

对不同结构样品进行识别的 PTR 算法研究及其应用*

王碧泉 黄汉明 范洪顺

(中国北京 100081 国家地震局地球物理研究所)

摘 要

阐述了 4 种模式识别方法. 在样品平面投影、典型样品分析和这几种模式识别方法的基础上, 提出了对不同结构样品进行识别的 PTR 方法. 文中还给出了两个应用实例. 结果表明, PTR 方法的效果较好.

关键词 模式识别; PTR 算法; 地震预报; 典型样品; 平面映射

1 引 言

模式识别在地震预报和潜在震源判定中的应用已很多, 有关模式识别方法及其应用的详述可参看王碧泉等(1989b)的论著. 除了数学家们发展了许多模式识别方法外, 地震学家们也提出了一些模式识别方法, 如王碧泉和王春珍(1988)提出的 ICHAM 方法, Keilis-Borok 和 Rotwain(1990a)提出的 CN 算法以及 Keilis-Borok 和 Kossobokov(1990b)提出的 M8 算法等. 显然, 方法研究和应用研究都是十分重要的.

实践中我们发现: 对某一组样品, 可能某一种方法分类效果较好, 但对另一组样品则可能另一种方法较好, 分类结果的好坏不仅与所选用的模式识别方法有关, 而且与被分类样品的结构有关.

本文目的旨在对几种模式识别方法研究的基础上, 结合平面映射结果, 给出对不同典型样品分别应采用的方法. 为此, 我们提出了“平面映射—典型样品分析—模式识别”一套方法(plane projection—typical sample analysis—pattern recognition), 简称为 PTR 算法. 文中还给出了应用 PTR 算法的两个实例.

2 几种模式识别方法

下面阐述几种典型的模式识别方法, 它们均属于分类判别方法.

* 地震科学联合基金会资助课题, 国家地震局地球物理研究所论著 93A0087.
1992 年 3 月 30 日收到初稿, 1993 年 1 月 11 日决定采用.

2.1 最小距离法

有关方法可参看王碧泉和陈祖荫(1989b)的论著.

2.2 Fisher 方法

对于两类问题, Fisher 方法仅要求已知两类样品的均值 \bar{X}_1 和 \bar{X}_2 以及两类样品的协方差矩阵 S_1 和 S_2 , 并不要求两类样品服从正态分布且协方差矩阵相等. Fisher 方法是一种线性分类方法. 对于两类问题线性分类器的一般形式为

$$h(X) = c_1 x_1 + c_2 x_2 + \cdots + c_n x_n + c_0 = C^T X + c_0 \quad (1)$$

若记 1 类样品为 $X_{11}, X_{12}, \cdots, X_{1N_1}$; 2 类样品为 $X_{21}, X_{22}, \cdots, X_{2N_2}$. 将它们代入上式得到

$$h(X_{11}), h(X_{12}), \cdots, h(X_{1N_1}) \quad h(X_{21}), h(X_{22}), \cdots, h(X_{2N_2})$$

2.2.1 求 $h(X)$ 的均值和方差

对两类样品, 分别求对应的 $h(X)$ 值的均值 η_p 和方差 σ_p^2 .

$$\begin{aligned} \eta_p &= \frac{1}{N_p} [h(X_{p1}) + h(X_{p2}) + \cdots + h(X_{pN_p})] \\ &= \frac{1}{N_p} [(C^T X_{p1} + c_0) + (C^T X_{p2} + c_0) + \cdots + (C^T X_{pN_p} + c_0)] \\ &= C^T \left[\frac{1}{N_p} (X_{p1} + X_{p2} + \cdots + X_{pN_p}) \right] + c_0 \\ &= C^T \bar{X}_p + c_0 \quad p = 1, 2 \end{aligned} \quad (2)$$

$$\begin{aligned} \sigma_p^2 &= \frac{1}{N_p} \sum_{i=1}^{N_p} [h(X_{pi}) - \eta_p]^2 \\ &= \frac{1}{N_p} \sum_{i=1}^{N_p} [C^T (X_{pi} + c_0 - (C^T \bar{X}_p + c_0))]^2 \\ &= \frac{1}{N_p} \sum_{i=1}^{N_p} C^T (X_{pi} - \bar{X}_p)^T (X_{pi} - \bar{X}_p) C \\ &= C^T S_p C \quad p = 1, 2 \end{aligned} \quad (3)$$

$\eta_1 - \eta_2$ 反映两样品经 $h(X)$ 作用后的分开程度, 而 σ_1^2, σ_2^2 反映两类样品经 $h(X)$ 作用后各自的密集程度.

