

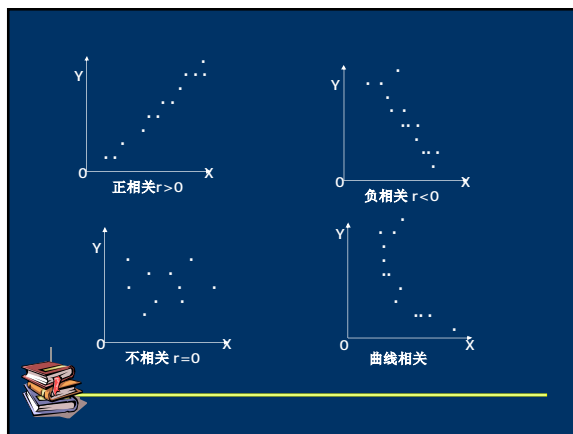
第四章 一元线性回归

§ 4.1 一元线性回归模型

一、变量之间的关系

函数关系 变量之间可以用数学公式表示。如圆的半径R与面积S等。

统计关系 两变量之间有一定依存关系，但没有严格的对应关系。如人的年龄和血压，身高和体重，储蓄额与居民收入等都是统计关系或相关关系。相关关系可以通过相关图表示出来。相关关系有线性相关和非线性相关（曲线相关）。



二、两变量之间的线性相关系数

1、定义：变量X与Y之间的线性相关程度可以用简单相关系数来度量。计算公式为：

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}} = \frac{L_{xy}}{\sqrt{L_{xx} L_{yy}}}$$

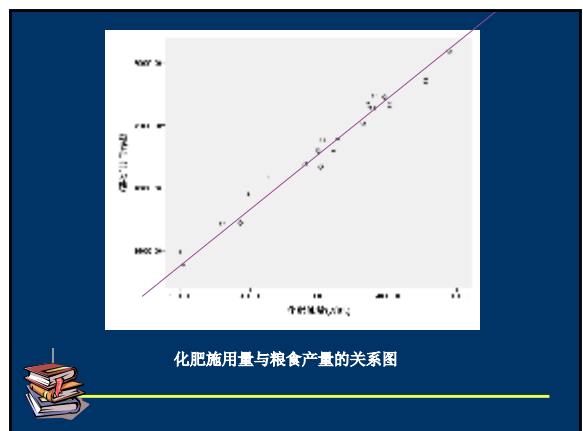
由样本数据得：

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

其中 $-1 \leq r \leq 1$ ，且 $r > 0$ 时称为正相关， $r < 0$ 时为负相关， $r = 0$ 时，两变量间不存在线性相关关系。

例 4.1 讨论化肥施用量与粮食产量之间的关系 数据如下表：

化肥施用量 x(万吨)	4541.05	3637.87	2287.49	3056.89	4883.7	3779.3	4021.09
粮食产量 y(万吨)	48526.69	45110.87	40753.79	43824.58	50890.11	46370.88	46577.91
化肥施用量 x(万吨)	2989.06	3021.9	3953.97	3212.13	3804.76	1598.28	1998.56
粮食产量 y(万吨)	42947.44	41673.21	47244.34	43061.53	47336.78	37127.89	39515.07
化肥施用量 x(万吨)	3710.56	3269.03	1017.12	1864.23	2797.24	1034.09	
粮食产量 y(万吨)	46598.04	44020.92	34866.91	37184.14	41864.77	33717.78	



2、相关系数的显著性检验

对样本相关系数 r 的相关性检验是确定两变量间的线性关系是否显著。

相关系数的检验过程：

- 1、假设 $H_0: r=0$; $H_1: r \neq 0$
- 2、根据给出的样本数据，计算样本相关系数 r 和统计量 t ；在原假设成立的条件下，检验统计量 t 服从 t 分布。

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \quad t \sim t(n-2)$$

- 3、根据给定的显著性水平 α ，查临界值 $t_{\alpha/2}(n-2)$ ；
- 4、若 $t \geq t_{\alpha/2}(n-2)$ ，说明 r 在统计上是显著的，即两个变量之间存在线性关系，若 $t < t_{\alpha/2}(n-2)$ ，则两个变量之间不存在线性相关。

r 为两变量的相关系数的真值



SPSS统计软件得到的化肥施用量与粮食产量的相关性输出表：

Correlations		
	化肥施用量(万吨)	粮食产量(万吨)
化肥施用量(万吨)	1	.989*
	Sig. (2-tailed)	.000
	N	20
粮食产量(万吨)	.989*	1
	Sig. (2-tailed)	.000
	N	20

**. Correlation is significant at the 0.01 level (2-tailed).

注意：在相关分析中两个变量的地位是平等的。



三、一元线性回归的数学模型

如果两个变量之间存在线性关系，设一个变量为自变量，另一个是因变量，则两个变量之间的关系可以用一元线性回归模型表达。

设自变量 X 为确定性变量（解释变量），因变量 Y （被解释变量）为随机变量，两者之间的数学结构式为

$$y = b_0 + b_1x + e$$

式中： b_0, b_1 是回归系数， e 是随机项，表示除了变量 X 之外其他因素对变量 Y 的影响。

线性方程由两部分组成，一部分由 X 的变化引起，另一部分是由随机因素引起。



基本假设

假设1：随机项 e 服从正态分布 $N(0, \sigma^2)$ ，即：

$$E(e) = 0, \quad Var(e) = \sigma^2$$

假设2：随机项 e_i 之间是相互独立的，并具有相同的方差，即：

$$Cov(e_i, e_j) = \begin{cases} \sigma^2 & i = j \\ 0 & i \neq j \end{cases} \quad i, j = 1, 2, \dots, n$$

