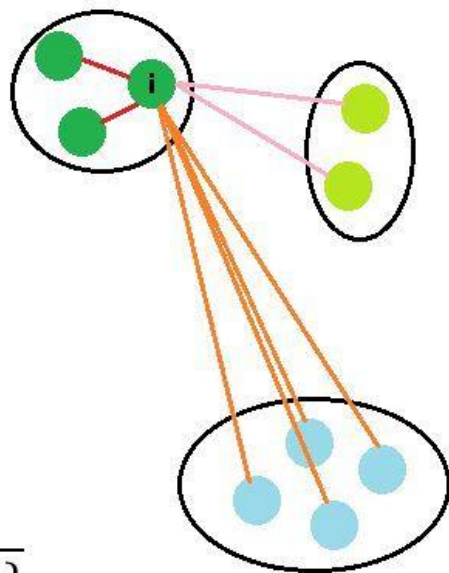
内聚度

$a(i) = \text{avg} \{ \text{red lines} \}$

分离度

$b(i) = \min \left\{ \begin{array}{l} \text{avg} \{ \text{pink lines} \} \\ \text{avg} \{ \text{orange lines} \} \end{array} \right.$

$$S(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$



$a(i) = \text{avg}\{i \text{ 到所属簇内所有其他点的距离} \}$

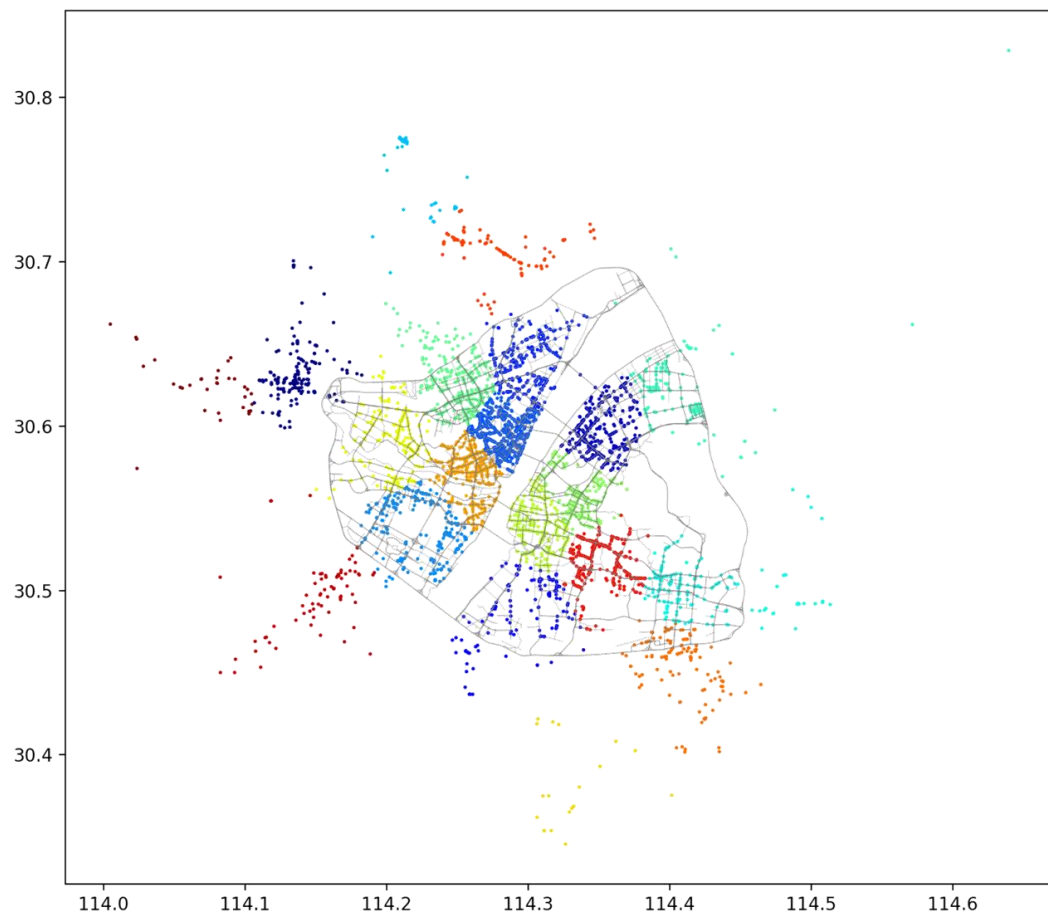
$b(i) = \min\{i \text{ 到其他某一簇内所有点的平均距离} \}$



