

系统工程导论

开课单位：清华大学自动化系
主讲教师：胡坚明 副教授

模块二：系统分析

聚类分析方法

系统工程导论

内容和重点

■本章内容

- 聚类问题的一般性描述
- K-均值聚类方法
- 系统聚类方法
- 动态聚类方法
- 基于自组织映射 (SOM) 的聚类方法
- 应用实例

3

聚类分析方法

■引言

物以类聚，人以群分！

聚类问题出现在我们社会生活的方方面面

- 世界——发达国家/发展中国家 (GDP)
- 中国——东部/中部/西部 (经济状况)
- 大学——重点大学/一般院校 (综合实力)
- 专业——文、理、工、医等 (学习内容)
- 学生——擅长数理、擅长推理等 (学习成绩)
-

集合

子集

测度

4

聚类分析方法

■引言

- 聚类问题实际上是将包含若干元素的集合，按照某种测度，划分成若干子类。
- 测度是指定义在每个类上的函数，我们的目标就是使其达到最大或最小。
- 聚类问题的本质是划分问题——属于NP难问题

5

聚类分析方法

■引言

举例——推研面试考核

- 自动化系每年大约有140个毕业生，每名学生推研前完成了45门课程的学习。不同老师如何根据这些课程成绩，考察其所关心的能力？
- 这实际上构成了一个140×45的矩阵，不同老师可以根据不同目的构造一个函数（评价指标、测度）、使得出现不同分类。

6

聚类分析方法——一般描述

■条件

给定一组对象，用以下指标集表示

$$J(M) = \{1, 2, \dots, M\}$$

给定一个对 $J(M)$ 的所有非空子集有定义的
实值函数

$$\rho(\varpi) \quad \forall \varpi \subseteq J(M), \varpi \neq \emptyset$$

其中 \emptyset 表示空集

7

聚类分析方法——一般描述

■任务

确定 $\Omega = \{\varpi_i \subseteq J(M), 1 \leq i \leq k\}$ 满足

$$\bigcup_{1 \leq i \leq k} \varpi_i = J(M), \quad \varpi_i \cap \varpi_j = \emptyset \quad \forall i \neq j$$

使下述目标函数最小(或最大)

$$\sum_{i=1}^k \rho(\varpi_i)$$

k —一般未知，要根据最优目标函数的变化情况确定

8

聚类分析方法——一般描述

■对象——向量聚类

$$x(t) = [x_1(t) \ x_2(t) \ \dots \ x_n(t)]^T$$

$1 \leq t \leq N$

$$\rho(\varpi) = \sum_{i \in \varpi} (x(i) - e_{\varpi}(x))^T (x(i) - e_{\varpi}(x))$$

$$\forall \varpi \subseteq J(N)$$

n —考试科目
 N —学生总数
 $x(t)$ — i 同学考试成绩

与中心点的距离，
即紧密程度

其中

$$e_{\varpi}(x) = \frac{1}{|\varpi|} \sum_{i \in \varpi} x(i)$$

子类中心

问题

$$\min_{\Omega} \sum_{i=1}^k \rho(\varpi_i)$$

尽可能紧密！

9

聚类分析方法——一般描述

■对象——变量聚类

$$X_i = [x_i(1) \ x_i(2) \ \dots \ x_i(N)]$$

$1 \leq i \leq n$

$$\rho(\varpi) = \min_{i \in \varpi, j \in \varpi} r(x_i, x_j) \quad \forall \varpi \subseteq J(n)$$

其中 $r(x_i, x_j)$

r 是两个变量的协方差，我们希望类内协方差的值越大越好。即：类内两个变量的夹角越小越好。保守准则是先找到类内夹角最大的那两个变量。

$$= \frac{1}{N-1} \sum_{i=1}^N \frac{(x_i(t) - \bar{x}(x_i))(x_i(t) - \bar{x}(x_j))}{\sqrt{\delta^2(x_i) \delta^2(x_j)}}$$

类间：我们希望整体分类结果是相关性越强越好！

$$\max_{\Omega} \sum_{i=1}^k \rho(\varpi_i)$$

按照“最大最不相关”
进行分类

10

系统工程导论

■例：分片线性逼近问题

对象 $(y(t), x(t)), y(t) \in R, x(t) \in R^n$
 $1 \leq t \leq N$

$$\rho(w) = \min_{\alpha \in R, \beta \in R^n} \sum_{t \in \Omega} \left(y(t) - (\alpha + \beta^T x(t)) \right)^2$$