假设3：样本数据是相互独立的。



由假设知，随机变量 y 也服从正态分布，由公式 $y = b_0 + b_1x + e$ 得出：

$$E(y) = b_0 + b_1x, \quad Var(y) = \sigma^2$$

在一般情况下，从研究的总体中抽取一个样本观察值 (x_i, y_i) ， $i=1, 2, \dots, n$ ，对于样本 X, Y 的每一组数，有

$$y_i = b_0 + b_1x_i + e_i \quad i = 1, 2, \dots, n$$

由假设条件知， $e_i \sim N(0, \sigma^2)$ ，且 $E(e_i) = 0, Var(e_i) = \sigma^2, i = 1, 2, \dots, n$

推导出观察值 $y_i (i=1, 2, \dots, n)$ 也是相互独立的正态随机变量，且

$$E(y_i) = b_0 + b_1x_i, \quad Var(y_i) = \sigma^2$$

线性模型在平均意义上表达了变量 Y 与 X 的统计规律性。



回归模型的矩阵表达式：

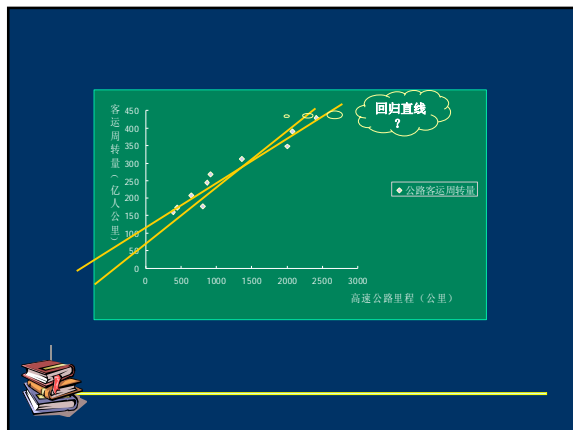
$$\begin{cases} Y = Xb + e \\ E(e) = 0 \\ Var(e) = \sigma^2 I_n \end{cases}$$

$$\text{其中: } Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad b = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}, \quad e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}, \quad X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

线性回归的任务是：通过样本观察值对回归系数进行估计，求出 b_0, b_1 的估计值 \hat{b}_0, \hat{b}_1 ，得出一元线性回归方程：

$$\hat{y} = \hat{b}_0 + \hat{b}_1x$$





§4.2 回归系数的最小二乘估计

一、普通最小二乘估计 (OLSE)

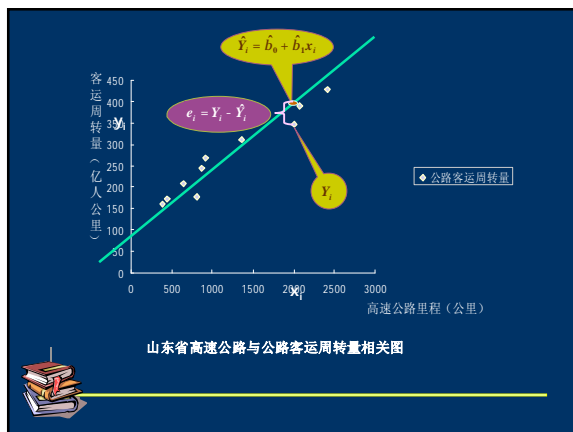
回归系数的估计是依赖于最小二乘法的基本思想, 考虑观察值 y 与回归值 \hat{y}_i 的差(称为残差, 记为 e_i). 通过使残差平方和为最小来估计回归系数。通过样本得到残差平方和为:

$$Q(b_0, b_1) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2 = \sum_{i=1}^n e_i^2$$

求极值得:

$$\begin{cases} \frac{\partial Q}{\partial b_0} = -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i) = 0 \\ \frac{\partial Q}{\partial b_1} = -2 \sum_{i=1}^n x_i (y_i - b_0 - b_1 x_i) = 0 \end{cases} \quad \leftarrow \begin{cases} \hat{a} e_i = 0 \\ \hat{a} e_i x_i = 0 \end{cases}$$

$\hat{y}_i = \hat{b}_0 + \hat{b}_1 x_i$



由上面方程组可得:

$$\begin{cases} \hat{a} y_i - n \hat{b}_0 - \hat{b}_1 \hat{a} x_i = 0 \\ \hat{a} x_i y_i - \hat{b}_0 \hat{a} x_i - \hat{b}_1 \hat{a} x_i^2 = 0 \end{cases}$$

$$\hat{b}_1 = \frac{L_{xy}}{L_{xx}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x} = \frac{1}{n} \sum_{i=1}^n y_i - \hat{b}_1 \frac{1}{n} \sum_{i=1}^n x_i$$