2.2.2 求分类器

Fisher 分类器的设计原则是: C 和 c_0 的选取方法应使 $(\eta_1 - \eta_2)^2$ 尽量大, 而 $\sigma_1^2 + \sigma_2^2$ 尽量小. 或使

$$f = (\eta_1 - \eta_2)^2 / (\sigma_1^2 + \sigma_2^2) \quad (4)$$

达到极大. 满足这一条件的分类器是线性的. 下面由此条件推导求 C 的公式.

求 f 的极大值相当于解方程

$$\frac{\partial f}{\partial C} = 0 \quad \frac{\partial f}{\partial c_0} = 0$$

即

$$\frac{\partial f}{\partial \sigma_1^2} \cdot \frac{\partial \sigma_1^2}{\partial C} + \frac{\partial f}{\partial \sigma_2^2} \cdot \frac{\partial \sigma_2^2}{\partial C} + \frac{\partial f}{\partial \eta_1} \cdot \frac{\partial \eta_1}{\partial C} + \frac{\partial f}{\partial \eta_2} \cdot \frac{\partial \eta_2}{\partial C}$$

$$= -\frac{(\eta_1 - \eta_2)^2}{(\sigma_1^2 + \sigma_2^2)^2} \left(\frac{\partial \sigma_1^2}{\partial C} + \frac{\partial \sigma_2^2}{\partial C} \right) + \frac{2(\eta_1 - \eta_2)}{\sigma_1^2 + \sigma_2^2} \left(\frac{\partial \eta_1}{\partial C} - \frac{\partial \eta_2}{\partial C} \right) = 0$$

$$\text{和} \quad -\frac{(\eta_1 - \eta_2)^2}{(\sigma_1^2 + \sigma_2^2)^2} \left(\frac{\partial \sigma_1^2}{\partial c_0} + \frac{\partial \sigma_2^2}{\partial c_0} \right) + \frac{2(\eta_1 - \eta_2)}{\sigma_1^2 + \sigma_2^2} \left(\frac{\partial \eta_1}{\partial c_0} - \frac{\partial \eta_2}{\partial c_0} \right) = 0$$

由式(2), (3)得到

$$\frac{\partial \eta_i}{\partial C} = \bar{X}_i, \quad \frac{\partial \eta_i}{\partial c_0} = 1, \quad \frac{\partial \sigma_i^2}{\partial C} = 2S_i C, \quad \frac{\partial \sigma_i^2}{\partial c_0} = 0$$

代入上式有

$$\frac{\eta_1 - \eta_2}{\sigma_1^2 + \sigma_2^2} (S_1 + S_2) C = \bar{X}_1 - \bar{X}_2$$

即

$$C = \frac{\sigma_1^2 + \sigma_2^2}{\eta_1 - \eta_2} (S_1 + S_2)^{-1} (\bar{X}_1 - \bar{X}_2)$$

其中, $(\sigma_1^2 + \sigma_2^2)/(\eta_1 - \eta_2)$ 为常数因子, 在 c_0 的表达式中也会出现同样的因子, 因此可约去, 则得到

$$C = (S_1 + S_2)^{-1} (\bar{X}_1 - \bar{X}_2) \quad (5)$$

c_0 可以这样计算: 把两类样品的共用均值

$$\bar{X} = (N_1 \bar{X}_1 + N_2 \bar{X}_2) / (N_1 + N_2) \quad (6)$$

代入 $h(X)$, 并使它取零值作为两类的分界, 即

$$h(\bar{X}) = C^T \bar{X} + c_0 = 0$$

$$\text{则得} \quad c_0 = -C^T \bar{X} = -(N_1 C^T \bar{X}_1 + N_2 C^T \bar{X}_2) / (N_1 + N_2) \quad (7)$$

2.2.3 判别准则

上面已求得了 Fisher 线性分类器 $h(X) = C^T X + c_0$. 则 Fisher 方法的分类判别准则为

$$\begin{cases} h(X) = C^T X + c_0 \geq 0 \rightarrow X \in \text{第1类} \\ h(X) = C^T X + c_0 < 0 \rightarrow X \in \text{第2类} \end{cases} \quad (8)$$

2.3 ICHAM 方法

ICHAM 方法是由王碧泉等人提出的, 它是改进的连续 Hamming 方法的简称. 有关方法参看王碧泉和王春珍(1988)的文章.

