

第4章 数据的概括性度量

PowerPoint



第 4 章 数据的概括性度量

- 4.1 集中趋势的度量
- 4.2 离散程度的度量
- 4.3 偏态与峰态的度量

学习目标

1. 集中趋势各测度值的计算方法
2. 集中趋势各测度值的特点及应用场合
3. 离散程度各测度值的计算方法
4. 离散程度各测度值的特点及应用场合
5. 偏态与峰态的测度方法
6. 用**Excel**计算描述统计量并进行分析

4.1 集中趋势的度量

4.1.1 分类数据：众数

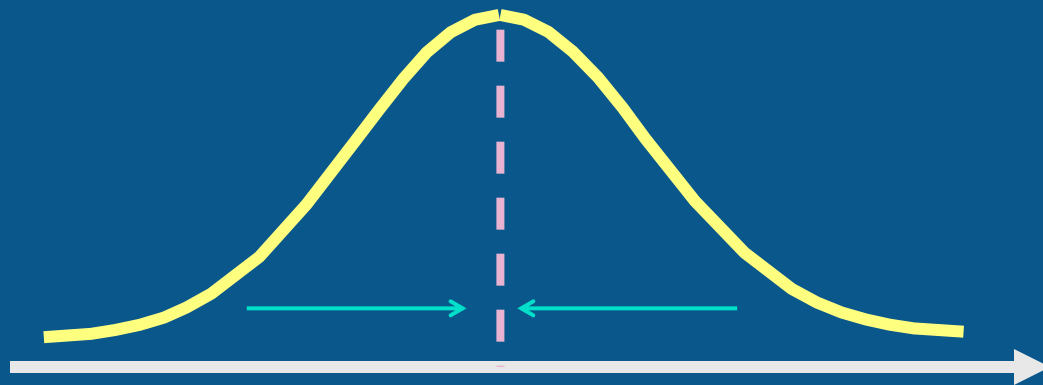
4.1.2 顺序数据：中位数和分位数

4.1.3 数值型数据：平均数

4.1.4 众数、中位数和平均数的比较

集中趋势 (central tendency)

1. 一组数据向其中心值靠拢的倾向和程度
2. 测度集中趋势就是寻找数据水平的代表值或中心值
3. 不同类型的数据用不同的集中趋势测度值
4. 低层次数据的测度值适用于高层次的测量数据，但高层次数据的测度值并不适用于低层次的测量数据



分类数据：众数

众数 (mode)

1. 一组数据中出现次数最多的变量值
2. 适合于数据量较多时使用
3. 不受极端值的影响
4. 一组数据可能没有众数或有几个众数
5. 主要用于分类数据，也可用于顺序数据和数值型数据

众数 (不惟一性)

无众数

原始数据: 10 5 9 12 6 8



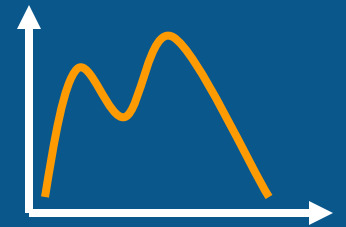
一个众数

原始数据: 6 **5** 9 8 **5** **5**



多于一个众数

原始数据: 25 **28** **28** 36 **42** **42**



分类数据的众数

(例题分析)

不同品牌饮料的频数分布

饮料品牌	频数	比例	百分比 (%)
果汁	6	0.12	12
矿泉水	10	0.20	20
绿茶	11	0.22	22
其他	8	0.16	16
碳酸饮料	15	0.30	30
合计	50	1	100

解：这里的变量为“饮料品牌”，这是个分类变量，不同类型的饮料就是变量值

所调查的50人中，购买碳酸饮料的人数最多，为15人，占总被调查人数的30%，因此众数为“可口可乐”这一品牌，即

$M_o = \text{碳酸饮料}$

顺序数据的众数

(例题分析)

甲城市家庭对住房状况评价的频数分布

回答类别	甲城市	
	户数 (户)	百分比 (%)
非常不满意	24	8
不满意	108	36
一般	93	31
满意	45	15
非常满意	30	10
合计	300	100.0

解：这里的数据为顺序数据。变量为“回答类别”

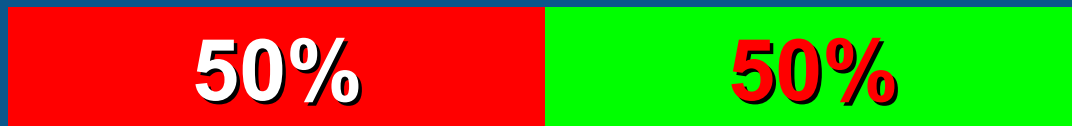
甲城市中，对住房表示不满意的户数最多，为108户，因此众数为“不满意”这一类别，即

$$M_0 = \text{不满意}$$

顺序数据：中位数和分位数

中位数 (median)