这种方法称普通最小二乘估计(OLSE), 是线性回归方程中回归系数求解的基本方法。这样由样本估计得到的回归方程称为一元线性经验回归方程 记为:

$$\hat{Y} = \hat{b}_0 + \hat{b}_1 X$$

例: 根据恩格尔定律得出食品支出 Y 与实际收入 X 的关系的一元线性回归模型来反映.用表的形式列出原始数据 X, Y 的值及相应的计算数据.见表

解: 根据计算表中的数据得出回归系数为:

$$\hat{b}_1 = \frac{15 \cdot 44632 - 1516 \cdot 423}{15 \cdot 163634 - 1516^2} = 0.1802$$

$$\hat{b}_0 = \frac{423}{15} - 0.1802 \cdot \frac{1516}{15} = 9.9872$$

所求的经验回归方程为:

$$\hat{y}_i = 9.99 + 0.1802 x_i$$

回归方程的**实际意义**是: 当收入每增加一个单位时, 食品支出会增加0.18单位, 即使在收入为0的情况下, 食品支出依然需要9.99单位。

编号	X	Y	XY	X ²	Y ²
1	102	27	2754	10404	729
2	96	26	2496	9216	676
3	97	25	2425	9409	625
4	102	28	2856	10404	784
5	91	27	2457	8281	729
6	158	36	5688	24964	1296
7	54	19	1026	2916	361
8	83	26	2158	6889	676
9	123	31	3813	15129	961
10	106	31	3286	11236	961
11	129	34	4386	16641	1156
12	138	38	5244	19044	1444
13	81	27	2187	6561	729
14	92	28	2576	8464	784
15	64	20	1280	4096	400
合计	1516	423	44632	163634	12311

§4.3 回归系数的性质及统计意义

一、 \hat{b}_0, \hat{b}_1 是 y_i 的线性组合。

例如：
$$\hat{b}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{L_{xx}} = \frac{\sum (x_i - \bar{x})y_i}{L_{xx}} = \sum \frac{(x_i - \bar{x})}{L_{xx}} y_i$$

$$\hat{b}_1 = \sum C_i y_i \quad C_i = \frac{x_i - \bar{x}}{L_{xx}} \quad i = 1, 2, \dots, n$$

同理得：
$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x} = \sum \left(\frac{1}{n} - C_i \bar{x} \right) y_i$$



二、 \hat{b}_0, \hat{b}_1 服从正态分布，且最小二乘估计的回归系数具有无偏性。

由于因变量 y_i 服从正态分布， \hat{b}_0 和 \hat{b}_1 是 y_i 的线性组合，故也服从正态分布。

即有：

$$E(\hat{b}_0) = b_0 \quad S_{\hat{b}_0}^2 = S^2 \left\{ \frac{1}{n} + \frac{\bar{x}^2}{L_{xx}} \right\} = S^2 \left\{ \frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right\}$$

$$E(\hat{b}_1) = b_1 \quad S_{\hat{b}_1}^2 = \frac{S^2}{L_{xx}} = \frac{S^2}{\sum (x_i - \bar{x})^2}$$

$$\hat{b}_0 \sim N\left(b_0, \left(\frac{1}{n} + \frac{\bar{x}^2}{L_{xx}}\right) S^2\right); \quad \hat{b}_1 \sim N\left(b_1, \frac{S^2}{L_{xx}}\right)$$



证明： $E(\hat{b}_1) = b_1$

因为：
$$\hat{b}_1 = C_1 y_1 + C_2 y_2 + \dots + C_n y_n$$

由于：
$$\sum C_i = \sum \frac{x_i - \bar{x}}{\sum (x_j - \bar{x})^2} = 0$$

$$\sum C_i x_i = \sum C_i (x_i - \bar{x}) = 1$$

$$\begin{aligned} E(\hat{b}_1) &= \sum C_i E(y_i) \\ &= \sum C_i (b_0 + b_1 x_i) \\ &= b_0 \sum C_i + b_1 \sum C_i x_i \\ &= b_1 \end{aligned}$$



同理得：
$$Var(\hat{b}_1) = \sum C_i^2 Var(y_i) = \sum C_i^2 S^2$$

由于：
$$\sum C_i^2 = \frac{\sum (x_i - \bar{x})^2}{(\sum (x_j - \bar{x})^2)^2} = \frac{1}{\sum (x_j - \bar{x})^2} = \frac{1}{L_{xx}}$$

于是有：
$$Var(\hat{b}_1) = S_{\hat{b}_1}^2 = \sum C_i^2 S^2 = \frac{S^2}{\sum (x_j - \bar{x})^2} = \frac{S^2}{L_{xx}}$$

同理，可以证明关于 \hat{b}_0 的相应的数字特征的结论。

注意：线性相关系数 r 与一元线性回归系数 \hat{b}_1 符号相同。因为

$$r = \frac{L_{xy}}{\sqrt{L_{xx} L_{yy}}} = \frac{L_{xy}}{L_{xx}} \sqrt{\frac{L_{xx}}{L_{yy}}} = \hat{b}_1 \sqrt{\frac{L_{xx}}{L_{yy}}}$$



在实际计算时，由于总体的方差 σ^2 常常是未知的，则用 S^2 近似地估计总体方差。故可得：

$$S_{\hat{b}_0}^2 = S_{\hat{b}_0}^2 = S^2 \left\{ \frac{1}{n} + \frac{\bar{x}^2}{L_{xx}} \right\} \quad S_{\hat{b}_1}^2 = S_{\hat{b}_1}^2 = \frac{S^2}{L_{xx}}$$

同理，还可以得到回归系数 \hat{b}_0, \hat{b}_1 的协方差为：

$$Cov(\hat{b}_0, \hat{b}_1) = -\frac{\bar{x}}{L_{xx}} S^2$$

上式表明只有在 $\bar{x} = 0$ 时，两回归系数才是相互独立的。

Ø 在基本假设条件下，在各种线性无偏估计量中，由最小二乘估计得到的回归系数具有最小的方差界。



三、方差 σ^2 的估计 S^2

在一元回归方程中，可以证明，未知参数 σ^2 的无偏估计 S^2 可以表示为：

$$S^2 = S^2 = \frac{\sum e_i^2}{n-2}$$

上式中，分子是残差平方和，
分母是自由度，
 n 是样本容量。

由于在估计回归系数时用到 $\sum e_i = 0, \sum e_i x_i = 0$ ，故可以得知自由度为 $n-2$ 。



§4.4 回归模型的检验

一、回归系数的显著性检验

首先，建立原假设： $H_0: b_1=0$

备选假设： $H_1: b_1 \neq 0$

原假设 $b_1=0$

其次，计算回归系数的检验统计量t值

$$t_{b_1} = \frac{\hat{b}_1 - b_1}{S_{\hat{b}_1}} = \frac{\hat{b}_1}{S_{\hat{b}_1}}$$