2.4 BEG 方法

BEG 方法的全名为基本事件形成法 (Basic Event Generation). 它适用于样品成多个团状分布的情形, 且可用于多类的识别. 这一算法的要点是形成一此“基本事件”, 然后依某样品是否“落入”基本事件来决定此样品的类别.

2.4.1 事件与基本事件

设样品分成 M 类, 以 c_i 表示样品的类别, $i=1, 2, \dots, M$, 且 $M \geq 2$. 每一样品以 n 个特征描述, $X = (x_1, x_2, \dots, x_n)^T$, 即每一样品是 n 维特征空间 R^n 中的一个点

$$R^n = \overbrace{R_1 \times R \times \dots \times R}^n$$

定义 n 维空间中的一个事件 E 为一超矩形

$$E = I_1 \times I_2 \times \cdots \times I_n$$

$$I_k \text{ 为一闭区间 } [a_k, b_k] \subset R^1$$

一个事件中至少包含一个样品. 若一个事件只包含一个样品时称为退化事件.

当某一事件包含了 c_i 类中的某些样品而不包含除 c_i 类以外的其它类的任何样品时, 则称这一事件为 c_i 类的一个基本事件. c_i 类的全部基本事件构成 c_i 类的基本事件集 E .

2.4.2 合并函数和基本事件的形成

合并函数 $M(E_1, E_2)$ 为任意一对事件 E_1 和 E_2 的组合

$$M(E_1, E_2) = E (\subset R^n)$$

其中

$$E_1 = I_{11} \times I_{12} \times \cdots \times I_{1n} \quad E_2 = I_{21} \times I_{22} \times \cdots \times I_{2n}$$

$$I_{ij} = [a_{ij}, b_{ij}] \quad i = 1, 2 \quad j = 1, 2, \cdots, n$$

$$E = I_1 \times I_2 \times \cdots \times I_n$$

$$I_k = [\min(a_{1k}, a_{2k}), \max(b_{1k}, b_{2k})] \quad k = 1, 2, \cdots, n$$

也就是说, 合并函数 $M(E_1, E_2)$ 是包含两事件 E_1, E_2 的最小超矩形, 它也是一个事件.

下面定义某样品与事件 E 的距离为

$$d(X|E) = \sum_{i=1}^n \varphi(x_i|I_i)$$

$$\varphi(x_i|I_i) = \begin{cases} 1 & \text{若 } x_i \notin I_i \\ 0 & \text{若 } x_i \in I_i \end{cases}$$

BEG 算法是通过合并过程来形成基本事件的. 其合并原则如下: 设 X_1, \cdots, X_{N_i} 是 c_i 类的训练样品, E 为 c_i 类的某个事件 (由 c_i 类样品形成的某个超矩形). $Y_1, Y_2, \cdots, Y_{N_j}$ 是除 c_i 类以外的其它各类的训练样品. 若选定某一正整数 T_A 为阈值, 当 c_i 类中的某个样品 X 满足

$$d[Y_k|M(E, X)] \geq T_A \quad k = 1, 2, \cdots, N_j$$

则认为 X 与事件 E 可以合并. 也就是说, 样品 X 与事件 E 先形成合并函数 $M(E, X)$, 它是包括 X 与 E 的最小超矩形. 若此超矩形与 c_i 类以外其它类的任一样品 Y_k 的距离都大于或等于 T_A 时, 则认为 X 与 E 可以合并.

2.4.3 判别准则

对于每类 $c_i (i=1, 2, \cdots, M)$ 都可按上述算法形成基本事件集 $e = E_{i1} \cup E_{i2} \cup \cdots \cup E_{iN_i}$. 其中, E_{ik} 为 c_i 类的第 k 个基本事件, $k=1, 2, \cdots, N_i$. 对未知类别的样品 X , 可视其与某类基本事件的距离决定它是否“落入”该类基本事件而判定这一样品的类别.

取定一非负整数 T_c 为阈值, 则判定样品 X 的类别的准则为: 若仅存在一个 k , 使满足

$$d(X|E_{ik}) \leq T_c$$

则判定 X 属于 c_i 类. 一般取 $T_c < T_A$. 当 $T_c = 0$ 时, 样品 X 在通常意义下落入超矩形. 当 $T_c > 0$ 时, 则为广义的“落入”.