1. 排序后处于中间位置上的值



M_e

2. 不受极端值的影响
3. 主要用于顺序数据，也可用数值型数据，但不能用于分类数据
4. 各变量值与中位数的离差绝对值之和最小，即

$$\sum_{i=1}^n |x_i - M_e| = \min$$

中位数

(位置和数值的确定)

位置确定

$$\text{中位数位置} = \frac{n+1}{2}$$

数值确定

$$M_e = \begin{cases} x_{\left(\frac{n+1}{2}\right)} & n \text{为奇数} \\ \frac{1}{2} \left\{ x_{\left(\frac{n}{2}\right)} + x_{\left(\frac{n}{2}+1\right)} \right\} & n \text{为偶数} \end{cases}$$

顺序数据的中位数

(例题分析)

甲城市家庭对住房状况评价的频数分布

回答类别	甲城市	
	户数 (户)	累计频数
非常不满意	24	24
不满意	108	132
一般	93	225
满意	45	270
非常满意	30	300
合计	300	—

解：中位数的位置为
 $(300+1)/2=150.5$

从累计频数看，
中位数在“一般”这一
组别中

中位数为

$$M_e = \text{一般}$$

数值型数据的中位数

(9个数据的算例)

【例】 9个家庭的人均月收入数据

原始数据:	1500	750	780	1080	850	960	2000	1250	1630
排 序:	750	780	850	960	1080	1250	1500	1630	2000
位 置:	1	2	3	4	5	6	7	8	9



$$\text{位置} = \frac{n+1}{2} = \frac{9+1}{2} = 5$$

中位数  1080

数值型数据的中位数

(10个数据的算例)

【例】：10个家庭的人均月收入数据

排	序:	660	750	780	850	960	1080	1250	1500	1630	2000
位	置:	1	2	3	4	5	6	7	8	9	10



$$\text{位置} = \frac{n+1}{2} = \frac{10+1}{2} = 5.5$$

$$\text{中位数} = \frac{960+1080}{2} = 1020$$

四分位数 (quartile)

1. 排序后处于25%和75%位置上的值

25%

25%

25%

25%

Q_L

Q_M

Q_U

2. 不受极端值的影响

3. 计算公式

$$\begin{cases} Q_L \text{位置} = \frac{n}{4} \\ Q_U \text{位置} = \frac{3n}{4} \end{cases}$$

顺序数据的四分位数 (例题分析)

甲城市家庭对住房状况评价的频数分布

回答类别	甲城市	
	户数 (户)	累计频数
非常不满意	24	24
不满意	108	132
一般	93	225
满意	45	270
非常满意	30	300
合计	300	—

解: Q_L 位置 = $(300)/4 = 75$

$$Q_U \text{位置} = (3 \times 300)/4 \\ = 225$$

从累计频数看, Q_L 在“不满意”这一组别中; Q_U 在“一般”这一组别中

四分位数为

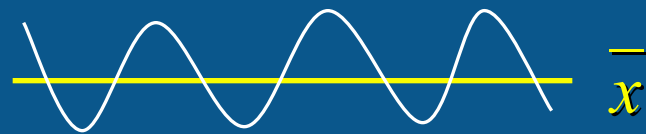
$$Q_L = \text{不满意}$$

$$Q_U = \text{一般}$$

数值型数据：平均数

平均数 (mean)

1. 也称为均值
2. 集中趋势的最常用测度值
3. 一组数据的均衡点所在
3. 体现了数据的必然性特征
4. 易受极端值的影响
5. 有简单平均数和加权平均数之分
6. 根据总体数据计算的，称为平均数，记为 μ ；根据样本数据计算的，称为样本平均数，记为 \bar{x}



简单平均数 (Simple mean)

设一组数据为： x_1, x_2, \dots, x_n (总体数据 x_N)

样本平均数 $\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n}$

总体平均数 $\mu = \frac{x_1 + x_2 + \dots + x_N}{N} = \frac{\sum_{i=1}^N x_i}{N}$

加权平均数 (Weighted mean)

设各组的组中值为: M_1, M_2, \dots, M_k

相应的频数为: f_1, f_2, \dots, f_k

样本加权平均 $\bar{x} = \frac{M_1 f_1 + M_2 f_2 + \dots + M_k f_k}{f_1 + f_2 + \dots + f_k} = \frac{\sum_{i=1}^k M_i f_i}{n}$