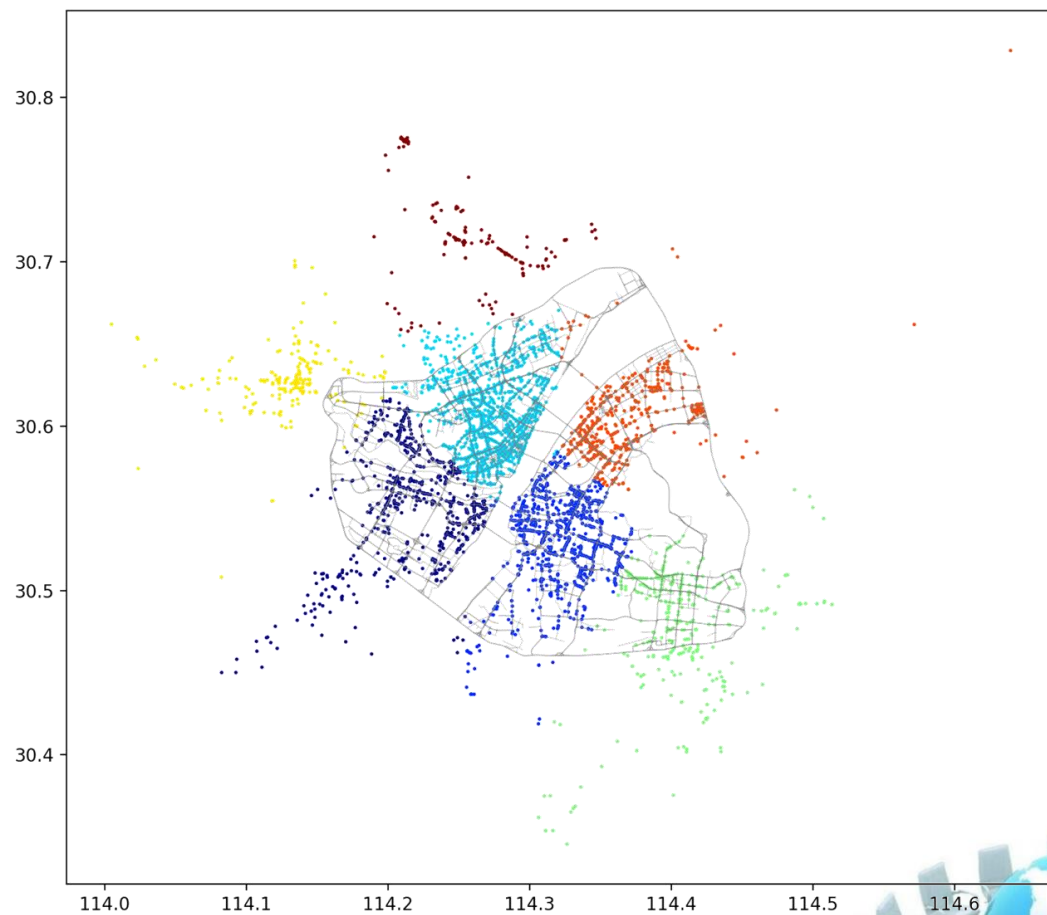
聚类结果



测绘与地理信息学院
COLLEGE OF SURVEYING AND GEO-INFORMATICS



$n = 20$



$n = 7$



1个核心思想：基于密度

2个算法参数：邻域半径 ϵ 和最少点数目MinPts

3种点的类别：核心点，边界点、噪声点

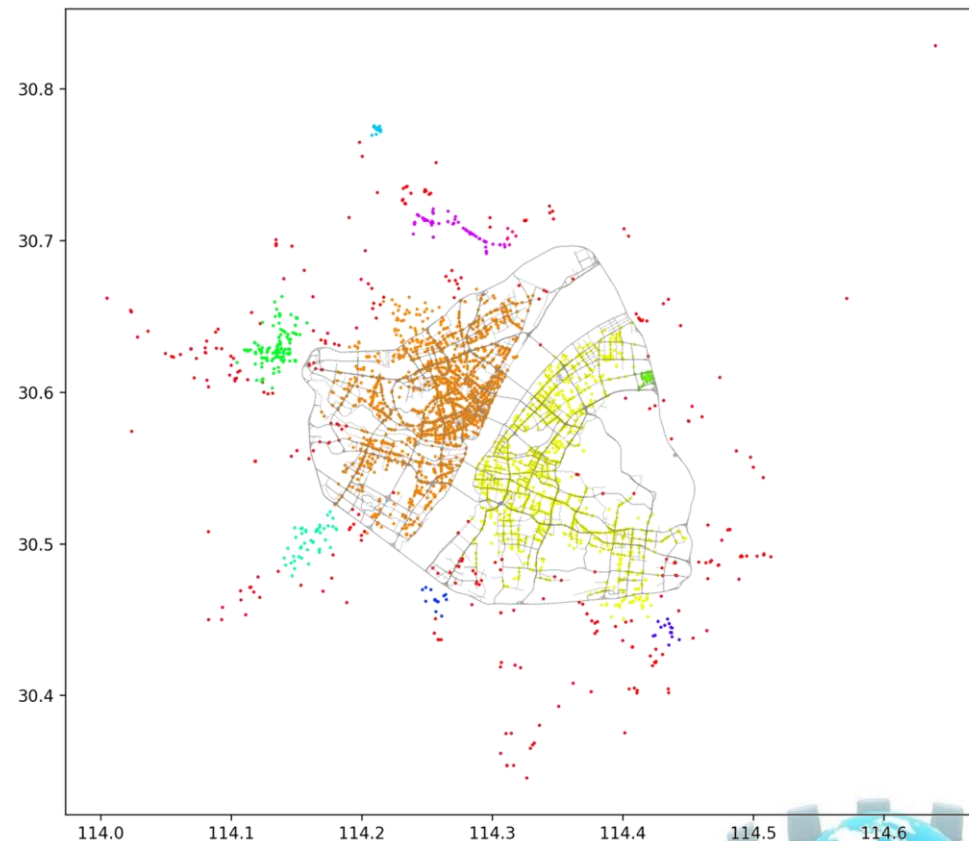
4种点的关系：密度直达，密度可达，密度相连，非密度相连

Step1. 从数据集中任意选取一个数据对象点 p ;

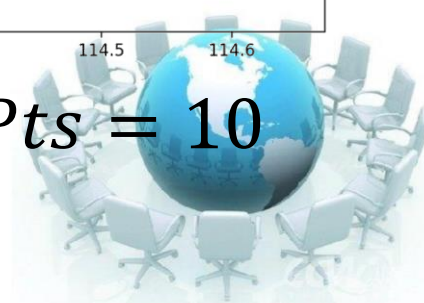
Step2. 如果对于参数 ϵ 和 MinPts，所选取的数据对象点 p 为**核心点**，则找出所有从 p **密度可达**的数据对象点，形成一个簇；