式中： $S_{\hat{b}_1}$ 是由样本计算的回归系数的估计标准差。

$S_{\hat{b}_1}$ 的计算公式：

$$S_{\hat{b}_1} = \sqrt{\hat{S}^2 / L_{xx}} = \sqrt{S^2 / L_{xx}}$$

最后，根据给定的显著性水平 α ，查表得临界值 $t_{\alpha/2}(n-2)$ ，

如果 $|t| < t_{\alpha/2}(n-2)$ ，则接受原假设 H_0 ，即认为 $b_1=0$ ，说明回归系数与0的差异不显著，回归模型无效。



同理，可检验系数 b_0 。设原假设： $H_0: b_0=0$
备选假设： $H_1: b_0 \neq 0$

再计算检验统计量t值

$$t_{b_0} = \frac{\hat{b}_0 - b_0}{S_{\hat{b}_0}} = \frac{\hat{b}_0}{S_{\hat{b}_0}}$$

$S_{\hat{b}_0}$ 的计算公式：

$$S_{\hat{b}_0} = \sqrt{S^2 \left\{ \frac{1}{n} + \frac{\bar{x}^2}{L_{xx}} \right\}}$$

最后，根据给定的显著性水平 α ，查表得临界值 $t_{\alpha/2}(n-2)$ ，

如果 $|t| < t_{\alpha/2}(n-2)$ ，则接受原假设 H_0 ，即认为 $b_0=0$ ，说明回归系数 b_0 与0的差异不显著。



三、回归方程的显著性检验。

对于线性回归方程整体进行检验判定两变量之间是否存在线性关系。

建立原假设： $H_0: b_1=0$

对立假设： $H_1: b_1 \neq 0$

构造检验统计量：

考虑每个 y_i 与平均值 \bar{y} 之间的差异及总的偏差平方和。先将差异分解后再取平方得：

$$y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$$

$$\dot{a}(y_i - \bar{y})^2 = \dot{a}(\hat{y}_i - \bar{y})^2 + \dot{a}(y_i - \hat{y}_i)^2$$

交叉项为0



$$\dot{a}(y_i - \bar{y})^2 = \dot{a}(\hat{y}_i - \bar{y})^2 + \dot{a}(y_i - \hat{y}_i)^2$$

$\dot{a}(y_i - \bar{y})^2$ 称总偏差平方和，记为SST。

$\dot{a}(\hat{y}_i - \bar{y})^2$ 称回归平方和，记作SSR。反映由X变化引起Y的波动。

$\dot{a}(y_i - \hat{y}_i)^2$ 称残差平方和，记为SSE。反映了随机误差引起的波动。

$$SST = \dot{a}(y_i - \bar{y})^2 = L_{yy}$$

$$SSR = \dot{a}(\hat{y}_i - \bar{y})^2 = \hat{b}_1^2 L_{xx}$$

$$SSE = \dot{a}(y_i - \hat{y}_i)^2 = \dot{a}e_i^2$$



可以证明： $E(SSE) = (n-2)S^2$
 $E(SSR) = S^2 + b_1^2 L_{xx}$

上式表明，当回归系数为0时，回归平方和只反映了随机误差引起的差异。故在原假设成立的条件下，有

$$E(SSE/(n-2)) = S^2; E(SSR) = S^2$$

在原假设成立的条件下还可以证明：

$$\frac{SSR}{S^2} \sim c^2(1); \quad \frac{SSE}{S^2} \sim c^2(n-2)$$



构造F统计量为：

$$F = \frac{SSR}{SSE / (n-2)}$$

根据给定的显著性水平 α 和两个自由度， $df_1=1$ ， $df_2=n-2$ ，查临界值 $F_{\alpha}(1, n-2)$ ，如果 $F > F_{\alpha}(1, n-2)$ ，则拒绝原假设 H_0 ，即回归效果是显著的；反之，回归效果是不显著的。

方差分析表

方差来源	平方和	自由度	均方	F值
回 归	SSR	1	$\frac{SSR}{1} = SSR$	$\frac{SSR}{SSE / (n-2)}$
残 差	SSE	n-2	$\frac{SSE}{n-2}$	
总 和	SST	n-1		

注意：一元线性回归方程的检验，回归系数的检验及相关系数的检验是等价的。



四、拟合优度检验（样本决定系数）

由于总离差平方和 $SST = SSR + SSE$ ，故在等式两边同除总离差平方和得：

$$1 = \frac{SSR}{SST} + \frac{SSE}{SST}$$

则样本决定系数 r^2 定义为：

$$r^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}; \text{ 或 } r^2 = \frac{\hat{\beta}^2 L_{xx}}{L_{yy}} = \frac{L_{xy}^2}{L_{xx} L_{yy}}$$

在一元线性回归模型中，线性相关系数 r 是决定系数 r^2 的平方根。

由定义知，样本决定系数是对回归方程拟合程度的综合测量， r^2 越接近1，拟合程度就越好。反之，样本决定系数越小，模型的拟合程度就越差。当 $r^2=0$ 时，有 $SSR=0$ ，变量 X 对 Y 没有关系。