2.5 几何解释

对于两类问题, 并假设特征数为 $n=2$. 图 1 中示出了上述几种模式识别方法的几何解释. 图中直线、圆或折线分别表示各方法的分类器. 由图可见, 最小距离法得到的两类

分界线是一条直线, 即两类样品均值连线的垂直平分线. Fisher 方法的分界线也是一条直线. ICHAM 方法的分界线是一个圆, 它以某一类的 Hamming 核为圆心. BEG 方法得到的分界线可以认为是折线, 但其中每条直线都平行于坐标轴.

不难看出, 最小距离法、Fisher 方法对可分离的两类样品, 其分类效果都将较好, 但最小距离法简单方便, Fisher 方法用式(4)式的原则求分类器其分类效果会更佳些. ICHAM 方法对于一类包含于另一类中的样品或两类交叉分布的样品其分类效果会较好. BEG 方法则适于样品成多个团状分布或彼此交叉的样品分布.

关于以上几种模式识别方法的详细叙述, 还可参看有关文献(王碧泉、王春珍, 1988; 王碧泉、陈祖荫, 1989b).

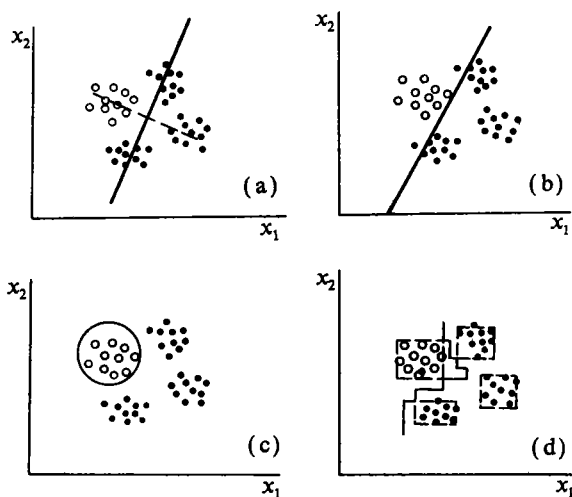


图 1 几种模式识别方法的几何解释. (a)最小距离法, (b)Fisher 方法, (c)ICHAM 方法, (d)BEG 方法

3 PTR 算法

在黄汉明和王碧泉(1993)的文章中, 提出了两种平面映射方法, 可将多维空间中的样品投影到平面上; 还提出了样品分布信息函数及 4 种典型样品; 并研究了典型样品和本文所述 4 种模式识别方法之间的关系, 得到各种典型样品的最佳分类识别方法. 其结果与本文下节对各方法识别效果的分析是完全一致的.

在该文和本文的基础上, 我们建立了一套算法. 它的核心是样品平面投影、典型样品分析和模式识别, 该方法简称为 PTR 算法. 此算法的要点是:

3.1 样品的平面投影

在“一次映射”和“逐次映射”(黄汉明、王碧泉, 1993)两种投影方法中任选一种, 将 n 维空间中的样品投影到平面上, 得到反映样品结构的样品分布图形, 便于直观地认识样品结构.

3.2 典型样品分析

按黄汉明和王碧泉(1993)提出的方法, 计算 A, B 两类样品的样品分布信息函数 $\mu(B|A)$ 和 $\mu(A|B)$, 并按下列准则确定所研究的样品属于哪一种典型样品:

(1) 若 $\mu(B|A) = \mu(A|B) = 0$, 或 $\mu(B|A) \leq \mu_1^T$ (如 0.35) 和 $\mu(A|B) \leq \mu_1^T$, 且 $S^2(A) \gg S^2(B)$ (或 $S^2(A) < S^2(B)$), 则为第 I 种典型样品.

(2) 若 $\mu(B|A) = 0$, 且 $\mu(A|B) = 1$, 或 $\mu(B|A) \leq \mu_1^T$ (如 0.35), 且 $\mu(A|B) \geq \mu_2^T$ (如 0.65) 时, (或 $\mu(A|B) \leq \mu_1^T$, 且 $\mu(B|A) \geq \mu_2^T$), 则为第 II 种典型样品.