总体加权平均 $\mu = \frac{M_1 f_1 + M_2 f_2 + \dots + M_k f_k}{f_1 + f_2 + \dots + f_k} = \frac{\sum_{i=1}^k M_i f_i}{N}$

加权平均数 (例题分析)

某电脑公司销售量数据分组表

按销售量分组	组中值(M_i)	频数(f_i)	$M_i f_i$
140~150	145	4	580
150~160	155	9	1395
160~170	165	16	2640
170~180	175	27	4725
180~190	185	20	3700
190~200	195	17	3315
200~210	205	10	2050
210~220	215	8	1720
220~230	225	4	900
230~240	235	5	1175
合计	—	120	22200



$$\begin{aligned}\bar{x} &= \frac{\sum_{i=1}^k M_i f_i}{n} \\ &= \frac{22200}{120} = 185\end{aligned}$$

几何平均数 (geometric mean)

1. n 个变量值乘积的 n 次方根
2. 适用于对比率数据的平均
3. 主要用于计算平均增长率
4. 计算公式为

$$G_m = \sqrt[n]{x_1 \times x_2 \times \cdots \times x_n} = \sqrt[n]{\prod_{i=1}^n x_i}$$

5. 可看作是平均数的一种变形

$$\lg G_m = \frac{1}{n} (\lg x_1 + \lg x_2 + \cdots + \lg x_n) = \frac{\sum_{i=1}^n \lg x_i}{n}$$

几何平均数 (例题分析)

【例】一位投资者购持有有一种股票，连续4年收益率分别为4.5%、2.1%、25.5%、1.9%。计算该投资者在这四年内的平均收益率

几何平均:

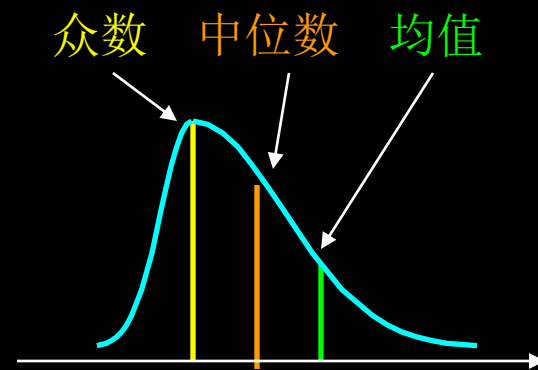
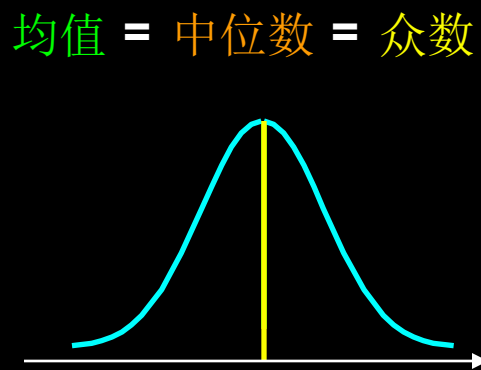
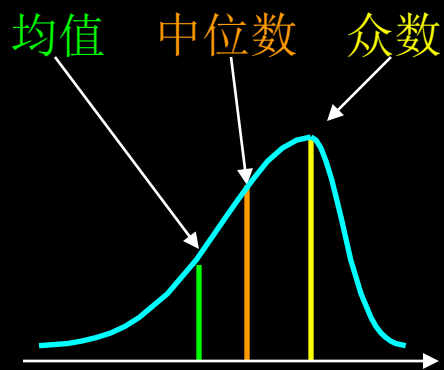
$$\begin{aligned}\bar{G} &= \sqrt[4]{104.5\% \times 102.1\% \times 125.5\% \times 101.9\%} - 1 \\ &= 8.0787\%\end{aligned}$$

算术平均:

$$\bar{G} = (4.5\% + 2.1\% + 25.5\% + 1.9\%) \div 4 = 8.5\%$$

众数、中位数和平均数的比较

众数、中位数和平均数的关系



众数、中位数、平均数的特点和应用

1. 众数

- 不受极端值影响
- 具有不惟一性
- 数据分布偏斜程度较大且有明显峰值时应用

2. 中位数

- 不受极端值影响
- 数据分布偏斜程度较大时应用

3. 平均数

- 易受极端值影响
- 数学性质优良
- 数据对称分布或接近对称分布时应用

4.2 离散程度的度量

4.2.1 分类数据：异众比率

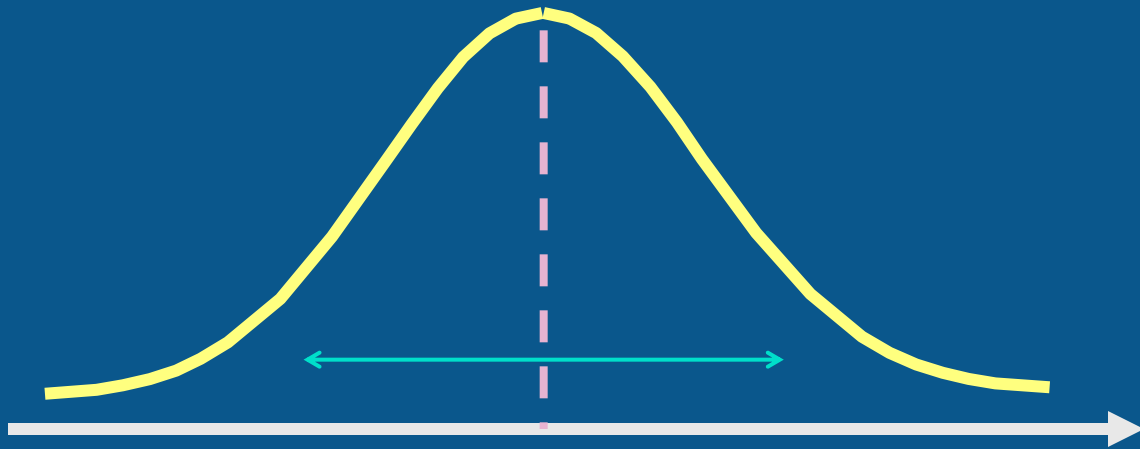
4.2.2 顺序数据：四分位差

4.2.3 数值型数据：方差和标准差

4.2.4 相对离散程度：离散系数

离中趋势

1. 数据分布的另一个重要特征
2. 反映各变量值远离其中心值的程度(离散程度)
3. 从另一个侧面说明了集中趋势测度值的代表程度
4. 不同类型的数据有不同的离散程度测度值



分类数据：异众比率

异众比率 (variation ratio)

1. 对分类数据离散程度的测度
2. 非众数组的频数占总频数的比例
3. 计算公式为

$$v_r = \frac{\sum f_i - f_m}{\sum f_i} = 1 - \frac{f_m}{\sum f_i}$$

4. 用于衡量众数的代表性

异众比率 (例题分析)

不同品牌饮料的频数分布

饮料品牌	频数	比例	百分比(%)
果汁	6	0.12	12
矿泉水	10	0.20	20
绿茶	11	0.22	22
其他	8	0.16	16
碳酸饮料	15	0.30	30
合计	50	1	100

解：

$$\begin{aligned}v_r &= \frac{50 - 15}{50} \\ &= 1 - \frac{15}{50} \\ &= 0.7 = 70\%\end{aligned}$$

在所调查的50人当中，购买其他品牌饮料的人数占70%，异众比率比较大。因此，用“碳酸饮料”代表消费者购买饮料品牌的状况，其代表性不是很好

顺序数据：四分位差

四分位差 (quartile deviation)

1. 对顺序数据离散程度的测度
2. 也称为内距或四分间距
3. 上四分位数与下四分位数之差

$$Q_d = Q_U - Q_L$$

4. 反映了中间50%数据的离散程度
5. 不受极端值的影响
6. 用于衡量中位数的代表性

四分位差 (例题分析)

甲城市家庭对住房状况评价的频数分布

回答类别	甲城市	
	户数 (户)	累计频数
非常不满意	24	24
不满意	108	132
一般	93	225
满意	45	270
非常满意	30	300
合计	300	—

解：设非常不满意为1, 不满意为2, 一般为3, 满意为4, 非常满意为5。已知

$$Q_L = \text{不满意} = 2$$

$$Q_U = \text{一般} = 3$$

四分位差为

$$\begin{aligned} Q_d &= Q_U - Q_L \\ &= 3 - 2 = 1 \end{aligned}$$

数值型数据：方差和标准差

极差 (range)

1. 一组数据的最大值与最小值之差
2. 离散程度的最简单测度值
3. 易受极端值影响
4. 未考虑数据的分布
5. 计算公式为

$$R = \max(x_i) - \min(x_i)$$

平均差 (mean deviation)

1. 各变量值与其平均数离差绝对值的平均数
2. 能全面反映一组数据的离散程度
3. 数学性质较差，实际中应用较少
4. 计算公式为

未分组数据 $M_d = \frac{\sum_{i=1}^n |x_i - \bar{x}|}{n}$

组距分组数据 $M_d = \frac{\sum_{i=1}^k |M_i - \bar{x}| f_i}{n}$

平均差

(例题分析)