Step3. 如果选取的数据对象点 p 是**边缘点**，选取另一个数据对象点；

Step4. 重复（2）、（3）步，直到所有点被处理。



$$\epsilon = 0.01, \text{MinPts} = 10$$

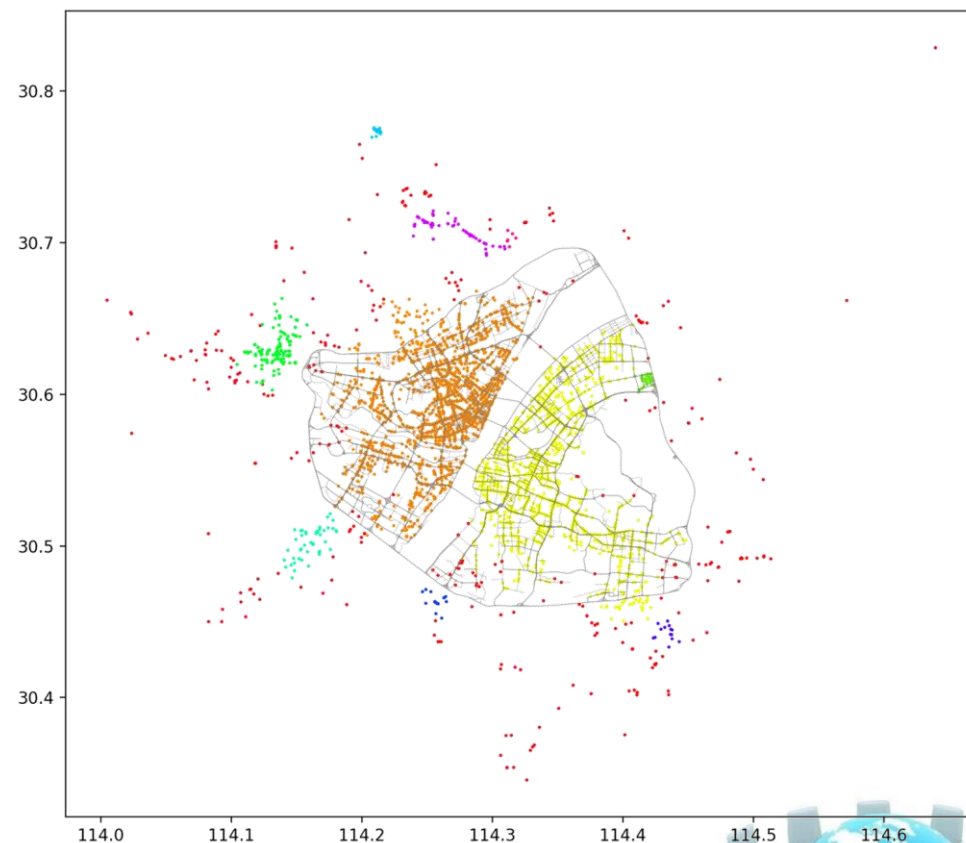


优点:

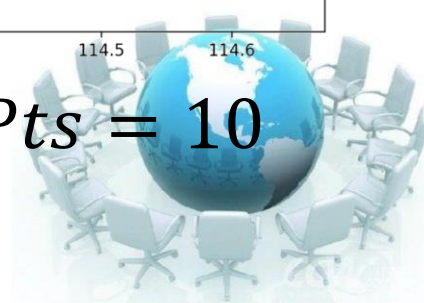
- 1) 可以自主计算聚类的数目，不需要人为指定
- 2) 不要求类的形状是凸的，可以是任意形状的
- 3) 对噪音不敏感
- 4) 算法应用参数少，只需要两个
- 5) 聚类结果几乎不依赖于节点的遍历顺序

缺点:

- 1) 如果样本集较大时，聚类收敛时间较长
- 2) 聚类效果依赖于距离公式的选取
- 3) 不适合数据集中密度差异很大的情形



$$eps = 0.01, MinPts = 10$$



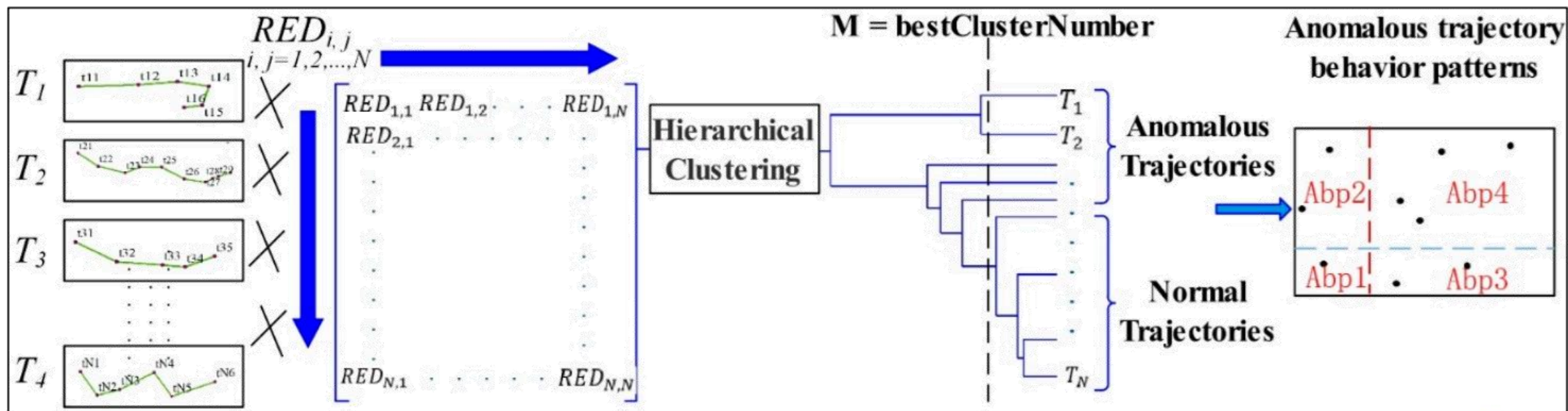
异常轨迹分析



测绘与地理信息学院
COLLEGE OF SURVEYING AND GEO-INFORMATICS

1. 计算轨迹距离/相似性矩阵
2. 进行层次聚类

3. 通过统计指标确定最佳聚类个数
4. 根据统计指标识别异常轨迹



Detecting anomalous trajectories and behavior patterns
using hierarchical clustering from taxi GPS data
(Wang, Y., Qin, K., Chen, Y., & Zhao, P.)