$\forall \Omega \subseteq J(N)$ 使优化问题有唯一解

$$\text{问题 } \min_{\Omega} \sum_{i=1}^k \rho(w_i)$$

11

系统工程导论

考虑向量聚类问题

$$\min_{\Omega} \sum_{i=1}^k \sum_{t \in \Omega_i} \left(x(t) - e_{w_i}(x) \right)^T \left(x(t) - e_{w_i}(x) \right)$$

$$\text{其中 } e_{w_i}(x) = \frac{1}{|\Omega_i|} \sum_{t \in \Omega_i} x(t)$$

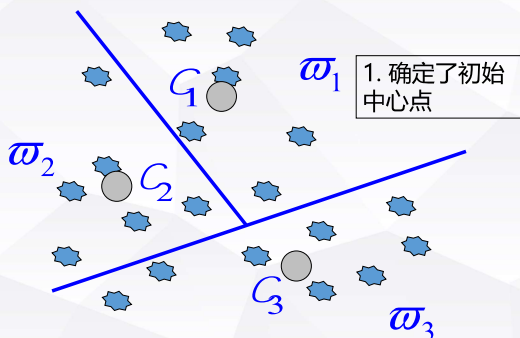
如果已有一个分类 $\Omega = \{\Omega_i, 1 \leq i \leq k\}$

可以先求出 k 个集合的中心点(样本均值)

$$c_i = e_{\Omega_i}(x) \in R^n, 1 \leq i \leq k$$

12

系统工程导论



13

系统工程导论

然后确定一个新的分类

$$\hat{\Omega} = \{\hat{\Omega}_i, 1 \leq i \leq k\}$$

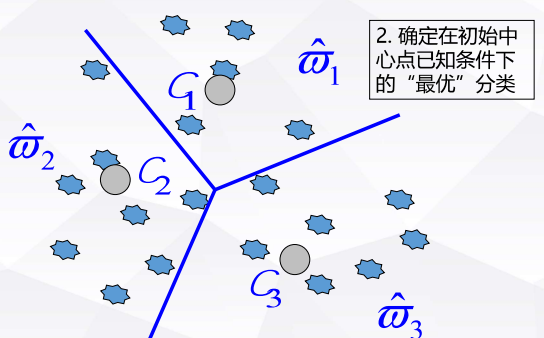
其中

$$\hat{\Omega}_i = \left\{ t \in J(N) \mid \left(x(t) - c_i \right)^T \left(x(t) - c_i \right) \leq \left(x(t) - c_j \right)^T \left(x(t) - c_j \right), \forall j \right\}$$

也就是说新的分类中每个 $x(t)$ 都距离 c_i 更近!

14

系统工程导论



15

系统工程导论

$$\text{由 } \hat{w}_i = \left\{ t \in J(N) \mid \left(x(t) - c_i \right)^T \left(x(t) - c_i \right) \leq \left(x(t) - c_j \right)^T \left(x(t) - c_j \right), \forall j \right\}$$

$$\text{可知 } \sum_{i=1}^k \sum_{t \in \hat{\Omega}_i} \left(x(t) - c_i \right)^T \left(x(t) - c_i \right) \leq \sum_{i=1}^k \sum_{t \in \Omega_i} \left(x(t) - c_i \right)^T \left(x(t) - c_i \right)$$

又因为 $c_i = e_{\Omega_i}(x) \in R^n, 1 \leq i \leq k$

$$\text{所以 } \sum_{i=1}^k \sum_{t \in \hat{\Omega}_i} \left(x(t) - c_i \right)^T \left(x(t) - c_i \right) \leq \sum_{i=1}^k \rho(w_i)$$

回忆向量分类的准则

16

系统工程导论

另外, 对于优化问题

$$\min_{\alpha \in R^n} \sum_{t \in \Omega_i} \left(x(t) - \alpha \right)^T \left(x(t) - \alpha \right)$$

$$\text{由于 } \frac{\partial \left(\sum_{t \in \Omega_i} \left(x(t) - \alpha \right)^T \left(x(t) - \alpha \right) \right)}{\partial \alpha} = -2 \sum_{t \in \Omega_i} \left(x(t) - \alpha \right)$$

$$\text{所以最优解为 } \hat{\alpha}_i = \frac{1}{|\hat{\Omega}_i|} \sum_{t \in \hat{\Omega}_i} x(t) = e_{\hat{\Omega}_i}(x)$$

3. 即最优解应该是该类的中心。

17

系统工程导论

$$\begin{aligned} \sum_{i=1}^k \rho(\hat{w}_i) &= \sum_{i=1}^k \sum_{t \in \hat{\Omega}_i} \left(x(t) - e_{\hat{\Omega}_i}(x) \right)^T \left(x(t) - e_{\hat{\Omega}_i}(x) \right) \\ &\leq \sum_{i=1}^k \sum_{t \in \Omega_i} \left(x(t) - c_i \right)^T \left(x(t) - c_i \right) \\ &\leq \sum_{i=1}^k \rho(w_i) \end{aligned}$$

4. 判断当前 c_i 是否是最优解的依据。只要有一个分类内部的中心点和前面已知的中心点不吻合, 就没有找到最优分类。

并且, 只要有一个 $c_i \neq e_{\hat{\Omega}_i}(x)$

$$\text{就一定有 } \sum_{i=1}^k \rho(\hat{w}_i) < \sum_{i=1}^k \rho(w_i)$$

18

系统工程导论

■步骤

相比前面做法, 此处进行了简化!

