

# 多元统计方法

方开泰 编

中国科学院应用数学研究所  
概率统计咨询服务部 印

1984.7

## 编 者 的 话

近年来，许多部门经常要求我们举办多元分析方法的学习讨班，为了给这一类活动提供适当的教材，特地选编了部分常用的方法，其中包括：多变量样本的图分析法，聚类分析，回归分析，判别分析，主分量分析和因子分析。由于这些材料来源不同，符号体例均无法统一，希读者谅解。

方开泰

1984·6·18

## 目 录

- 1、多变量样本的图分析法(一) 方开泰..... 1
- 2、多变量样本的图分析法(二) 方开泰..... 19
- 3、聚类分析(I) 方开泰..... 36
- 4、聚类分析(II) 方开泰..... 68
- 5、回归分析简介 方开泰..... 88
- 6、判别分析 方开泰..... 161
- 7、主分量 M·肯德尔..... 255
- 8、因子分析 M·肯德尔..... 281

## 多变量样本的图分析法(一)

中国科学院应用数学研究所 方开泰

图形是帮助人们思维和判断的重要工具,当样本只有两个特性(变量或指标)时,可以用通常的直角坐标在平面上点图,当样本有三个变量时,虽然可以在三维的笛卡儿坐标里点图,但也是很很不方便的。当变量数大于三时,用通常的方法已不能点图了。在多元分析中,样本的变量数一般均大于三,探讨多变量的点图法是长期来一直为人们所关注的研究课题,这里介绍一些有关的方法,特别是近十年来发展的一些方法。

作者非常感谢刘璋温副教授,因为他提供了许多有用的文献,并仔细校阅和修改了此文。

### 一、雷达图(radar chart)

设有  $n$  个样本  $X_1, X_2, \dots, X_n$ , 各有  $p$  个变量  $x_1, x_2, \dots, x_p$ , 用  $x_{ij}$  ( $i=1, 2, \dots, n; j=1, 2, \dots, p$ ) 表示第  $i$  个样本的第  $j$  个变量的值, 通常数据  $\{x_{ij}\}$  列成表 1.1 的形式。

表 1.1 数据矩阵

变 量 样本	$x_1$	$x_2$	...	$x_j$	...	$x_p$
$X_1$	$x_{11}$	$x_{12}$	...	$x_{1j}$	...	$x_{1p}$
$X_2$	$x_{21}$	$x_{22}$	...	$x_{2j}$	...	$x_{2p}$
⋮	⋮	⋮	...	⋮	...	⋮
$X_i$	$x_{i1}$	$x_{i2}$	...	$x_{ij}$	...	$x_{ip}$
⋮	⋮	⋮	...	⋮	...	⋮
$X_n$	$x_{n1}$	$x_{n2}$	...	$x_{nj}$	...	$x_{np}$

雷达图也称蜘蛛网图，其作图步骤如下：划一个圆，由  $P$  个点把圆周分成相等的  $P$  个部分，由圆心连接这  $P$  个点，得  $P$  个辐射状的半径，这  $P$  个半径就作为  $P$  个变量的坐标轴。根据各变量波动的大小，将相应的坐标轴标上刻度。将每个变量的值点在各自的坐标轴上，依次连接起来得一  $P$  边形。对每个样本都这样做，可得  $n$  个  $P$  边形，形成一个雷达图。

例 1、表 1.2 列出了邓阜仙岩体的部分化学成分，试用雷达图来表示它们。

运用上述作雷达图的步骤得图 1.1，从图上看到： $X_1$  与  $X_2$ ； $X_3, X_4$  与  $X_5$ ； $X_6$  与  $X_7$  图形分别相似，可以分成三类，这正好是三种类型的花岗岩。

表 1.2 邓阜仙岩体化学成分(部分)

岩 体	SiO <sub>2</sub>	TiO <sub>2</sub>	FeO	CaO	K <sub>2</sub> O
X <sub>1</sub> : 中粗粒斑状黑云母花岗岩	75.20	0.14	1.86	0.91	5.21
X <sub>2</sub> : 中粗粒斑状黑云母花岗岩	75.15	0.16	2.11	0.74	4.93
X <sub>3</sub> : 中粒二云母花岗岩	72.19	0.13	1.52	0.69	4.65
X <sub>4</sub> : 中粒二云母花岗岩	72.35	0.13	1.37	0.83	4.87
X <sub>5</sub> : 中粒二云母花岗岩	72.74	0.10	1.41	0.72	4.99
X <sub>6</sub> : 细粒白云母花岗岩	73.29	0.033	1.07	0.17	3.15
X <sub>7</sub> : 细粒白云母花岗岩	73.72	0.033	0.77	0.28	2.78