某电脑公司销售量数据平均差计算表

按销售量分组	组中值(M_i)	频数(f_i)	$ x - \bar{x} $	$ x - \bar{x} f$
140~150	145	4	40	160
150 ~ 160	155	9	30	270
160 ~ 170	165	16	20	320
170 ~ 180	175	27	10	270
180 ~ 190	185	20	0	0
190 ~ 200	195	17	10	170
200 ~ 210	205	10	20	200
210 ~ 220	215	8	30	240
220 ~ 230	225	4	40	160
230 ~ 240	235	5	50	250
合计	—	120	—	2040

平均差 (例题分析)

$$M_d = \frac{\sum_{i=1}^k |M_i - \bar{x}| f_i}{n} = \frac{2040}{120} = 17(\text{台})$$

含义：每一天的销售量平均数相比，
平均相差17台



方差和标准差

(variance and standard deviation)

1. 数据离散程度的最常用测度值
2. 反映了各变量值与均值的平均差异
3. 根据总体数据计算的，称为总体方差(标准差)，记为 $\sigma^2(\sigma)$ ；根据样本数据计算的，称为样本方差(标准差)，记为 $s^2(s)$

样本方差和标准差

(sample variance and standard deviation)

方差的计算公式

未分组数据

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

组距分组数据

$$s^2 = \frac{\sum_{i=1}^k (M_i - \bar{x})^2 f_i}{n-1}$$

注意：

样本方差用自由度 $n-1$ 去除！



标准差的计算公式

未分组数据

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$$

组距分组数据

$$s = \sqrt{\frac{\sum_{i=1}^k (M_i - \bar{x})^2 f_i}{n-1}}$$

自由度

(degree of freedom)

1. 自由度是指数据个数与附加给独立的观测值的约束或限制的个数之差
2. 从字面涵义来看，自由度是指一组数据中可以自由取值的个数
3. 当样本数据的个数为 n 时，若样本平均数确定后，则附加给 n 个观测值的约束个数就是1个，因此只有 $n-1$ 个数据可以自由取值，其中必有一个数据不能自由取值
4. 按着这一逻辑，如果对 n 个观测值附加的约束个数为 k 个，自由度则为 $n-k$

自由度

(degree of freedom)

1. 样本有3个数值，即 $x_1=2$ ， $x_2=4$ ， $x_3=9$ ，则 $\bar{x}=5$ 。当 $\bar{x}=5$ 确定后， x_1 ， x_2 和 x_3 有两个数据可以自由取值，另一个则不能自由取值，比如 $x_1=6$ ， $x_2=7$ ，那么 x_3 则必然取2，而不能取其他值
2. 为什么样本方差的自由度为什么是 $n-1$ 呢？因为在计算离差平方和时，必须先求出样本均值 \bar{x} ，而 \bar{x} 则是附件给离差平方和的一个约束，因此，计算离差平方和时只有 $n-1$ 个独立的观测值，而不是 n 个
3. 样本方差用自由度去除，其原因可从多方面解释，从实际应用角度看，在抽样估计中，当用样本方差 s^2 去估计总体方差 σ^2 时，它是 σ^2 的无偏估计量

样本标准差

(例题分析)

某电脑公司销售量数据平均差计算表

按销售量分组	组中值(M_i)	频数(f_i)	$(M_i - \bar{x})^2$	$(M_i - \bar{x})^2 f_i$
140~150	145	4	40	160
150 ~ 160	155	9	30	270
160 ~170	165	16	20	320
170 ~180	175	27	10	270
180 ~ 190	185	20	0	0
190 ~ 200	195	17	10	170
200 ~ 210	205	10	20	200
210 ~220	215	8	30	240
220 ~230	225	4	40	160
230 ~240	235	5	50	250
合计	—	120	—	55400

样本标准差 (例题分析)

$$s = \sqrt{\frac{\sum_{i=1}^k (M_i - \bar{x})^2 f_i}{n-1}}$$
$$= \sqrt{\frac{55400}{120-1}} = 21.58(\text{台})$$



含义：每一天的销售量与平均数相比，
平均相差21.58台

总体方差和标准差

(Population variance and Standard deviation)

方差的计算公式

未分组数据

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

组距分组数据

$$\sigma^2 = \frac{\sum_{i=1}^K (M_i - \mu)^2 f_i}{N}$$

标准差的计算公式

未分组数据

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$

组距分组数据

$$\sigma = \sqrt{\frac{\sum_{i=1}^K (M_i - \mu)^2 f_i}{N}}$$

相对位置的度量：标准分数

标准分数 (standard score)