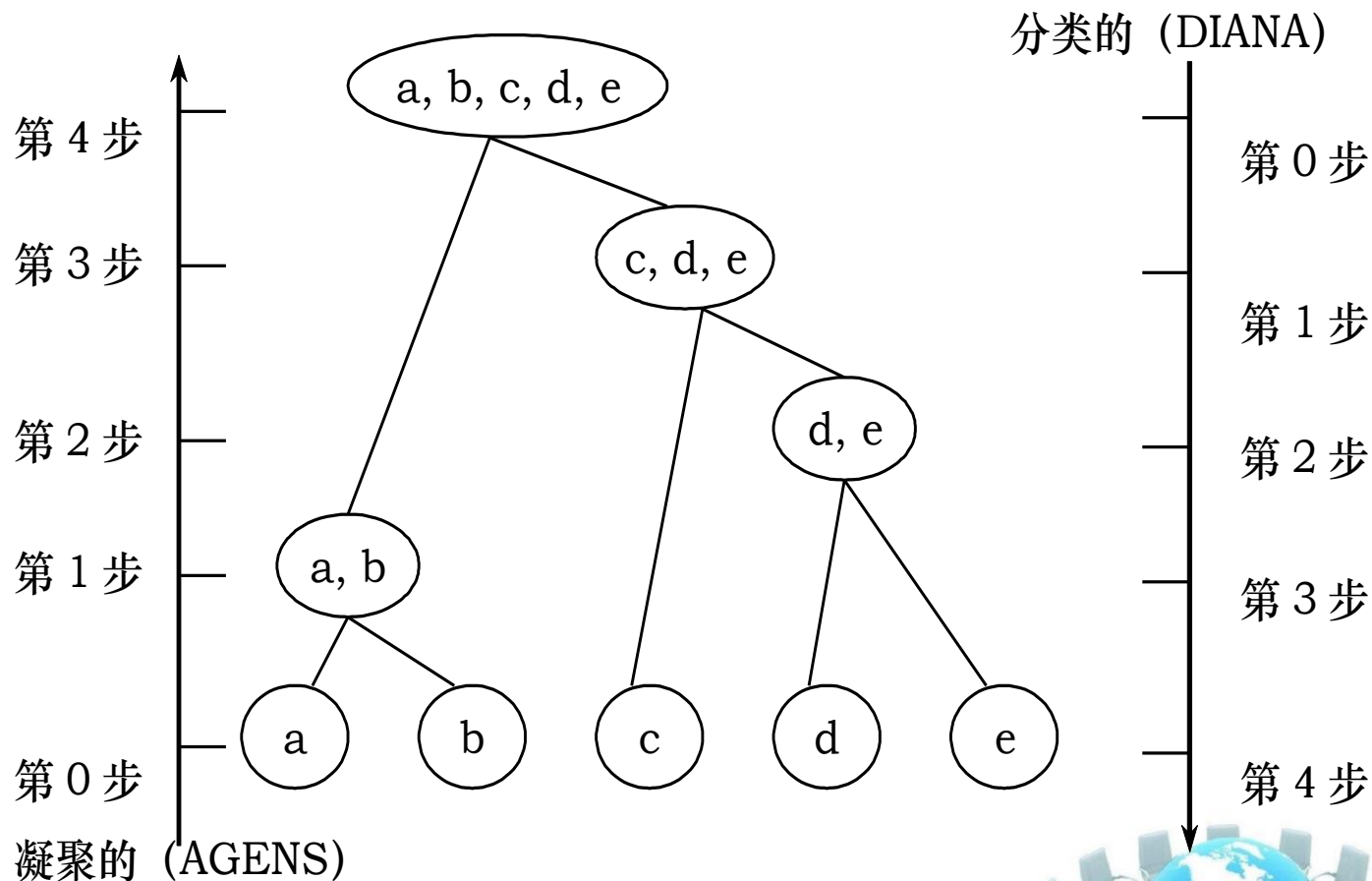
层次聚类



层次聚类按数据分层建立簇，形成一棵以簇为节点的树，称为**聚类图**。

按**自底向上**层次分解，则称为**凝聚**的层次聚类。

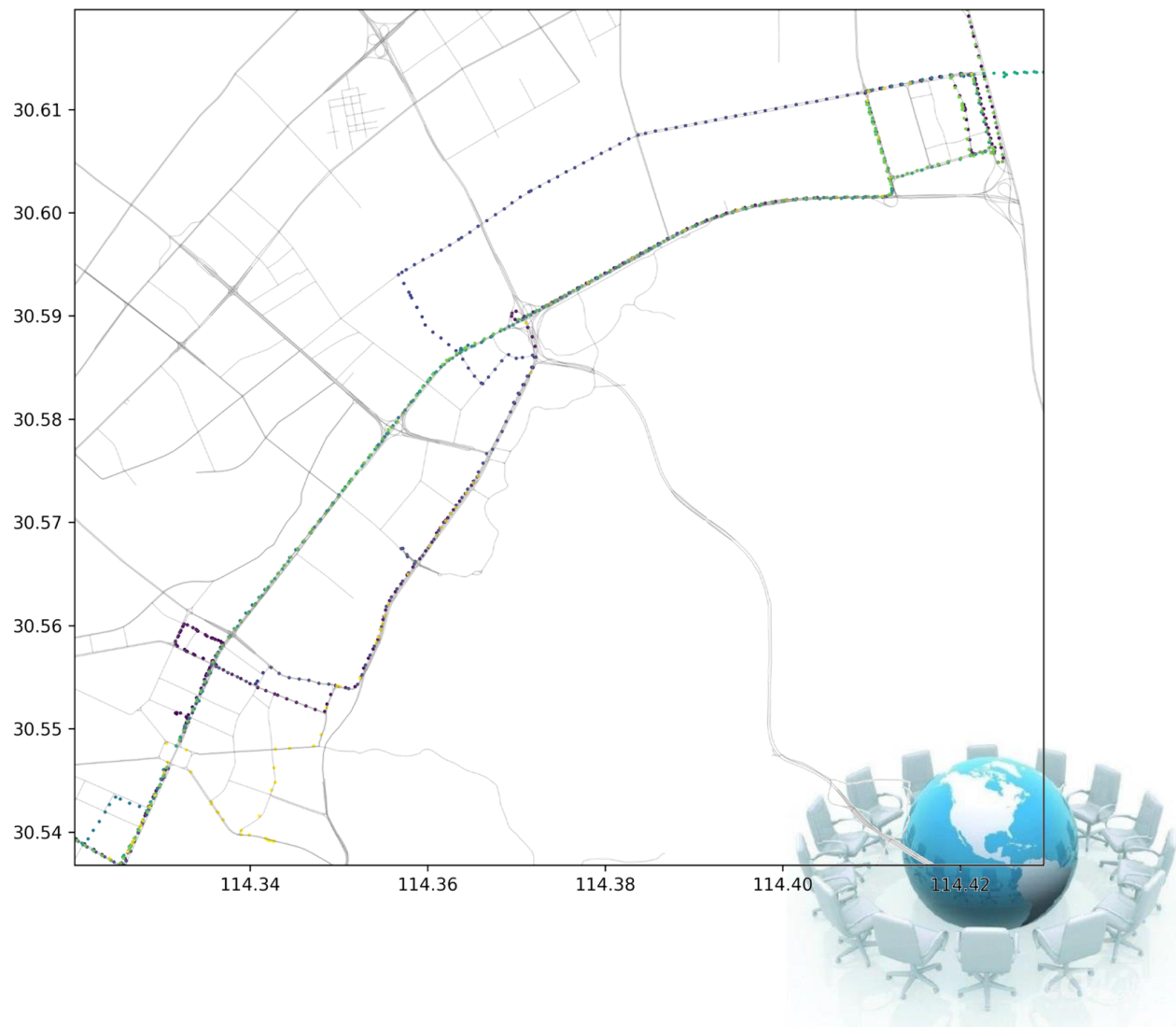
按**自顶向下**层次分解，就称为**分裂**的层次聚类。



以凝聚的层次聚类为例：

假设有 n 个待聚类的样本：

1. 初始化, 将每个样本都视为一个聚类；
2. 计算各个聚类之间的**相似度**；
3. 寻找最近的两个聚类，将他们归为一类；
4. 重复2，3；直到所有样本归为一类。