1) 在 $x(t), t \in J(N)$ 中随机选取 $c_i, 1 \leq i \leq k$

2) 确定分类 $\Omega = \{\Omega_i, 1 \leq i \leq k\}$, 其中

$$\Omega_i = \left\{ t \in J(N) \mid \left(x(t) - c_i \right)^T \left(x(t) - c_i \right) \leq \left(x(t) - c_j \right)^T \left(x(t) - c_j \right), \forall j \right\}$$

3) 如果 $e_{\Omega_i}(x) = c_i, 1 \leq i \leq k$, 停止, 否则令

$$c_i = e_{\Omega_i}(x), 1 \leq i \leq k, \text{ 回到2) 继续迭代}$$

是否收敛?

19

系统工程导论

■进一步解释

① 选中心点

首先确定分类数目 k , 然后在所有样本中挑选 k 个作为初始中心点,

例如, 可简单地取

$$c(i) = x(i), i = 1, 2, \dots, k$$

② 逐个利用每个样本修改中心点

对所有的样本, 顺序进行下述计算:

第一、将其归入与其最近的中心点所在的类;

第二、重新计算该类的重心, 并用新的重心替换中心点

20

系统工程导论

■进一步解释

③ 停止准则

如果上述步骤开始时和结束后的中心点差别很小，停止分类

可以证明，该迭代算法是收敛的

21

系统工程导论

系统聚类算法是一种贪婪算法！

考虑变量聚类问题

$$\max_{\Omega} \sum_{i=1}^k \rho(w_i)$$

其中

$$\rho(w) = \min_{i \in w, j \in w} r(x_i, x_j) \quad \forall w \subseteq J(n)$$

22

系统工程导论

■基本步骤

- (1) 首先将每个变量视为一类，得到 n 类变量
- (2) 每次选择最相关的两个类合并，顺序得到 $n-1, n-2, n-3, \dots$ 直至一类变量
- (3) 记录合并过程生成聚类谱系图
- (4) 设定阈值，根据聚类谱系图决定最终分类

23

系统工程导论

假定所有变量已经合并成了 m 类

$$w_i \subseteq J(n), 1 \leq i \leq m$$

在确定下一步合并哪两类时，选择

$$1 \leq \alpha < \beta \leq m$$

满足

$$\rho(w_\alpha \cup w_\beta) \geq \rho(w_i \cup w_j) \quad \forall 1 \leq i < j \leq m$$

将 w_α 和 w_β 合并，从而得到 $m-1$ 类

24

系统工程导论

例 四个变量依相关性聚类

$$R = \begin{bmatrix} 1 & 0.334 & 0.950 & 0.597 \\ & 1 & 0.217 & -0.117 \\ & & 1 & 0.569 \\ & & & 1 \end{bmatrix} \begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{matrix}$$

森林面积

森林覆盖率

林木蓄积量

草原面积

$$w_1 = \{1\} \quad w_2 = \{2\} \quad w_3 = \{3\} \quad w_4 = \{4\}$$

25

系统工程导论

$$\rho = \begin{bmatrix} 1 & 0.334 & 0.950 & 0.597 \\ & 1 & 0.217 & -0.117 \\ & & 1 & 0.569 \\ & & & 1 \end{bmatrix} \begin{matrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{matrix}$$

$$\rho = \min R(x_i, x_j)$$

聚类标准

$$\rho(w_\alpha \cup w_\beta) \geq \rho(w_i \cup w_j)$$

聚类: $w_1 \cup w_3$

26

系统工程导论

$$\rho = \begin{bmatrix} 1 & 0.334 & 0.950 & 0.597 \\ & 1 & 0.217 & -0.117 \\ & & 1 & 0.569 \\ & & & 1 \end{bmatrix} \begin{matrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{matrix}$$

$$\rho = \min R(x_i, x_j)$$

$$w_5 = w_1 \cup w_3 \quad w_6 = w_2 \quad w_7 = w_4$$

$$\rho(w_5 \cup w_6) = \min\{0.334, 0.217\} = 0.217$$

$$\rho(w_5 \cup w_7) = \min\{0.597, 0.569\} = 0.569$$

$$\rho(w_6 \cup w_7) = -0.117$$

27

系统工程导论

■计算方法

$$\rho = \begin{bmatrix} 1 & 0.334 & 0.950 & 0.597 \\ & 1 & 0.217 & -0.117 \\ & & 1 & 0.569 \\ & & & 1 \end{bmatrix} \begin{matrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{matrix}$$

$$\rho = \min R(x_i, x_j)$$

$$w_5 = w_1 \cup w_3 \quad w_6 = w_2 \quad w_7 = w_4$$

$$\rho(w_5 \cup w_6) = \rho(w_1 \cup w_3 \cup w_2)$$

$$= \min\{R(x_1, x_2), R(x_3 \cup x_2)\}$$

$$= \min\{0.334, 0.217\}$$

$$= 0.217$$

28

系统工程导论

■计算方法

$$\rho = \begin{bmatrix} 0.95 & 0.217 & 0.569 \\ & 1 & -0.117 \\ & & 1 \end{bmatrix} \begin{matrix} w_5 \\ w_6 \\ w_7 \end{matrix}$$