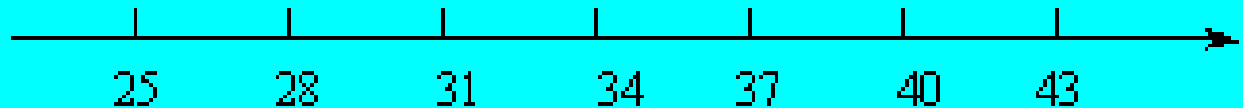
1. 也称标准化值
2. 对某一个值在一组数据中相对位置的度量
3. 可用于判断一组数据是否有离群点(outlier)
4. 用于对变量的标准化处理
5. 计算公式为

$$z_i = \frac{x_i - \bar{x}}{s}$$

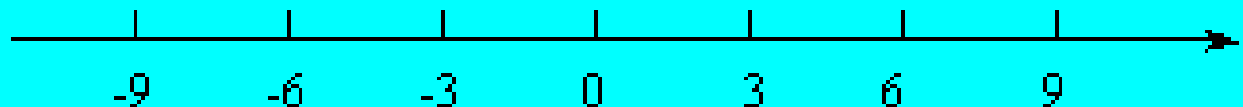
标准分数 (性质)

z分数只是将原始数据进行了线性变换，它并没有改变一个数据在该组数据中的位置，也没有改变该组数分布的形状，而只是使该组数据均值为0，标准差为1

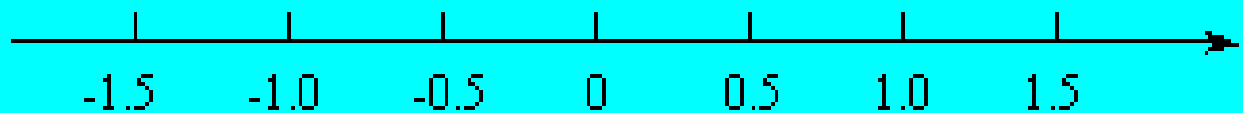
原始数据：



减去 34：



除以 6：



标准分数 (例题分析)

9个家庭人均月收入标准化值计算表

家庭编号	人均月收入 (元)	标准化值 z
1	1500	0.695
2	750	-1.042
3	780	-0.973
4	1080	-0.278
5	850	-0.811
6	960	-0.556
7	2000	1.853
8	1250	0.116
9	1630	0.996

经验法则

- 经验法则表明：当一组数据对称分布时
- 约有**68%**的数据在平均数加减1个标准差的范围之内
 - 约有**95%**的数据在平均数加减2个标准差的范围之内
 - 约有**99%**的数据在平均数加减3个标准差的范围之内

切比雪夫不等式 (Chebyshev's inequality)

1. 如果一组数据不是对称分布，经验法则就不再适用，这时可使用切比雪夫不等式，它对任何分布形状的数据都适用
2. 切比雪夫不等式提供的是“下界”，也就是“所占比例至少是多少”
3. 对于任意分布形态的数据，根据切比雪夫不等式，至少有 $1 - 1/k^2$ 的数据落在平均数加减 k 个标准差之内。其中 k 是大于 1 的任意值，但不一定是整数

切比雪夫不等式 (Chebyshev's inequality)

- 对于 $k=2, 3, 4$, 该不等式的含义是
1. 至少有**75%**的数据落在平均数加减**2**个标准差的范围之内
 2. 至少有**89%**的数据落在平均数加减**3**个标准差的范围之内
 3. 至少有**94%**的数据落在平均数加减**4**个标准差的范围之内

相对离散程度：离散系数

离散系数

(coefficient of variation)

1. 标准差与其相应的均值之比
2. 对数据相对离散程度的测度
3. 消除了数据水平高低和计量单位的影响
4. 用于对不同组别数据离散程度的比较
5. 计算公式为

$$v_s = \frac{s}{\bar{x}}$$

离散系数 (例题分析)

【例】某管理局抽查了所属的8家企业，其产品销售数据如表。试比较产品销售额与销售利润的离散程度

某管理局所属8家企业的产品销售数据

企业编号	产品销售额（万元）	销售利润（万元）
	x_1	x_2
1	170	8.1
2	220	12.5
3	390	18.0
4	430	22.0
5	480	26.5
6	650	40.0
7	950	64.0
8	1000	69.0