五、估计标准误差（估计标准误差）

变量 Y 的观察值 y 与回归值 \hat{y} 的差异程度以下式表示：

$$s_y = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-2}} = \sqrt{\frac{\sum e_i^2}{n-2}}$$

上式为最小二乘残差 e_i 标准差。被称估计标准误差。

估计标准误差是 y 值与回归直线变差的测度。可以用来判别回归方程的回归效果。

例：根据表提供的数据，建立某地区居民对某产品的需求量与居民收入的回归方程。

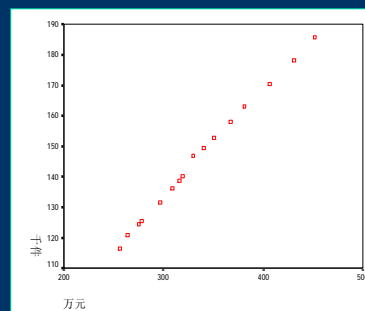
解：令居民收入为自变量 X ，需求量为因变量 Y ，根据表中数据绘制散点图，根据已经给出的数据计算：



序号	Y	X	Y ²	X ²	XY
1	116.5	255.7	13572.25	65382.49	29789.05
2	120.8	263.3	14592.64	69326.89	31806.64
3	124.4	275.4	15475.36	75845.16	34259.76
4	125.5	278.3	15750.25	77450.89	34926.65
5	131.7	296.7	17344.89	88030.89	39075.39
6	136.2	309.3	18550.44	95666.49	42126.66
7	138.7	315.8	19237.69	99729.64	43801.46
8	140.2	318.8	19656.04	101633.44	44695.76
9	146.8	330	21550.24	108900	48444
10	149.6	340.2	22380.16	115736.04	50893.92
11	153	350.7	23409	122990.49	53657.1
12	158.2	367.3	25027.24	134909.29	58106.86
13	163.2	381.3	26634.24	145389.69	62228.16
14	170.5	406.5	29070.25	165242.25	69308.25
15	178.2	430.8	31755.24	185588.64	76768.56
16	185.9	451.5	34558.81	203852.25	83933.85
合计	2339.4	5371.6	348564.74	1855674.54	803822.07



根据数据得到相关图：



解：令居民收入为自变量 X ，需求量为因变量 Y ，根据表中数据绘制散点图，根据已经给出的数据采用最小二乘法计算回归系数：

$$\hat{\beta}_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} = 0.3524$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 27.9123$$

得一元线性回归方程：

$$\hat{y} = 27.9123 + 0.3524 X$$

并计算出各离差平方和为：

$$SST = \sum (y_i - \bar{y})^2 = 5191.35 \quad SSE = \sum (y_i - \hat{y}_i)^2 = 17.57$$

$$SSR = \sum (\hat{y}_i - \bar{y})^2 = \hat{\beta}_1^2 L_{xx} = 5173.78$$



回归系数检验

建立假设： $H_0: \beta_1=0; H_1: \beta_1 \neq 0$

$$\text{计算检验统计量 } t_{\beta_1} = \frac{\hat{\beta}_1}{S_{\hat{\beta}_1}}$$

$$S^2 = s^2 = \frac{\sum e_i^2}{n-2} = \frac{17.57}{14} = 1.225$$

$$S_{\hat{\beta}_1} = \sqrt{\frac{S^2}{\sum (x_i - \bar{x})^2}} = \sqrt{\frac{1.225}{40495.87}} = 0.0055$$

$$t_{\beta_1} = \frac{\hat{\beta}_1}{S_{\hat{\beta}_1}} = \frac{0.3524}{0.0055} = 64.2069$$

根据显著性水平 $\alpha=0.05$ ，查 $t_{0.025}(14)=2.1448$ 。显然有 $t_{\beta_1} > t_{0.025}(14)=2.1448$ ，表明回归系数是显著不等于0。



Ø 回归方程的显著性检验

建立假设: $H_0: \beta_1=0$; $H_1: \beta_1 \neq 0$

$$F = \frac{\hat{\alpha} (\hat{y}_i - \bar{y})^2}{\hat{\alpha} (y_i - \hat{y}_i)^2 / n - 2} = \frac{SSR / 1}{SSE / n - 2}$$

$$Q SSR = 5173.78$$

$$SSE = 17.57$$

$$\therefore F = \frac{5173.78 / 1}{17.57 / 14} = 4122.53$$

根据给定的显著性水平 $\alpha=0.05$, 两个自由度, $df_1=1, df_2=14$. 查临界值 $F_{0.05}(1, 14)=4.60$ 显然有

$$F=4122.53 > F_{0.05}(1, 14)=4.60$$

F检验通过, 可以认为回归方程的回归效果是显著的;



Ø 拟合程度测定

$$\text{计算样本可决系数: } r^2 = \frac{SSR}{SST} = \frac{5173.78}{5191.35} = 0.9966$$

r^2 接近于1, 表明回归直线与样本点的拟合程度很高。

六、回归系数的置信区间

根据区间估计的计算方法, 给定置信水平 $1-\alpha$, 则可以求出回归系数 b_0, b_1 的置信区间分别为:

$$(\hat{b}_0 - \hat{s}_{b_0} t_{\alpha/2}(n-2), \hat{b}_0 + \hat{s}_{b_0} t_{\alpha/2}(n-2))$$

$$(\hat{b}_1 - \hat{s}_{b_1} t_{\alpha/2}(n-2), \hat{b}_1 + \hat{s}_{b_1} t_{\alpha/2}(n-2))$$