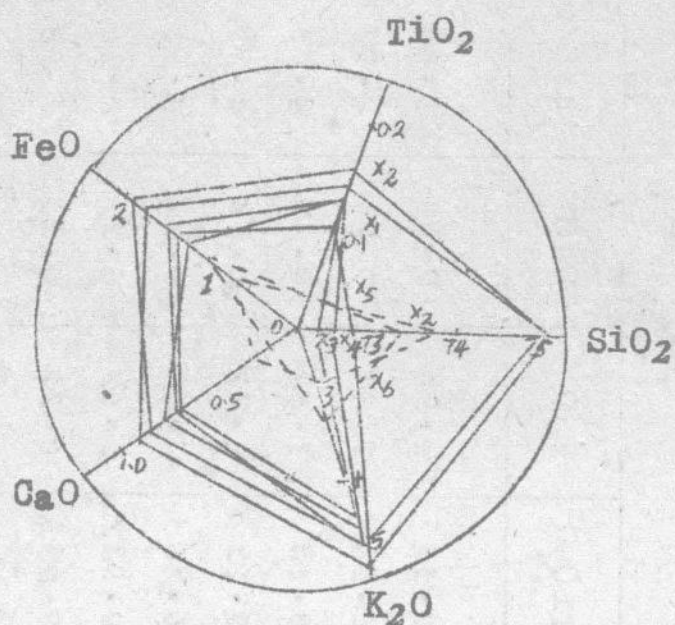


图1.1 雷达图

在雷达图上，每个样本对应一个P边形。当P较大或样本较多时，图上线段太多，看起来很不方便，这时可以每个样本画一张图，如图1.2。如果样本的分类已很清楚，这时可以将同类的计算均值，用每类的均值来作雷达图，这时图像既清晰，每类特征又很明显。例如对例1，我们求得三类的均值列表1.3，用三组均值作成雷达图1.3，我们看到图像清晰，三类特点一目了然。

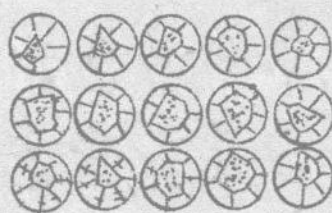


图1.2

表 1.3 类均值

类	SiO <sub>2</sub>	TiO <sub>2</sub>	FeO	CaO	K <sub>2</sub> O
I: {X <sub>1</sub> , X <sub>2</sub> }	75.18	0.15	1.99	0.83	5.07
II: {X <sub>3</sub> , X <sub>4</sub> , X <sub>5</sub> }	72.43	0.12	1.43	0.75	4.84
III: {X <sub>6</sub> , X <sub>7</sub> }	73.51	0.033	0.92	0.23	2.97