离散系数 (例题分析)

$$\bar{x}_1 = 536.25(\text{万元})$$

$$s_1 = 309.19(\text{万元})$$

$$v_1 = \frac{309.19}{536.25} = 0.577$$

$$\bar{x}_2 = 32.5215(\text{万元})$$

$$s_2 = 23.09(\text{万元})$$

$$v_2 = \frac{23.09}{32.5215} = 0.710$$

结论： 计算结果表明， $v_1 < v_2$ ，说明产品销售额的离散程度小于销售利润的离散程度

4.3 偏态与峰态的度量

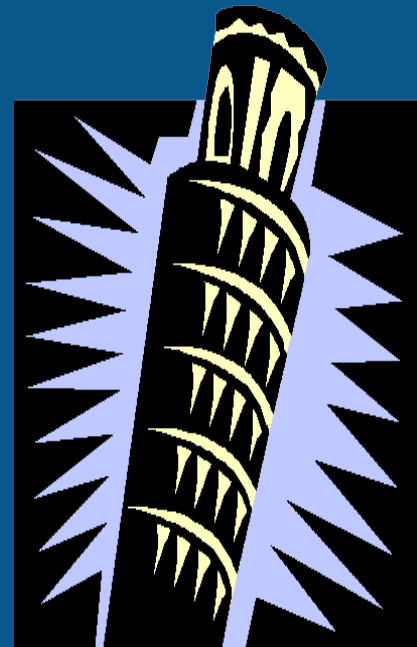
4.3.1 偏态及其测度

4.3.2 峰态及其测度

偏 态

偏态 (skewness)

1. 统计学家Pearson于1895年首次提出
2. 数据分布偏斜程度的测度
2. 偏态系数=0为对称分布
3. 偏态系数 > 0 为右偏分布
4. 偏态系数 < 0 为左偏分布
5. 偏态系数大于1或小于-1，被称为高度偏态分布；偏态系数在0.5~1或-1~-0.5之间，被认为是中等偏态分布；偏态系数越接近0，偏斜程度就越低



偏态系数 (coefficient of skewness)

1. 根据原始数据计算

$$SK = \frac{n \sum (x_i - \bar{x})^3}{(n-1)(n-2)s^3}$$

2. 根据分组数据计算

$$SK = \frac{\sum_{i=1}^k (M_i - \bar{x})^3 f_i}{ns^3}$$

偏态系数 (例题分析)

某电脑公司销售量偏态及峰度计算表

按销售量份组(台)	组中值(M_i)	频数 f_i	$(M_i - \bar{x})^3 f_i$	$(M_i - \bar{x})^4 f_i$
140 ~ 150	145	4	-256000	10240000
150 ~ 160	155	9	-243000	7290000
160 ~ 170	165	16	-128000	2560000
170 ~ 180	175	27	-27000	270000
180 ~ 190	185	20	0	0
190 ~ 200	195	17	17000	170000
200 ~ 210	205	10	80000	1600000
210 ~ 220	215	8	216000	6480000
220 ~ 230	225	4	256000	10240000
230 ~ 240	235	5	625000	31250000
合计	—	120	540000	70100000

偏态系数 (例题分析)

$$\begin{aligned} SK &= \frac{\sum_{i=1}^k (M_i - \bar{x})^3 f_i}{ns^3} = \frac{\sum_{i=1}^{10} (M_i - 185)^3 f_i}{120 \times (21.58)^3} \\ &= \frac{540000}{120 \times (21.58)^3} = 0.448 \end{aligned}$$

结论：偏态系数为正值，但与0的差异不大，说明电脑销售量为轻微右偏分布，即销售量较少的天数占据多数，而销售量较多的天数则占少数

峰 态

峰态 (kurtosis)

1. 统计学家Pearson于1905年首次提出
2. 数据分布扁平程度的测度
3. 峰态系数=0扁平峰度适中
4. 峰态系数<0为扁平分布
5. 峰态系数>0为尖峰分布



峰态系数

(coefficient of kurtosis)

1. 根据原始数据计算

$$K = \frac{n(n+1) \sum (x_i - \bar{x})^4 - 3 \left[\sum (x_i - \bar{x})^2 \right]^2 (n-1)}{(n-1)(n-2)(n-3)s^4}$$

2. 根据分组数据计算

$$K = \frac{\sum_{i=1}^k (M_i - \bar{x})^4 f_i}{ns^4} - 3$$

峰态系数 (例题分析)

$$K = \frac{\sum_{i=1}^k (M_i - \bar{x})^4 f_i}{ns^4} - 3 = \frac{70100000}{120 \times (21.58)^4} - 3$$
$$= 2.694 - 3 = -0.306$$