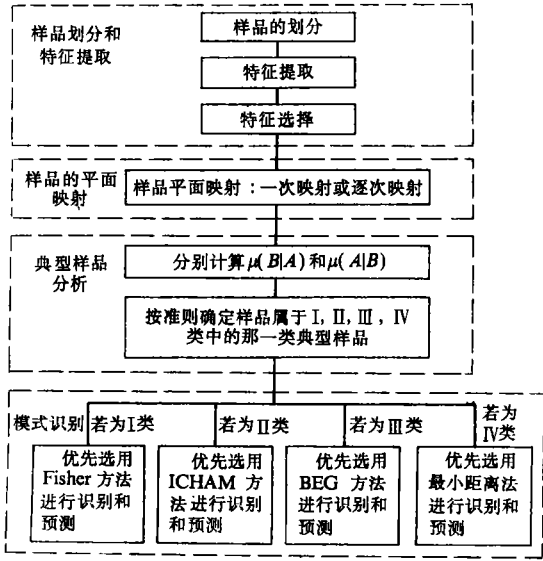


图 2 PTR 算法流程图

(3) 对第Ⅲ种典型样品, 选用 BEG 方法;

(4) 对第Ⅳ种典型样品, 选用最小距离法. PTR 算法的流程如图 2 所示.

4 PTR 算法的应用实例

4.1 实例一

4.1.1 样品的确定

取华北地区(107.5°E—120.5°E, 29.0°N—44.0°N)在 1400—1983 年期间发生的 $M \geq 7$ 的地震作为“有强震”, 即作为 D 类样品, 共有 $N_1=16$ 个(表 1). N 类样品则是按随机原则抽取的. 在上述时间、地域范围内, 形成一系列年、月、经度、纬度的随机数, 作为样品的时间和位置. 若在此时间、地点确未发生 $M \geq 7$ 的地震, 则得到一个“无强震”(即 N 类样品). 考虑到每次强震有一定的孕震时间和孕震区以及强震的震后效应, 去掉了在“有强震”前后 40 年且与该强震的距离小于 600 km 的那些“无强震”, 从而得到 $N_2=19$ 个 N 类样品(表 1). 样品总数 $N=35$.

4.1.2 特征提取

对地震活动水平、 b 值、空区、地震条带分布、地震的集中性等地震活动前兆, 分别建立了函数 N, b, g, BE, C . 根据对震例(D 类样品及 N 类样品)的分析, 对各函数选取参数, 提出了确定前兆的 10 条准则:

准则 1: $N(t, \varphi, \lambda|300, 5, 4.75, 14, 3) \geq 2$ (简记为 $N_1 \geq 2$).

准则 2: $N(t, \varphi, \lambda|300, 5, 4.75, 32, 3) \geq 2$ (简记为 $N_2 \geq 2$).

准则 3: $N(t, \varphi, \lambda|300, 5, 5.0, 12, 1) \geq 2$ (简记为 $N_3 \geq 2$).

准则 4: $N(t, \varphi, \lambda|300, 5, 5.0, 32, 3) \geq 2$ (简记为 $N_4 \geq 2$).

准则 5: $N(t, \varphi, \lambda|300, 5, 5.5, 21, 4) \geq 1$ (简记为 $N_5 \geq 1$).

(3) 若 $\mu(B|A)=1$, 且 $\mu(A|B)=1$, 或 $\mu(B|A) \geq \mu_2^2$ (如 0.65) 且 $\mu(A|B) \geq \mu_2^2$, 则为第Ⅲ种典型样品.

(4) 若 $\mu(B|A)=\mu(A|B)=0$, 或 $\mu(B|A) \leq \mu_1^2$ (如 0.35), 且 $\mu(A|B) \leq \mu_1^2$, 且 $|S^2(A)-S^2(B)| \leq T^s$, 则为第Ⅳ种典型样品.

3.3 模式识别

按下列原则选择最适于该样品结构的模式识别方法, 对样品进行分类识别(对未知类别样品的识别即为预测), 以得到最好的结果.

(1) 对第Ⅰ种典型样品, 选用 Fisher 方法;

(2) 对第Ⅱ种典型样品, 选用 ICHAM 方法;

准则 6: $b(t, \varphi, \lambda|300, 2, 4.75, 8, 1) \geq 1.5$ (简记为 $b_1 \geq 1.5$).

准则 7: $b(t, \varphi, \lambda|300, 2, 4.75, 15, 2) \leq 0.7$ (简记为 $b_2 \leq 0.7$).

准则 8: $g(t, \varphi, \lambda|4.75, 30, e, 200) \geq 2/3$ (简记为 $g \geq 2/3$).

准则 9: $BE(t, \varphi, \lambda|4.75, 20, 360, NE45^\circ) \geq 1000$ (简记为 $BE \geq 1000$).

准则 10: $C(t, \varphi, \lambda|550, 4.75, 30) \geq 8$ (简记为 $C \geq 8$).