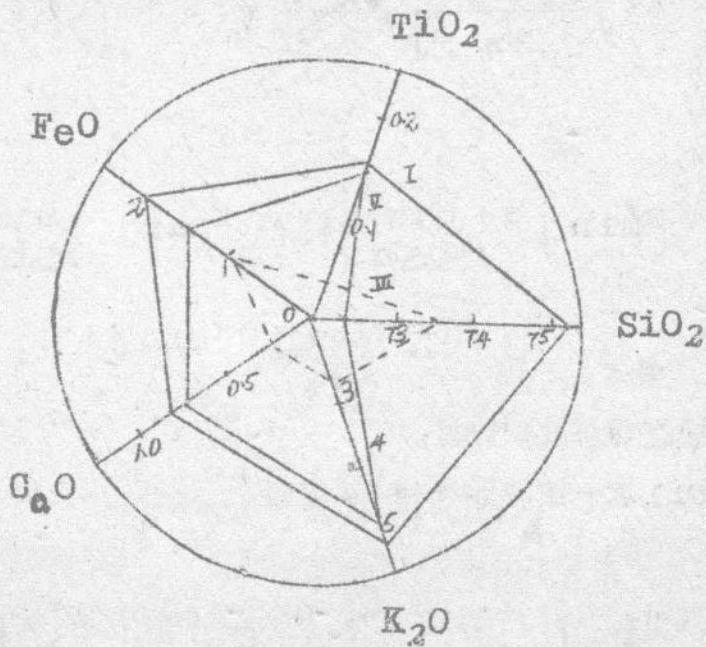


图 1.3 均值的雷达图

## 二、星座图 (constellation graph)

上节我们看到, 当  $n$  很大时, 用雷达图来表达每一个样本是不方便的, 这时可以用星座图来描述。所谓星座图, 就是将样本



点在一个半圆内，一个样本用一颗星来表示，同类的样本组成一个星座，不同类样本组成不同的星座，很象天文学上表示星座的图象。

星座图的作图步骤如下：

(i) 先将数据  $\{x_{ij}\}$  作一线性变换，使变换后的数据  $\{\varepsilon_{ij}\}$  落到  $(0, \pi)$  内，常取

$$\varepsilon_{ij} = \frac{x_{ij} - x_{\min j}}{R_j} \pi, \quad (1)$$

其中

$$x_{\min j} = \min_{1 \leq i \leq n} x_{ij}, \quad x_{\max j} = \max_{1 \leq i \leq n} x_{ij},$$

$$R_j = x_{\max j} - x_{\min j},$$

这就是极差标准化的办法。

(ii) 取一组权  $\{\omega_j\}$ ，使得

$$\omega_j \geq 0, \quad j=1, 2, \dots, p, \quad \sum_{j=1}^p \omega_j = 1.$$

重要的变量相应的权可取得大一点，或者也可以取等权，即

$$\omega_1 = \omega_2 = \dots = \omega_p = 1/p.$$

(iii) 画一个半径为 1 的上半圆及半圆底边的直径，使每个样本对应一个星，这些星就落在这个半圆内。比如为求  $X_1$  对应的星，先以圆心  $O$  为圆心， $\omega_1$  为半径划一上半圆，在这半圆周上对应弧度为  $\varepsilon_{11}$  的点记作  $O_1$ ，接着以  $O_1$  为心， $\omega_2$  为半

径作一上半圆，在这半圆周上对应弧度为  $\varepsilon_{12}$  的点记作  $O_2$ ；再以  $O_2$  为心、 $\omega_3$  为半径作一上半圆，取弧度为  $\varepsilon_{13}$  的点，记为  $O_3$ ；…依次一直求到  $O_p$ ，这最后的  $O_p$  就是  $X_1$  所对应的星的位置，记作  $z_1$ 。由于  $\sum \omega_j = 1$ ， $\omega_j \geq 0$ ，星必然落在半圆内。由  $O$  点通过上述作图步骤，到达星的路径（即  $OO_1O_2 \cdots O_p$ ）称做该星的路径，通过星的位置和路径就全面的刻划了该样本的特征。对其余样本作图是类似的，详见图 2.1。

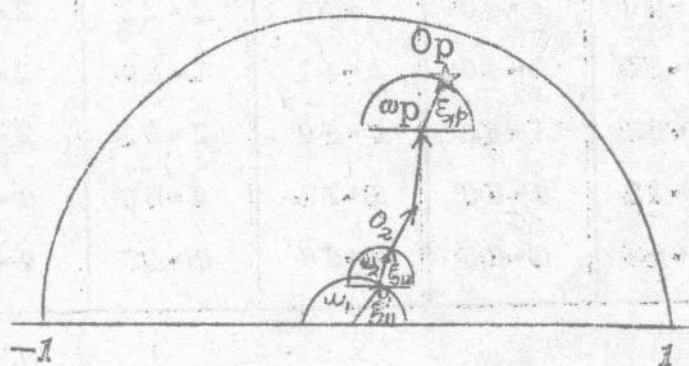


图 2.1 星座图

通过作图来求路径和星的位置是不方便的，最好借助于计算来实现。不难看出  $X_a$  对应的  $O_1, O_2, \dots, O_{p-1}, O_p$  的坐标为

$$\left( \sum_{j=1}^l \omega_j \cos \varepsilon_{aj}, \sum_{j=1}^l \omega_j \sin \varepsilon_{aj} \right), \quad l=1, 2, \dots, p.$$