结论：偏态系数为负值，但与0的差异不大，说明电脑销售量为轻微扁平分布

用Excel计算描述统计量

用Excel计算描述统计量

☞ 将120的销售量的数据输入到Excel工作表中，然后按下列步骤操作

第1步：选择【工具】下拉菜单

第2步：选择【数据分析】选项

第3步：在分析工具中选择【描述统计】，然后选择【确定】

第4步：当对话框出现时

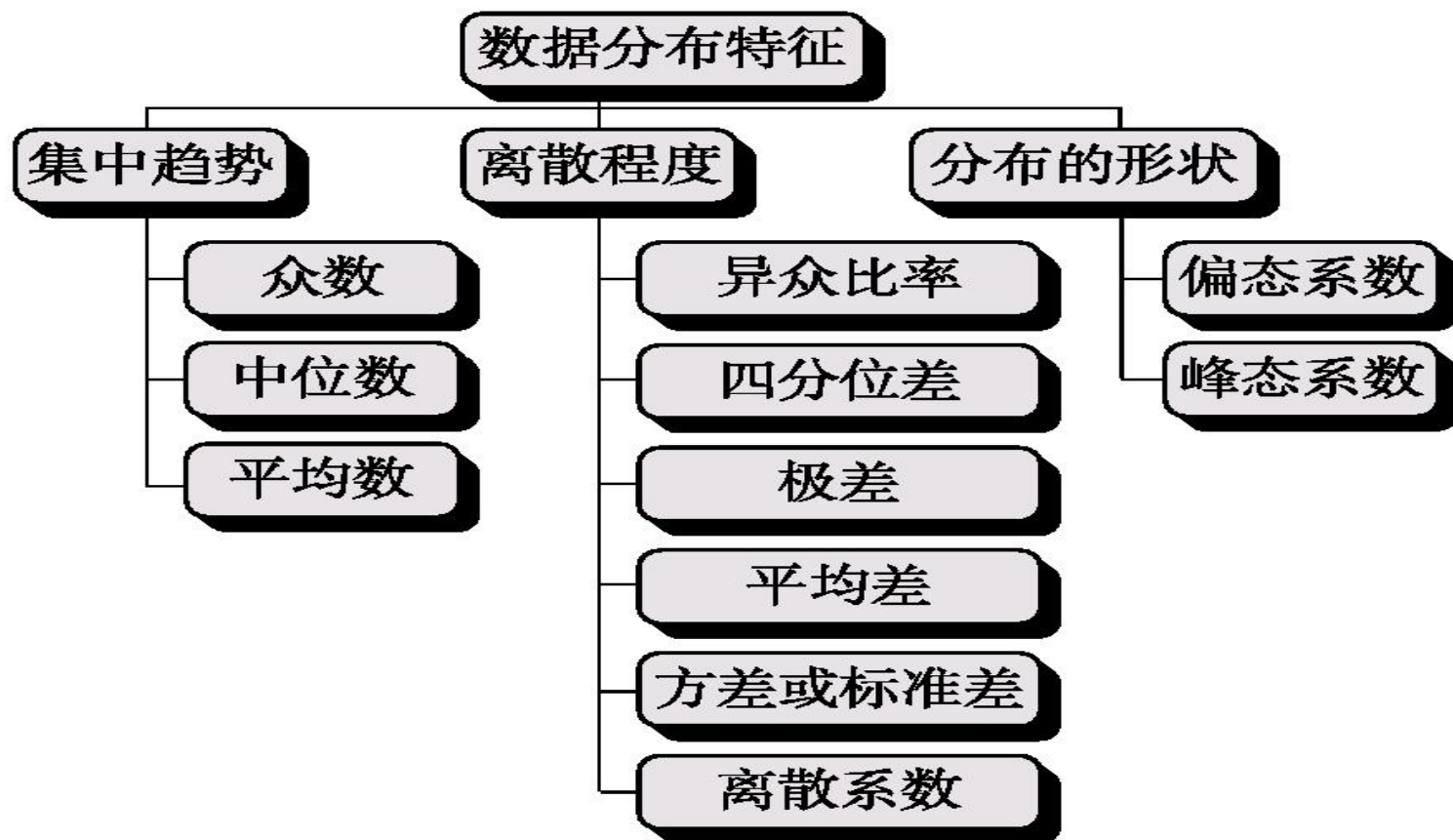
在【输入区域】方框内键入数据区域

在【输出选项】中选择输出区域

选择【汇总统计】

选择【确定】

数据分布特征和描述统计量



本章小节

1. 数据水平的概括性度量
2. 数据离散程度的概括性度量
3. 数据分布形状的度量
4. 用**Excel**计算描述统计量

结 束



THANKS