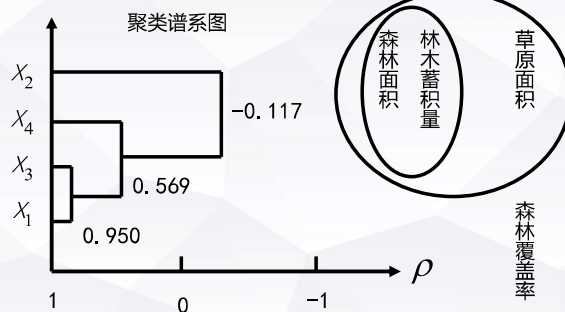
$$w_8 = w_5 \cup w_7 \quad w_9 = w_6$$

$$\rho(w_8 \cup w_9) = \min\{0.217, -0.117\} = -0.117$$

$$\text{最后 } w_{10} = w_8 \cup w_9 \quad \rho(w_{10}) = -0.117$$

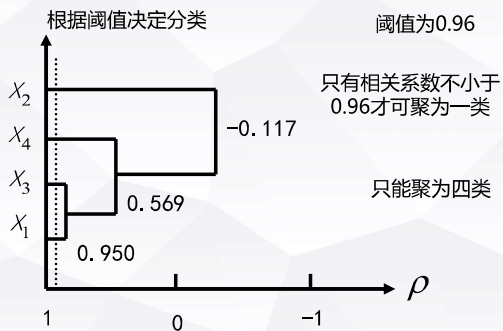
29

系统工程导论



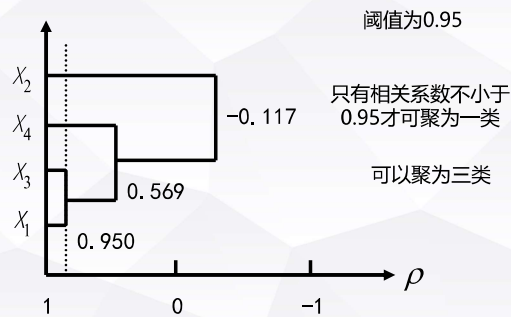
30

系统工程导论



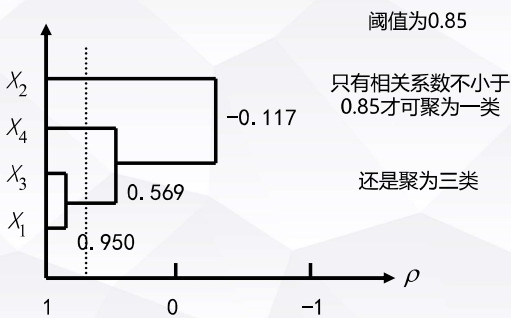
31

系统工程导论



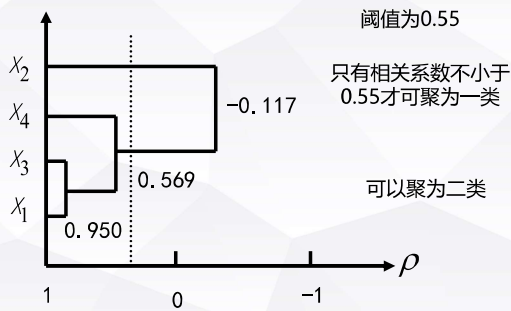
32

系统工程导论



33

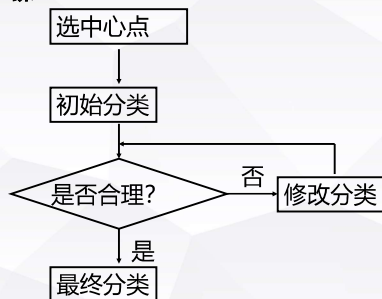
系统工程导论



34

系统工程导论

基本步骤



35

系统工程导论

例：对21个工厂的产品按两个质量指标分类

t	$x_1(t)$	$x_2(t)$	t	$x_1(t)$	$x_2(t)$
1	0	6	12	-2	2
2	0	5	13	-3	2
3	2	5	14	-3	0
4	2	3	15	-5	2
5	4	4	16	1	1
6	4	3	17	0	-1
7	5	1	18	0	-2
8	6	2	19	-1	-1
9	6	1	20	-1	-3
10	7	0	21	-3	-5
11	-4	3			

36

系统工程导论

按批修改方法

1) 选中心点 (密度法)

取阈值 $a_1 = 2$ ，对每个样本 $x(t)$ ，用

$\eta(t)$ 表示下述集合元素个数 (密度)

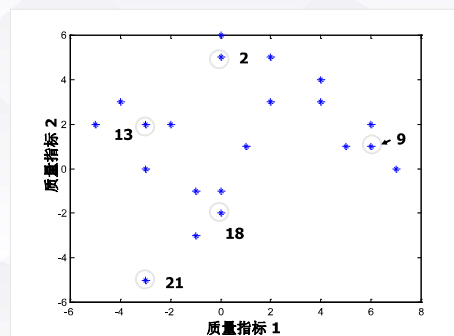
$$\{x(k), 1 \leq k \leq 21, k \neq t \mid d_2(x(k), x(t)) \leq a_1\}$$