每条准则是描述样品的一个特征,共有 $n=10$ 个特征.若某样品的地震活动满足某准则,则该准则所对应的样品特征取值为 1,否则取值为 0.关于此实例的样品和特征提取的详情参看王碧泉等(1986)的文章.

表 1 实例一的样品参数表

强震目录(D类对象的参数)						“无强震”(N类对象的参数)			
序号	发震时间	震中位置		M	地点	序号	时间	位置	
	年-月-日	$\varphi(^{\circ}N)$	$\lambda(^{\circ}E)$				年-月	$\varphi(^{\circ}N)$	$\lambda(^{\circ}E)$
1	1501-01-19	34.8	110.1	7.0	朝 邑	1	1454-01	38.80	111.01
2	1548-09-13	38.0	121.0	7.0	渤 海	2	1460-09	37.78	117.31
3	1556-01-23	34.5	109.7	8.0	华 县	3	1487-05	38.45	119.86
4	1597-10-06	38.5	120.0	7.0	渤 海	4	1526-09	30.60	118.75
5	1626-06-28	39.4	114.2	7.0	灵 丘	5	1556-07	39.72	122.81
6	1668-07-25	35.3	118.6	8.5	郟 城	6	1623-11	33.54	110.14
7	1679-09-02	40.0	117.0	8.0	三 河	7	1646-05	29.37	121.17
8	1683-11-22	38.7	112.7	7.0	原 平	8	1746-01	36.21	119.93
9	1695-05-18	36.0	111.5	8.0	临 汾	9	1746-09	29.62	115.35
10	1830-06-12	36.4	114.2	7.5	磁 县	10	1767-01	32.20	111.37
11	1888-06-13	38.5	119.0	7.5	渤海湾	11	1782-05	35.39	116.63
12	1937-08-01	35.4	115.1	7.0	荷 泽	12	1795-10	40.09	120.40
13	1966-03-08	37.5	115.1	7.2	邢 台	13	1800-10	39.28	121.58
14	1969-07-18	38.2	119.4	7.4	渤 海	14	1862-04	12.11	111.27
15	1975-02-04	40.7	122.7	7.3	海 城	15	1871-12	32.76	122.63
16	1976-07-28	39.4	118.0	7.8	唐 山	16	1881-07	43.34	115.64
						17	1913-12	32.21	110.61
						18	1927-07	42.93	113.28
						19	1950-01	32.81	119.58

注:“无强震”指在该时间、地点上未发生 $M \geq 7$ 的地震.

表 2 样品分布信息函数与误识率(实例一)

样 品	样品分布 信息函数		样品 类别	Fisher 方法			ICHAM 方法			最小距离方法		
	$\mu(D N)$	$\mu(N D)$		错 分	总 误	错 分	总误识率	错 分	总 误			
				样品数	误识率	识 率	样 品数	误识率	样品数	误识率	识 率	
原样品	0.38	0.10	D	1	0.06	0.11	1	0.06	0.14	3	0.19	0.14
			N	3	0.16		4	0.21		2	0.11	
第Ⅰ种 方 法 投影后	0	0	D	0	0	0	0	0	0.03	0	0	0.03
			N	0	0		1	0.05		1	0.05	
第Ⅱ种 方 法 投影后	0.06	0.10	D	1	0.06	0.09	1	0.06	0.09	1	0.06	0.09
			N	2	0.11		2	0.11		2	0.11	

4.1.3 用 PTR 算法的识别结果

对此实例应用 PTR 算法进行分析计算：(1) 用一次映射方法投影后的样品平面分布如图 3 所示, D , N 两类分得较开；(2) 计算样品分布信息函数, 得到 $\mu(D|N)=\mu(N|D)=0$ (表 2), 可知此组样品为我们定义的第 I 种典型样品；(3) 按 PTR 算法, 对第 I 种典型样品, 最好选用 Fisher 方法. 我们采用 Fisher 方法对这组样品进行识别, 得到 D 类和 N 类

样品的错分类数目均为零, 因此总误识率亦为零, 其分类效果确实很好, 详见表 2.

为比较, 对原样品(未作映射)及用两种方法分别投影后的样品, 分别用 Fisher 方法、ICHAM 方法及最小距离法进行了模式识别, 其结果列于表 2 中. 由表可见:

1) 各种结果中, 以 PTR 算法所得结果最好(误识率最小), 即用一次映射后, 按 PTR 方法确定为第 I 种典型样品后用 Fisher 方法识别的结果最好. ICHAM 方法和最小距离的结果次之.