由此得出例4.1中回归系数 b_1 的95%的置信区间为:

$$(0.3524 - 2.1448 \cdot 0.0055, 0.3524 + 2.1448 \cdot 0.0055)$$

$$= (0.3406, 0.3642)$$



§4.5 残差分析

一、残差的概念与残差图

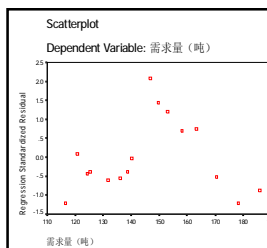
残差指实际观察值与回归值之差 e_i , 是随机项 ε_i 的估计。

$$e_i = y_i - \hat{y}_i = y_i - \hat{b}_0 - \hat{b}_1 x_i$$

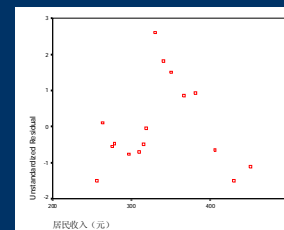
$$e_i = y_i - b_0 - b_1 x_i$$

以自变量X或因变量Y为横轴, 以残差为纵轴可以得到残差图。

如果一个回归模型满足给定的基本假设条件, 则残差应当在 $e=0$ 附近的带状区域内随机排列。否则, 则表明模型存在一定的问題。书上给出了几种典型的残差图的实例。



例題的残差图



二、残差的性质

性质1 $E(e_i) = 0$

性质2 $Var(e_i) = \frac{\hat{\sigma}^2}{n} \cdot \frac{1}{n} - \frac{(x_i - \bar{x})^2}{L_{xx}} \cdot \frac{\hat{\sigma}^2}{n} = (1 - h_{ii}) \hat{\sigma}^2$

式中 $h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{L_{xx}}$ 称为杠杆值 (Leverage Value),

有 $0 < h_{ii} < 1$. 说明 x_i 远离平均值时杠杆值比较大, 则方差较小。

性质3 残差满足约束条件: $\sum \hat{a} e_i = 0; \sum \hat{a} e_i x_i = 0$

即残差项不是独立的。



三、标准化残差与学生化残差

残差分析中的异常值是指超过 $\pm 3s$ 的残差。由于普通的残差的方差不等, 故引入标准化残差和学生化残差。

Ø 标准化残差

$$ZRE_i = \frac{e_i}{\hat{s}}$$

Ø 学生化残差

$$SRE_i = \frac{e_i}{\hat{s} \sqrt{1 - h_{ii}}}$$

标准化残差使残差具有可比性, 一般 $|ZRE| > 3$ 的相应观测值判为异常值。但没有解决方差不等的问题。而学生化残差则进一步解决了标准化残差没有解决的方差不等的问题, 在寻找异常值 $|SRE| > 3$ 方面有优越性。



§4.6 回归方程的应用——预测与控制

一、预测值及预测区间

点预测 预测某一点的值。以自变量的预测值 $X=X_0$ 代入得到 y_0 的预测值 \hat{y}_0 。

$$\hat{y}_0 = \hat{b}_0 + \hat{b}_1 X_0$$

区间预测

由于是随机抽样， \hat{y} 与 y_0 有一定的抽样误差，且不同的样本可以求出不同的回归系数值。因此，有必要对 y_0 值所在的区间作预测。可以证明：

$$\hat{y}_0 - y_0 \sim N\left[0, S^2\left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2}\right)\right]$$



因此， y_0 的概率为 $(1-\alpha)$ 预测区间为：

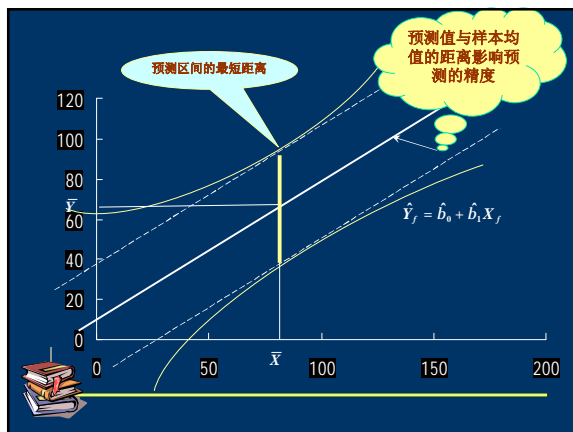
$$[\hat{y}_0 \pm t_{\alpha/2}(n-2) \times \hat{S} \sqrt{1+h_{11}}]$$

其中， h_{11} 为杠杆值，表明预测值远离平均值时，预测区间变大，预测精度降低。

$$h_{11} = \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2}$$

当 n 较大时， t 分布可以近似地用标准正态分布代替，根号下的值约等于1，则有 y_0 的预测区间近似为：

$$\hat{y}_0 \pm Z_{\alpha/2} \hat{S}$$



根据前面食品支出的例题，如果回归模型有效，可以预测当收入达到200单位时，相应的食品支出为：

$$\hat{y}_0 = 9.99 + 0.1802 \times 200 = 46.03$$

而预测值的95%的置信区间为：

$$\begin{aligned} \hat{y}_0 \pm t_{\alpha/2}(n-2) \times \hat{S} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2}} \\ = 46.03 \pm t_{0.025}(13) \times 3.334 \times \sqrt{1 + \frac{1}{15} + \frac{200 - 101.07}{20224.74}} \\ = (42.069, 49.988) \end{aligned}$$