例如 $\eta(2) = 2, \eta(9) = 3, \eta(13) = 4$

$\eta(18) = 3, \eta(21) = 0$

37

系统工程导论



38

系统工程导论

按批修改方法

取密度最大的点为第一中心点

$$c(1) = x(13) \quad (\eta(13) = 4)$$

取阈值 $a_2 = 2a_1 = 4$ ，考虑密度第二大的点

$$\eta(9) = 3$$

由于 $d_2(x(9), c(1)) > a_2$ ，取为第二中心点

$$c(2) = x(9)$$

反之不用该点为中心点

39

系统工程导论

按批修改方法

考虑密度第三大的点 $\eta(18) = 3$

由于 $d_2(x(18), c(1)) > a_2$

$$d_2(x(18), c(2)) > a_2$$

取为第三中心点 $c(3) = x(18)$

最终得到

$$c(1) = x(13) \quad c(2) = x(9) \quad c(3) = x(18)$$

$$c(4) = x(2) \quad c(5) = x(21)$$

40

系统工程导论

■按批修改方法

2) 初始分类 (最近中心点原则)

距某个中心点最近的点划归该中心点对应的类

如果对任意的 $j \neq i$

$$d_2(x(i), c(i)) \leq d_2(x(i), c(j))$$

将 $x(i)$ 划归第 i 类

41

系统工程导论

■按批修改方法

3) 修改分类 (修改中心点)

用每类的重心替换中心点

重心: 每类样本的均值向量

例如, 若某类样本为

$$x(7), x(8), x(9), x(10)$$

其重心为

$$\frac{1}{4}(x(7) + x(8) + x(9) + x(10))$$

42

系统工程导论

■按批修改方法

4) 停止准则

如果当前的中心点和对应的重心点都很接近, 停止分类

可以证明, 该迭代算法是收敛的。这就是说, 只要迭代次数足够多, 上述中心点和重心点可以任意接近

43

系统工程导论

■说明

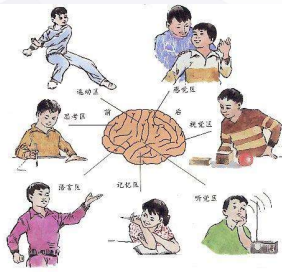
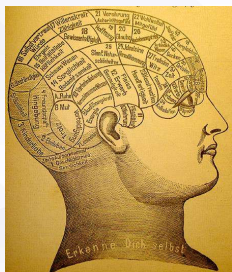
- 1) 动态聚类方法的结果一般和初始中心点有关, 选择初始中心点很重要
- 2) 聚类分析是实践性很强的问题, 所介绍的基本方法有很多变形和发展
- 3) 同系统聚类法对应的还有一大类逐渐分解的分类方法

44

系统工程导论

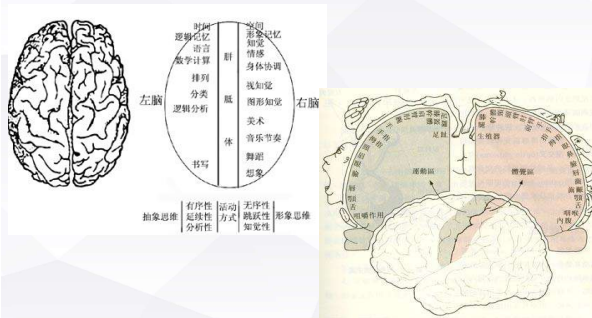
■基于自组织映射 (SOM) 的聚类

神经系统的自组织功能



45

系统工程导论



46

系统工程导论

自组织映射: Self-organizing Map, SOM网络是一种有效的**聚类**和聚类结果可视化工具。

它定义了从 n 维输入空间到规则的二维输出节点矩阵的非线性映射。

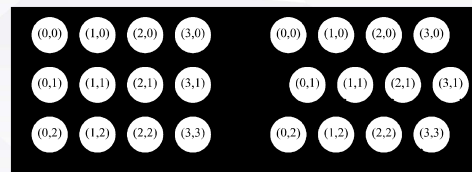
通常SOM是一个无监督的两层神经网络。

- 第一层是输入层
- 第二层是二维输出层。

输出层的点阵结构可以是矩形, 也可以是六边形。

47

系统工程导论



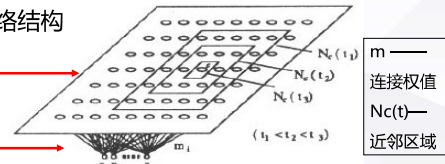
SOM网络结构

二维

输出层

输入

层



48

系统工程导论

SOM网络结构

- 神经元呈平面分布
- 输入向量的每一维连接到每个结点
- 结点间依分布位置关系而有相互作用

神经元计算特性

- 向量匹配, 与输入向量最佳匹配者称为Winner

SOM特性

- 经过适当训练后, 结点对输入的响应表现出一定的规律性, 这些规律性反映了输入样本集内部的一些特性, 即样本分布密度及样本间拓扑关系。

49

系统工程导论

■SOM网络中几个重要概念

- **象**: 随着不断学习, 对于某个输入向量, 对应Winner结点逐步趋于固定, 该结点称作该样本的象
- **原象**: 若向量 x 的象为结点 i , 则 x 为 i 的一个原象
- **象密度**: 同一结点上原象的数目, 即有多少样本映射到该结点。
- **象密度图**: 按节点位置把象密度的相对值画到一张图上, 这张图称为象密度图。