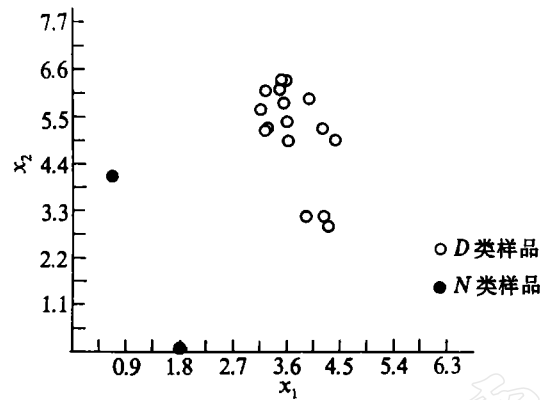


图 3 实例一的样品用一次映射方法投影到平面后的样品分布图

表 3 实例二的 D 类样品参数表

样本 序号	发 震 时 间		震 中 位 置		震 级	参 考 地 名
	年-月-日	时	$\varphi(^{\circ}\text{N})$	$\lambda(^{\circ}\text{E})$		
1	1966-03-08	05	37.35	114.92	7.0	河北宁晋
2	1967-03-27	16	38.57	116.62	6.5	河北河间大城
3	1967-07-28	13	40.65	115.77	5.8	河北怀来东北
4	1967-12-18	22	36.50	110.83	5.7	山西蒲县
5	1969-07-18	13	38.20	119.40	7.5	渤海湾
6	1970-11-10	20	40.75	109.42	5.5	内蒙乌拉特前旗东
7	1974-04-22	08	31.43	119.28	5.8	江苏溧阳
8	1975-02-04	19	40.67	122.83	7.4	辽宁海城
9	1975-09-02	20	32.83	121.83	5.7	黄海郎家沙
10	1976-04-06	00	40.28	112.18	6.4	内蒙和林格尔
11	1976-07-28	03	39.63	118.18	7.9	河北唐山
12	1976-10-06	09	35.33	124.23	5.8	山东青岛以东海中
13	1979-05-22	06	31.10	110.47	5.5	湖北秭归
14	1979-06-19	12	37.10	111.90	5.6	山西介休
15	1979-07-09	18	31.42	119.22	6.3	江苏溧阳
16	1979-08-25	00	41.17	108.07	6.3	内蒙五原
17	1981-08-13	11	40.58	113.40	5.9	内蒙丰镇
18	1983-11-07	05	35.20	115.30	6.4	山东菏泽
19	1984-05-21	23	32.70	121.60	6.5	黄海
20	1987-02-17	11	33.70	120.70	5.6	黄海

2) 对应用 Fisher 方法而言, 用一次映射后结果最好, 用逐次映射后的结果次之, 原样品的结果最差. 可见, 映射可以改进识别的效果.

4.2 实例二

4.2.1 样品的确定

样品取法类似于实例一, 只是范围为华北地区 (107.5°E — 125.0°E , 29.0°N — 45.0°N), 时间为 1966—1987 年, 强震取 $M \geq 5.5$ 的地震 (D 类样品). N 类样品仍为随机抽取, 只是去掉了强震前后 5 年内且与其距离小于 300 km 的 N 类样品和相互之间的时间小于 4 年且距离小于 200 km 的 N 类样品. 如此得到 D 类样品 $N_1 = 20$ 个 (表 3), N 类样品 $N_2 = 25$ 个 (其参数不一示出), 总样品数 $N = 45$.

4.2.2 特征提取

类似地, 从地震活动水平、 b 值、空区、地震的条带分布、地震的集中性等地震活动前兆中提取了 10 个特征, 这 10 个特征的取法和含意与实例一不同, 而且, 除特征 8 和 9 外其它特征均为连续值. 有关样品和特征的详情可参看王碧泉等 (1989a) 的文章. 各特征的含意简述如下:

- (1) x_1 , 二级地震的最高频度, 记为 $NM2H$.
- (2) x_2 , 二级地震的频度增强, 记为 $DNM2$.
- (3) x_3 , 三级地震的最高频度, 记为 $NM3H$.
- (4) x_4 , 三级地震的频度增强, 记为 $DNM3$.
- (5) x_5 , 四级地震的频度增强, 记为 $DNM4$.
- (6) x_6 , 高 b 值, 记为 BH .
- (7) x_7 , b 值下降, 记为 DB .
- (8) x_8 , 地震空区, 记为 GAP .
- (9) x_9 , 地震的条带分布, 记为 $BELT$.
- (10) x_{10} , 地震的集中性, 记为 CON .