上式的实际意义，当收入达到200单位时，以95%的把握程度认为，食品支出大约在42.069至49.988之间。



二、控制问题

控制问题是预测的逆问题，与预测有着密切的关系，如果要求因变量 y 控制在一定范围内，即范围内以概率 $1-\alpha$ 保证目标值 y 控制在 (T_1, T_2) 范围内，用数学表达式得：

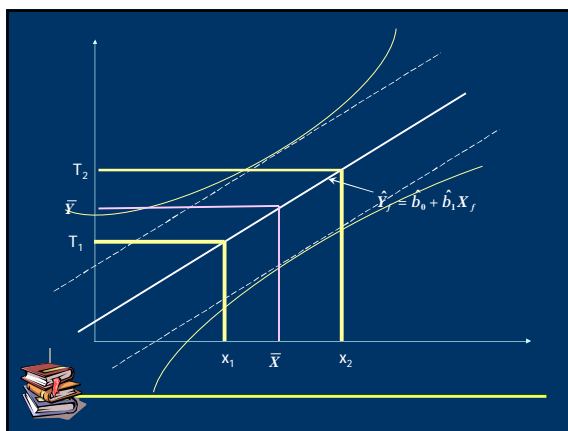
$$P\{T_1 < y < T_2\} = 1 - \alpha$$

从上式中反解出 x 值，就可以得到相应的 x 的取值范围。近似解法为：

$$\begin{aligned} \hat{y}(x) - 2\hat{S} &= \hat{y}(x) - 2\hat{S} > T_1 \\ \hat{y}(x) - 2\hat{S} &= \hat{y}(x) + 2\hat{S} < T_2 \end{aligned}$$

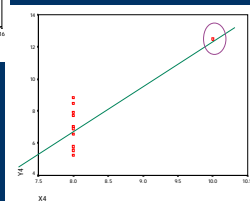
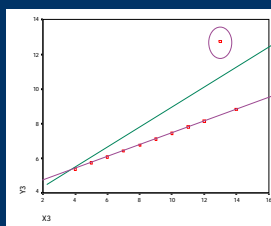
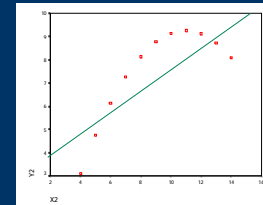
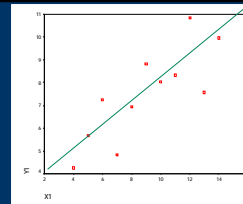
上式中： S 为 σ 的估计值。

2近似等于标准正态分布95%的置信概率。



回归模型的进一步讨论：见例题 P 101.

第一组		第二组		第三组		第四组	
X	Y	X	Y	X	Y	X	Y
4	4.26	4	3.1	4	5.39	8	6.58
5	5.68	5	4.74	5	5.73	8	5.76
6	7.24	6	6.13	6	6.08	8	7.71
7	4.82	7	7.26	7	6.44	8	8.84
8	6.95	8	8.14	8	6.77	8	8.47
9	8.81	9	8.77	9	7.11	8	7.04
10	8.04	10	9.14	10	7.46	8	5.25
11	8.33	11	9.26	11	7.81	8	5.56
12	10.84	12	9.13	12	8.15	8	7.91
13	7.58	13	8.74	13	12.74	8	6.89
14	9.96	14	8.1	14	8.84	10	12.5



Model Summary ^b					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Durbin-Watson
1	.816 ^a	.667	.629	1.23660	2.788

a. Predictors: (Constant), X1

Model Summary ^b					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Durbin-Watson
1	.816 ^a	.666	.629	1.23721	.385

a. Predictors: (Constant), X2

Model Summary ^b					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Durbin-Watson
1	.816 ^a	.666	.629	1.23618	2.638

a. Predictors: (Constant), X3

Model Summary ^b					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Durbin-Watson
1	.817 ^a	.667	.630	1.23570	1.296

a. Predictors: (Constant), X4

b. Dependent Variable: Y4

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	27.510	1	27.510	17.990	.002 ^a
Residual	13.763	9	1.529		
Total	41.273	10			

a. Predictors: (Constant), X1

b. Dependent Variable: Y1

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	27.500	1	27.500	17.966	.002 ^a
Residual	13.776	9	1.531		
Total	41.276	10			

a. Predictors: (Constant), X2

b. Dependent Variable: Y2

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	27.430	1	27.430	17.950	.002 ^a
Residual	13.753	9	1.528		
Total	41.183	10			

a. Predictors: (Constant), X3

b. Dependent Variable: Y3

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	27.490	1	27.490	18.003	.002 ^a
Residual	13.742	9	1.527		
Total	41.232	10			

a. Predictors: (Constant), X4

b. Dependent Variable: Y4

Coefficients ^a					
Model		Unstandardized Coefficients		Standardized Coefficients	Sig.
		B	Std. Error	Beta	
1	(Constant)	3.000	1.125		.026
	X1	.500	.118	.816	.002

a. Dependent Variable: Y1

Coefficients ^a					
Model		Unstandardized Coefficients		Standardized Coefficients	Sig.
		B	Std. Error	Beta	
1	(Constant)	3.001	1.125		.026
	X2	.500	.118	.816	.002

a. Dependent Variable: Y2

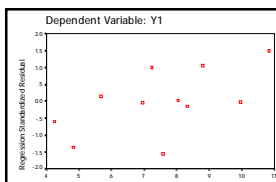


Coefficients ^a					
Model		Unstandardized Coefficients		Standardized Coefficients	Sig.
		B	Std. Error	Beta	
1	(Constant)	3.008	1.124		.025
	X3	.499	.118	.816	.002