决定了  $O_1, O_2, \dots, O_p$  的位置后，不难用计算机来打图。

现在我们仍用例 1 来说明星座图的应用。由表 1.2 算得  $\{\varepsilon_{ij}\}$  列于表 2.1，并取

$$\omega_1 = \omega_2 = \omega_3 = \omega_4 = \omega_5 = 1/5,$$

由此作成图 2·2，我们看到组成了三个星座，对应于三类岩体，似乎比雷达图表达得更为直观和醒目。

表 2·1  $\varepsilon_{ij}$  值

	SiO <sub>2</sub>	TiO <sub>2</sub>	FeO	CaO	K <sub>2</sub> O
X <sub>1</sub>	3·14	2·65	2·56	3·14	3·14
X <sub>2</sub>	3·09	3·14	3·14	2·42	2·78
X <sub>3</sub>	0·00	2·40	1·76	2·21	2·42
X <sub>4</sub>	0·17	2·40	1·41	2·80	2·70
X <sub>5</sub>	0·57	1·66	1·50	2·33	2·86
X <sub>6</sub>	1·15	0·00	0·70	0·00	0·48
X <sub>7</sub>	1·60	0·00	0·00	0·47	0·00

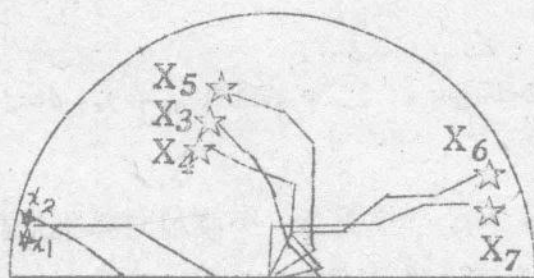


图 2·2

当各类分得很清楚(如图2·2)时,只需画星的位置不一定画相应的路径,使图更加清晰明了。但有时也会出现异途同归的情况,星的路径很不一样,而星的位置很接近,这时如在图中去掉路径,就会发生误解。为了能够去掉路径,最好适当调整 $\{\omega_i\}$ ,使得异途不能同归。

星座图最初是由Wakimoto和Taguri<sup>[3]</sup>提出来的,并在医学、行为科学、保险公司、学校教育等方面取得不少应用。笔者在国内介绍后,已在地质、气象、考古等方面得到了应用。如下例是关于考古学的。

例2、为了研究人类的进化,测了现代男人和古代雄猿的下颚和牙齿的十七个部位,有关数据来自[4],作成星座图2·3(图中·表示一颗星),其中编号1-20的为人类,21-38的为猿,人全落在半圆的右边,猿在左边,两类泾渭分明。

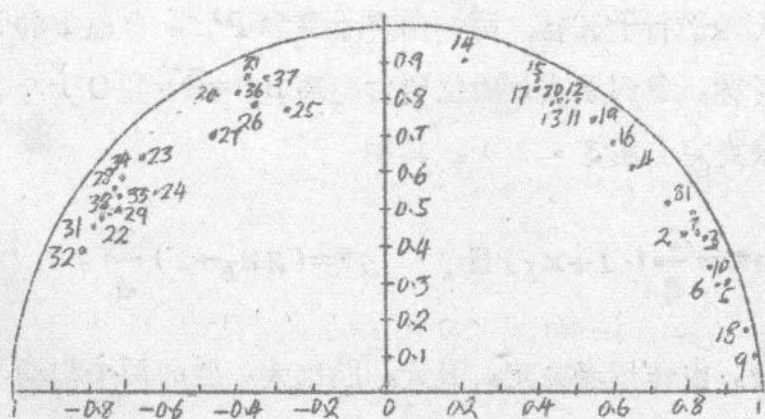


图2·3

### 三、脸谱图 (face graph)

脸谱是一个非常形象的东西，脸谱的胖瘦、喜怒哀乐给人留下深刻的印象，京剧就在脸谱的设计上下了很大功夫。用脸谱来表达样本首先是由统计学家 Chernoff<sup>[1]</sup> 提出来的，他将样本的特性用人脸的某一部位的形状或大小来表达，一个样本用一个脸谱来表达。相似的样本，它们构成的脸谱也很相象。Chernoff 首先把脸谱图用于聚类分析之中。日本很快的引进了脸谱图，并用于多元分析的各种应用之中。