4.2.3 用 PTR 算法的识别结果

(1) 采用一次映射方法投影后的样品平面分布如图 4 所示. 可见 D 、 N 两类有一定程度混杂.

(2) 得到样品分布信息函数 $\mu(D|N) = 0.70$, $\mu(N|D) = 0.04$ (表 4), 表明 N 类中只有少量包含于 D 类. 可知此组样品为我们定义的第 II 种典型样品.

(3) 按 PTR 算法, 采用 ICHAM 方法识别, 得到 D 、 N 类错分类的样品数分别为 0 和 1, 总误识率为 0.02, 分类效果很好, 详见表 4.

为比较, 亦计算了两种映射后及原样品用 3 种方法的识别结果 (表 4), 由表可见: 1)

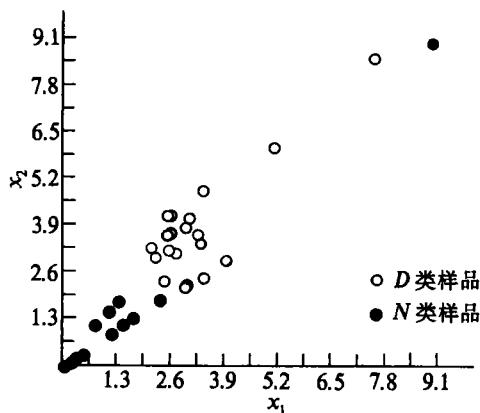


图 4 实例二的样品用一次映射方法投影到平面后的样品分布图

各结果中,以按 PTR 方法选取 ICHAM 方法进行识别的效果最好,其次才是最小距离法和 Fisher 方法;2) 用一次映射后的结果最好,用逐次映射后与原样品的结果次之. 以上结论与实例一是一致的.

表 4 样品分布信息函数与误识率(实例二)

样 品	样品分布 信息函数		样品 类别	Fisher 方法			ICHAM 方法			最小距离方法		
	$\mu(D N)$	$\mu(N D)$		错 分 样品数	误识率	总 误 识 率	错 分 样品数	误识率	总误识率	错 分 样品数	误识率	总 误 识 率
原样品	0.90	0.44	D	10	0.50	0.33	7	0.35	0.33	9	0.45	0.31
			N	5	0.20		8	0.32		5	0.20	
第 I 种 方 法 投影后	0.70	0.04	D	6	0.30	0.16	0	0	0.02	4	0.20	0.11
			N	1	0.04		1	0.04		1	0.04	
第 II 种 方 法 投影后	1	0.48	D	11	0.55	0.33	4	0.20	0.33	9	0.45	0.29
			N	4	0.16		11	0.44		4	0.16	

5 结 论

本文详细阐述了 4 种模式识别方法. 在样品平面映射、典型样品分析和本文所述模式识别方法的基础上,我们提出了 PTR 算法,该算法可以对不同结构的样品采用不同的模式识别方法进行分类识别,既方便、较快,又可改进分类效果. 两个应用实例表明,PTR 算法的效果较好. 此外还可看到,平面映射可以改善识别结果.

参 考 文 献

黄汉明、王碧泉,1993. 样品平面映射方法及典型样品研究. 地震学报, 15,增刊.

王碧泉、马秀芳、杨锦英、王春珍,1986. 应用模式识别方法综合分析多项震兆. 地震研究, 9, 643—657.

王碧泉、王春珍,1988. 对连续 Hamming 方法的改进及其在潜在震源判定中的应用. 地震学报,10, 113—123.

王碧泉、范洪顺、杨锦英、王春珍、陈佩燕,1989a. 模式识别方法应用于测震前兆的综合预测. 地震预报方法实用化研究文集,测震学专辑,514—526. 学术书刊出版社,北京.

王碧泉、陈祖荫,1989b. 模式识别——理论、方法和应用, 375pp. 地震出版社,北京.

Keilis-Borok, V. I. and Rotwain, I. M., 1990a. Diagnosis of time of increased Probability of Strong earthquakes in different regions of the world: Algorithm CN. *Phys. Earth. Planet. Inter.*, 61, 57—72.

Keilis-Borok, V. I. and Kossolokov, V. G., 1990b. Premonitory activation of earthquake flow: Algorithm M8. *Phys. Earth. Planet. Inter.*, 61, 73—83.