50

系统工程导论

系统工程导论

53

54

55

56

57

站号	站名、站址名称	始发站	长春大学	卫星广场	农研	客车厂	保养厂	湖大路	工农广场	紫荆花	自由大路
1	始发站	0	0	0	0	0	0	0	0	0	0
2	长春大学	23.0	0	0	0	0	0	0	0	0	0
3	卫星广场	28.30	43.90	0	0	0	0	0	0	0	0
4	农研	49.30	40.30	40.70	0	0	0	0	0	0	0
5	客车厂	49.50	52.90	41.20	50.20	0	0	0	0	0	0
6	保养厂	40.40	30.60	52.40	43.00	55.40	0	0	0	0	0
7	湖大路	42.80	37.30	63.50	58.70	53.90	41.60	0	0	0	0
8	工农广场	31.10	30.90	36.80	42.30	53.90	32.90	44.30	0	0	0
9	紫荆花	36.00	44.40	46.30	62.70	46.20	54.80	21.50	53.80	0	0
10	自由大路	25.40	32.20	37.10	41.10	65.30	41.80	37.20	24.80	36.60	0
11	丰华路	60.90	62.30	51.10	51.50	57.30	52.20	51.20	53.90	45.10	53.30
12	惠康路	39.40	52.60	52.60	34.70	58.90	41.20	38.00	41.90	46.10	40.50
13	解放大路	39.40	53.60	54.60	58.10	66.60	45.10	36.10	60.90	39.70	46.10
14	人民广场	25.50	30.50	45.60	52.60	61.40	36.40	45.00	31.70	45.30	47.00
15	国贸广场	28.70	19.90	35.70	27.50	42.20	30.80	55.10	28.60	58.10	32.10
16	医学院	28.80	38.20	32.90	41.50	64.80	29.80	49.60	49.00	31.40	30.40
17	解放大路	42.80	62.80	44.30	30.30	64.30	43.70	46.90	45.20	48.60	53.10
18	胜利大街	18.90	30.90	33.90	33.00	50.60	62.90	39.20	52.50	40.20	46.20
19	长春站	14.90	33.10	51.90	46.60	66.60	44.40	55.30	36.80	58.60	45.90
20	终点站	28.80	26.90	44.00	36.80	47.60	30.40	36.30	19.70	41.50	25.90

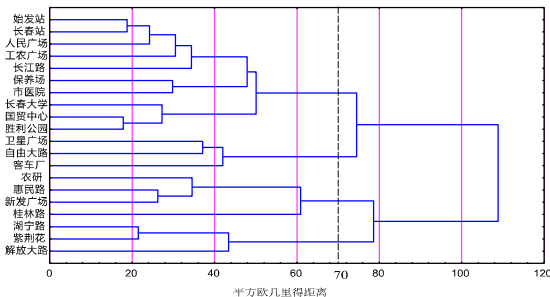
站点间平方欧几里得距离 (续)

站号	站点名称	桂林路	惠民路	解放大路	人民广场	国贸中心	市医院	新发广场	胜利公园	长江路	长春站
11	桂林路	0									
12	惠民路	54.4	0								
13	解放大路	55.0	31.6	0							
14	人民广场	51.6	34.9	41.2	0						
15	国贸中心	59.3	47.6	58.0	50.2	0					
16	市医院	38.1	40.9	46.1	34.6	42.8	0				
17	新发广场	41.7	26.3	34.6	40.9	42.9	42.8	0			
18	胜利公园	63.5	47.8	52.0	38.7	17.9	38.0	37.8	0		
19	长江路	53.3	29.1	41.0	30.4	35.6	44.2	41.0	33.8	0	
20	长春站	50.2	30.7	50.1	20.2	26.8	28.8	30.8	28.8	30.6	0

61

系统工程导论

■站点客流集散量聚类研究——聚类谱系图



62

系统工程导论

■站点客流集散量聚类研究——聚类结果

我们选择平方欧几里得距离为70将公交站点进行分类。分类结果为：

- 第一类：始发站、长春站、人民广场、工农广场、长江路开发区、保养场、市医院、长春大学、国贸中心和胜利公园；
- 第二类：卫星广场、自由大路和公交客车厂；
- 第三类：农研、惠民路和新发广场；
- 第四类：湖宁路、紫荆花酒店和解放大路。