脸谱图各部分用 18 个变量 ( $x_1, x_2, \dots, x_{18}$ ) 来构成。当变量数  $P < 18$  时，可将脸中某几个部位固定，当  $P > 18$  时，可以设法在脸中再添加一些部位。脸谱由六部分来构成：脸的轮廓、鼻、嘴、眼、眼珠和眉。现在分别来介绍每一部分的画法。

1、脸的轮廓：它由上下两个椭圆来构成，它们的短轴均在 Y 轴上，长轴平行于 X 轴，两椭圆交于 P 和 P'，自然 P 和 P' 关于 Y 轴对称，P 到原点 O 的位置由距离  $h^* = \overline{OP}$  和 OP 与 X 轴夹角  $\theta^*$  来决定 (图 3.1)，其中

$$h^* = \frac{1}{2}(1+x_1)H, \quad \theta^* = (2x_2-1)\frac{\pi}{4},$$

H 为一常数，由作图者决定，H 大，脸也大，脸的两个椭圆与 Y 轴分别交于 U 和 L，且

$$h \triangleq \overline{OU} = \overline{OL}, \quad h = \frac{1}{2}(1+x_3)H$$

由上述条件还不能唯一确定脸的两个椭圆，进一步规定，脸

的上半部椭圆的离心率为  $x_4$ ，下半部椭圆的离心率为  $x_5$ 。

我们知道，椭圆的标准方程为

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} = 1,$$

$a$ 、 $b$  分别为它的长轴和短轴。离心率  $e$  的定义是

$$e = \sqrt{a^2 - b^2} / a,$$

$a$ 、 $b$ 、 $e$  中只要知道了其中任两个就可决定另一个。利用离心率，椭圆的标准方程可以表为

$$(1 - e^2)x^2 + y^2 = b^2. \quad (2)$$

给了上述条件后如何来作轴的两个椭圆呢？现以下半椭圆来说明之。记下半椭圆的长轴，短轴分别为  $a$  和  $b$ ，记  $d \cong h - b$ ， $P$  点的坐标为  $(X^*, Y^*)$ ，显见

$$X^* = h^* \cos \theta^* \quad Y^* = h^* \sin \theta^*.$$

椭圆中心的坐标为  $(0, -d)$ ，故椭圆方程为

$$x^2(1 - e^2) + (y + d)^2 = b^2 = (h - d)^2$$

它过  $P$  点，应有

$$X^{*2}(1 - e^2) + (Y^* + d)^2 = (h - d)^2,$$

解得

$$d = \frac{(h^2 - Y^{*2}) - (1 - e^2)X^{*2}}{2(Y^* + h)}$$

$e$  为离心率 ( $= x_5$ )，上式  $h$ 、 $Y^*$ 、 $X^*$ 、 $e$  均为已知，故可求出

d. 从而便可得到脸下半部的椭圆方程。上半部的椭圆方程求法是类似的。

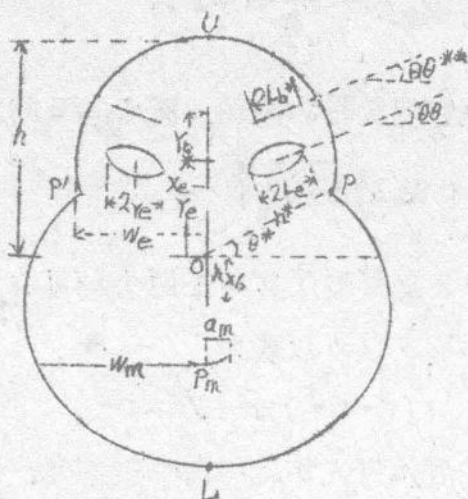


图3-1 脸谱图

2、鼻：以O为中心在Y轴上，上下各取长度 $hx_6$ 画一条粗线。

3、嘴：在O点下方， $h[x_7+(1-x_7)x_6]$ 的位置，用半径为 $h/|x_8|$ 的圆弧来描述，规定： $x_8$ 为正，圆弧向上； $x_8$ 为负，圆弧向下。嘴的大小由 $a_m$ 来决定， $a_m=x_9(h/|x_8|)$ ，嘴的圆弧关于Y轴对称。如果嘴太大，越过脸的轮廓，这时用 $x_9w_m$ 来代替，其中 $w_m$ 为点 $P_m$ 到脸轮廓的水平距离。