层次聚类



相似度确定

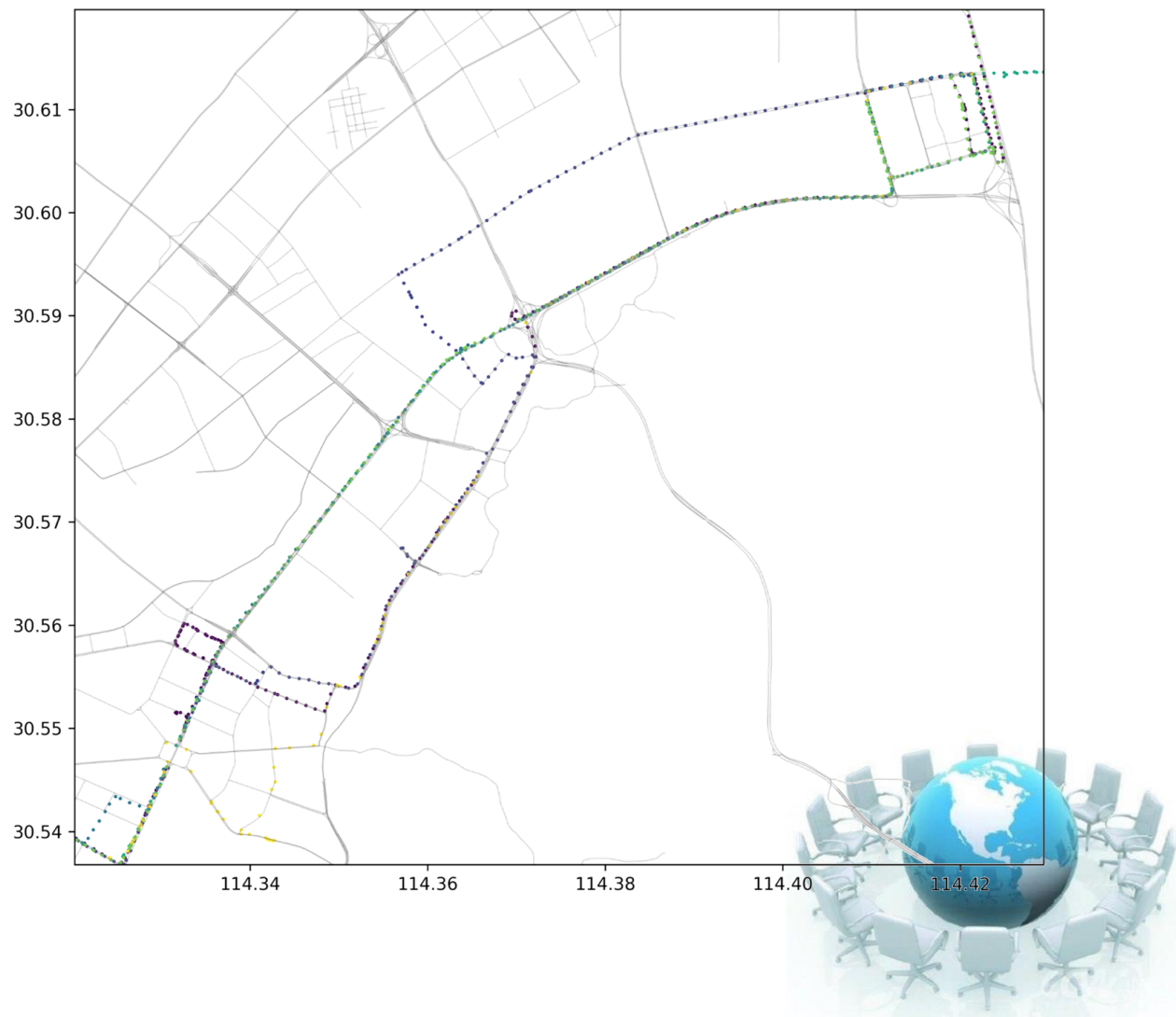
原则：最大化类间距离、最小化类内距离

$$SSW(M) = \max_t \left\{ \max_{i,j} (1 - IED(T_i, T_j)_{T_i \neq T_j \in C_t}) \right\} + \sum_{|C_t|=1} 1$$

$$SSB(M) = \sum_{t=1}^M \sum_{s>t}^M \min \left(1 - IED(T_i, T_j)_{T_i \in C_t, T_j \in C_s} \right)$$

其中IED为两轨迹直接距离的描述,M为设置的聚类数。

SSW越小则类内距离越小,SSB越大则类间距离越大,聚类效果越好。



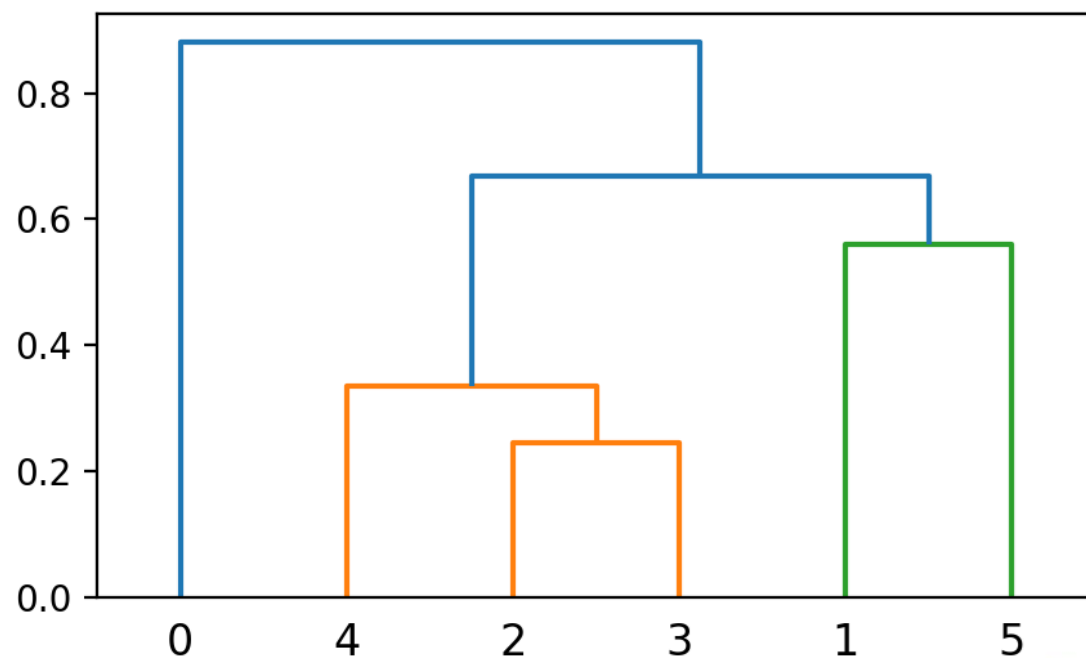
统计指标 — 聚类数确定

$$\text{WB-index} = M * \text{SSW}(M) / \text{SSB}(M)$$

$$\text{CH-index} = \frac{\text{SSB}(M) / M - 1}{\text{SSW}(M) / N - M}$$

$$\text{Xu-index} = \log \sqrt{(\text{SSW}(M)) / N^2} + \log M$$

	WB-index	CH-index	Xu-index
2	3.83517	0.41871	-1.17908
3	4.31683	0.28088	-0.77362
4	5.75577	0.22085	-0.48593
5	7.19472	0.18148	-0.26279



异常轨迹探测结果



测绘与地理信息学院
COLLEGE OF SURVEYING AND GEO-INFORMATICS