63

系统工程导论

■基于SOM网络聚类方法——研究背景

- 网络交通流时变、非线性特征
- 网络交通状态演变规律研究的重要性
- 网络交通流短时预测的重要性
- 交通流数据的海量特性

问题：如何从中挖掘出交通状态的演变规律，对网络层次的交通流进行预测？

64

系统工程导论

■基于SOM网络聚类方法——研究方法

- 将网络交通流数据转化为N维时间序列

即： $F(t) = [f_1(t), f_2(t), \dots, f_d(t)]^T$

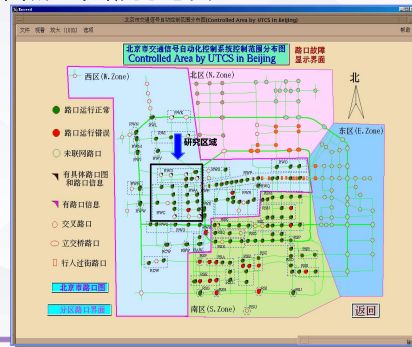
其中： $f_i(t)$ 表示 t 时段第 i 个路段的交通流量。 d 表示所研究交通网络的路段总数。

- 将上述数据作为SOM网络的输入，输出是一个 10×10 结点的二维平面。
- SOM聚类
- 基于聚类结果的网络交通状态分析与短时交通流预测研究。

65

系统工程导论

北京市西城区某路网示意图



66

■基于SOM网络聚类方法

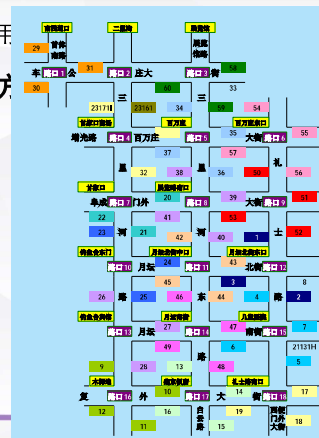
47个路段

2002年3月1日至3月30日一个月的交通流数据 (每个路口每15分钟一个数据)

即整个交通数据构成一个47维，长度为2786个有效数据

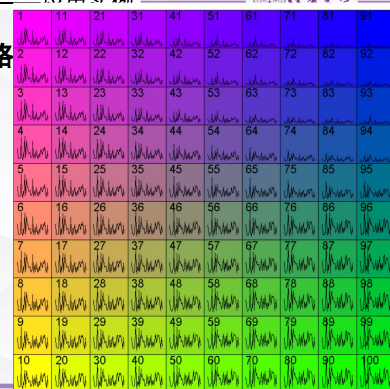
2500个用于训练

286个用于预测



■基于SOM网络

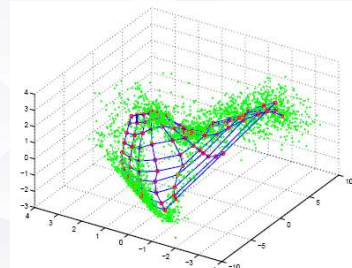
每一类用该类中所有流率向量的平均值来表示。对于每个结点，X轴表示路段名序列，Y轴表示流率 (veh/h)。左上角数字即聚类号。不同类用不同颜色表示，相邻结点具有相似颜色。



68

■基于SOM网络聚类方法——聚类结果

绿色点：样本点；红色点：SOM聚类中心结点

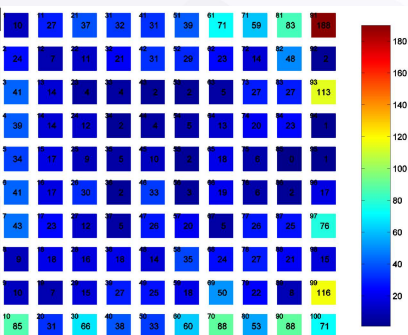


69

系统工程导论

■基于SOM网

左上角数字为类号，中间数字为该类中所包含的向量数目。颜色越趋向暖色表示结点数据多，反之少。



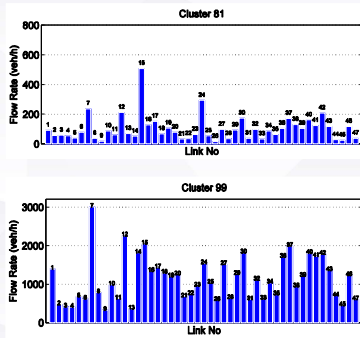
70

系统工程导论

典型交通流分布
模式类

X轴：路段号

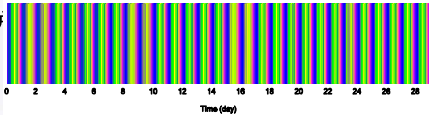
Y轴：流率

类81：包含83
个向量类99：包含116
个向量

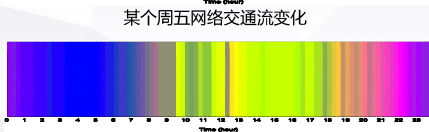
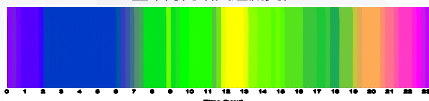
71