4、眼：用两个椭圆来表达，椭圆的中心为： $(x_e, y_e)$ ， $(-x_e, y_e)$ ，其中规定

$$y_e=h[x_{10}+(1-x_{10})x_6], \quad x_e=w_e(1+2x_{11})/4$$

其中 $w_e$ 为点 $(0, y_e)$ 至脸轮廓的水平距离。两眼椭圆的长轴

与X轴的夹角分别为 $\theta$ 和 $\pi - \theta$ ，其中 $\theta = (2x_{12} - 1)\pi/5$ ，其长轴为 $L_e$ ， $L_e = x_{14} \min(x_e, w_e - x_e)$ 。椭圆的离心率为 $x_{13}$ 。

5、眼珠：从眼的椭圆中心，沿着椭圆的长轴至 $\pm r_e(2x_{15} - 1)$ 的位置，其中

$$r_e = (\cos^2 \theta + \sin^2 \theta / x_{13}^2)^{-1/2} L_e,$$

$r_e$ 为眼的长轴在X轴上投影的长度。在作图时要注意务使眼珠为左右对称。

6、眉：从眼的椭圆中心向上至 $y_b$ 的高度决定了眉的中心，其长度为 $2L_b$ ，它与水平的夹角为 $\theta^{**}$ ，其中

$$\theta^{**} = \theta + 2(x_{17} - 1)\pi/5, \quad y_b = 2(x_{16} + 0.3)L_e x_{13},$$

$$L_b = r_e(2x_{18} + 1)/2.$$



表 3·1 脸谱图各变量的意义

变 量	描述脸的变量	变 换 公 式	数 值 范 围	备 注
X <sub>1</sub>	h*	$h^* = \frac{1}{2} (1 + X_1) H$	0-1	OP之长度, H为脸的大小比例数
X <sub>2</sub>	$\theta^*$	$\theta^* = (2X_2 - 1) \pi / 4$	0-1	X轴与OP的角度
X <sub>3</sub>	h	$h = \frac{1}{2} (1 + X_3) H$	0-1	OU=OL之长度
X <sub>4</sub>	X <sub>4</sub>		0.2-0.8	脸的上半椭圆离心率
X <sub>5</sub>	X <sub>5</sub>		0.2-0.8	脸的下半椭圆离心率
X <sub>6</sub>	X <sub>6</sub>	$2hX_6$	0.1-0.7	鼻的长度
X <sub>7</sub>	P <sub>m</sub>	$P_m = h [X_7 + (1 - X_7) X_6]$	0-1	嘴的位置
X <sub>8</sub>	X <sub>8</sub>		-5-5	嘴的曲率 (半径 = $h /  X_8 $ )
X <sub>9</sub>	a <sub>m</sub>	$a_m = X_9 (h /  X_8 )$ 或 $X_9 W_m$	0-1	嘴的大小
X <sub>10</sub>	Y <sub>e</sub>	$Y_e = h [X_{10} + (1 - X_{10}) X_6]$	0-1	眼的位置 (纵坐标)
X <sub>11</sub>	X <sub>e</sub>	$X_e = W_e (1 + 2X_{11}) / 4$	0-1	眼的位置 (横坐标)
X <sub>12</sub>	$\theta$	$\theta = (2X_{12} - 1) \pi / 5$	0-1	眼的倾斜角
X <sub>13</sub>	X <sub>13</sub>		0.4-0.8	眼的椭圆离心率
X <sub>14</sub>	L <sub>e</sub>	$L_e = X_{14} \min(X_e, W_e - X_e)$	0-1	眼的长轴之长度
X <sub>15</sub>	X <sub>15</sub>		0-1	眼珠的位置
X <sub>16</sub>	Y <sub>b</sub>	$Y_b = 2 (X_{16} + 0.3) L_e X_{13}$	0-1	眼到眉的高度
X <sub>17</sub>	$\theta^{**}$	$\theta^{**} = \theta + 2(1 - X_{17}) \pi / 5$	0-1	眉的倾斜角
X <sub>18</sub>	L <sub>b</sub>	$L_b = L_e (2X_{18} + 1) / 2$	0-1	眉的长度 (2L <sub>b</sub> )