系统工程导论

聚类分析方法

交通流向量时间
序列分析每个时间段的数据
都用它所在的
类来表示。

X轴：时间点



72

系统工程导论

聚类分析方法——应用实例

SOM实质上起到了非线性降维的作用，由47维降到2维。我们还可以将其用于网络交通流“批量”预测。具体方法（kNN法）包含三个步骤：

(1) 创建 t 时刻的统计量，预测 $t+1$ 时刻流率

$$r(t) = [c_1(t), c_2(t), c_1(t-1), c_2(t-1), \dots, c_1(t-p+1), c_2(t-p+1)]^T$$

(2) 计算当前时间的统计量与前面所有统计量的欧拉距离

$$d(s) = \|r(t) - r(s)\|, s = 1, 2, 3, \dots, l$$

(3) 确定 k 个与当前统计量最为接近的统计量，下一时刻的流量的预测值即为用这 k 个统计量对应时刻的流率平均值。

$$F(t+1) = \frac{1}{k} [F(s_1) + F(s_2) + \dots + F(s_k)]$$

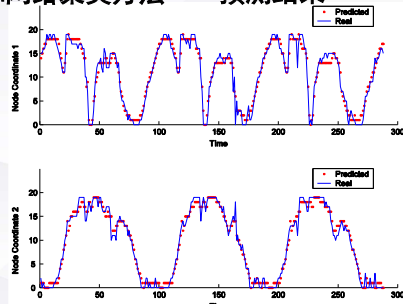
$t+1$ 时刻对应的结点坐标值也可以用前 k 个时段时的结点坐标的均值来计算。

73

系统工程导论

聚类分析方法——应用实例

■基于SOM网络聚类方法——预测结果

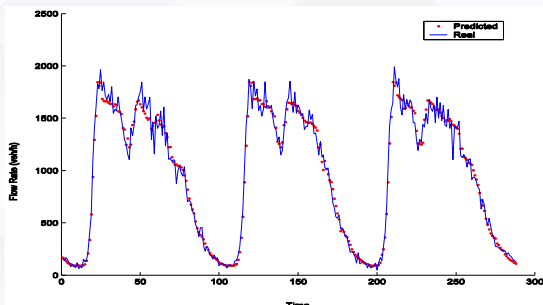
取 $p=4$ 

74

系统工程导论

聚类分析方法——应用实例

■基于SOM网络聚类方法——预测结果



75

系统工程导论

聚类分析方法——应用实例

■基于SOM网络聚类方法——预测结果

预测结果比较（ K 为最近邻数目）

k	kNN with SOM	kNN with PCA
	MAPE	MAPE
5	0.1619	0.1734
10	0.1414	0.1717
20	0.1433	0.1705
30	0.1519	0.1703
40	0.1393	0.1707

$$mape = \frac{\sum_{d=1}^D \sum_{t=1}^T |f_d(t) - \hat{f}_d(t)| / f_d(t)}{d \times l}$$

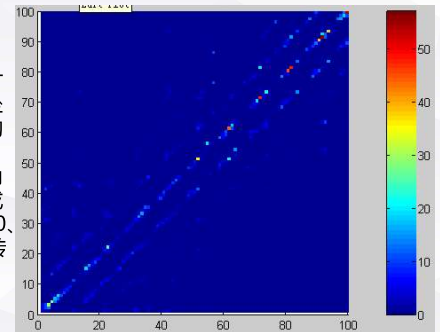
76

系统工程导论

聚类分析方法——应用实例

模式转移研究

横坐标是前一时
刻的标号，纵坐
标是下一时刻的
标号
可见集中在对
角线上（停留）
或行列标号相差
10、20的位置（只
转移到相邻的
pattern）

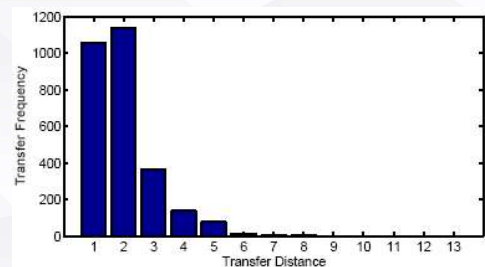


77

系统工程导论

聚类分析方法——应用实例

模式转移分布

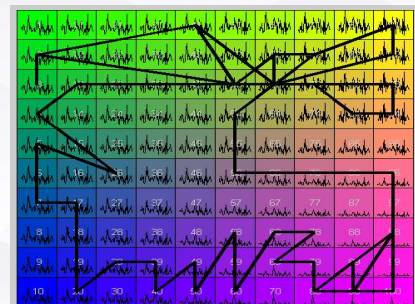


78

系统工程导论

聚类分析方法——应用实例

第25天内的转移路径

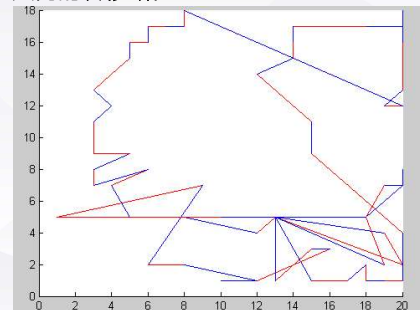


79

系统工程导论

聚类分析方法——应用实例

第23天内的转移路径



80

系统工程导论