a. Dependent Variable: Y3

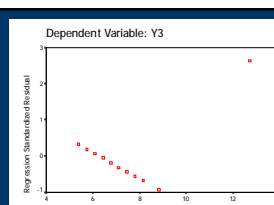
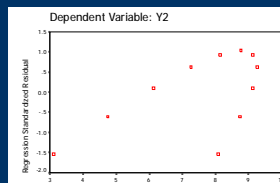
Coefficients ^a					
Model		Unstandardized Coefficients		Standardized Coefficients	Sig.
		B	Std. Error	Beta	
1	(Constant)	-14.995	5.315		.020
	X4	2.750	.648	.817	.002

a. Dependent Variable: Y4



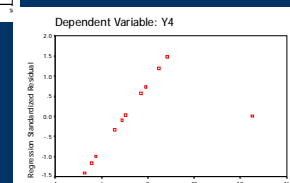
第一组数据的残差图有序列自相关的嫌疑。

第二组数据的残差图有曲线形式



第三组和第四组数据中都明显存在一个异常值。

这四个模型都需要进一步的分析。



一元线性回归应用步骤:

1. 首先根据研究目的确定因变量Y和自变量X;
2. 收集数据后, 首先进行相关分析, 大致确定回归方程(线性或通过变换得到线性模型);
3. 建立一元线性回归模型, 根据样本建立数据表;
4. 用SPSS统计软件对样本数据进行分析, 得到相应的输出结果。
5. 对照输出结果对回归模型进行检验: 包括拟合优度的检验(R平方)、回归系数的检验P值、回归方程的检验(方差分析表)、残差分析、寻找异常值等。
6. 通过检验的回归模型称为有效方程, 可以在实际中应用, 也可以用于预测和控制问题。



§4.7 可化为线性的曲线函数

一、几种可化为线性的曲线函数

1. 抛物线函数 (Quadratic 二次函数) $y = a + bX + cX^2$ 经变换得:

$$y = a + bX_1 + cX_2$$

其中: $X_1 = X$, $X_2 = X^2$, 故方程变换后成为一个线性方程。

2. 双曲线函数 (Inverse 逆函数) $y = a + b/X$ 做变换得:

$$y = a + bX_1 \quad X_1 = 1/X$$

3. 幂函数 (Power) $y = aX^b$

$$\ln y = \ln a + b \ln X$$

即: $y^c = a^c + bX^c$ 其中 $y^c = \ln y$, $a^c = \ln a$, $X^c = \ln X$

参数b是y对于自变量X的弹性, 是X变动1%时, y变动的百分比。



4. 指数函数 (Exponential)

$$y = ae^{bx} \quad \ln y = \ln a + bx$$

5. 逆函数 (Inverse)

$$y = b_0 + (b_1/x) \quad y = b_0 + b_1 x^c \quad x^c = 1/x$$

6. 对数函数 (Logarithmic)

$$y = a + b \log x$$

7. 逻辑函数 (Logistic)

$$y = \frac{L}{1 + ae^{-bx}} \quad (L, a, b > 0)$$

$$ae^{-bx} = \frac{L}{y} - 1 \quad \ln\left(\frac{L}{y} - 1\right) = \ln a - bx$$

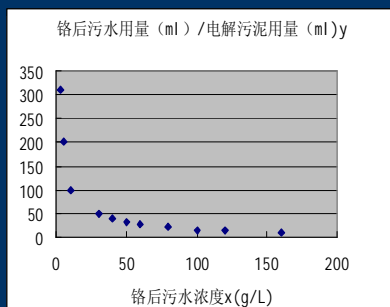


实例分析：某厂表面处理车间试验将铬后污水同电解污泥混合生成无毒溶液。但实际排污的浓度不完全相同，且一定浓度的定量铬后污水只有同定量的电解污泥混合后，才能反应完全。现通过试验，找出铬后污水用量与电解污泥用量之比对于铬后污水浓度之间的关系，试验数据如下：

序号	铬后污水浓度x (g/L)	铬后污水用量 (ml) / 电解污泥用量 (ml) y
1	3	310
2	5	200
3	10	100
4	30	49
5	40	40
6	50	32
7	60	28
8	80	23
9	100	16
10	120	14
11	160	10



解：首先根据数据得到散点图如下：



根据散点图的形状，试用幂函数进行曲线拟合。

假设曲线为幂函数，

$$y = aX^b \quad \text{对数变换} \quad \ln y = \ln a + b \ln X$$

$$\text{取} y^c = \ln y \quad x^c = \ln x \quad a^c = \ln a$$

以变换后计算得：

$$L_{x^c y^c} = 17.462, \quad L_{y^c y^c} = 12.003, \quad L_{x^c y^c} = -14.429$$

由此得回归系数为：

$$\hat{b} = \frac{L_{x^c y^c}}{L_{x^c x^c}} = -0.8263, \quad a^c = 6.6417$$

$$\hat{a} = e^{a^c} = 766.4$$

$$\text{回归方程为: } Y = 766.4x^{-0.8263}$$



对于线性函数应用相关系数表达其线性相关程度，对于非线性函数也可以用相关指数来反映一元非线性回归方程的优劣。

相关指数的定义为：

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{a}_i^2 e_i^2}{L_{yy^c}}$$

由例中的计算得相关指数为：

$$\hat{a}_i e_i^2 = \hat{a} (y_i - \hat{y}_i)^2$$

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{a}_i^2 e_i^2}{L_{yy^c}} = 1 - \frac{0.080}{12.003} = 0.993$$

$$L_{yy} = \hat{a} (y_i - \bar{y})^2$$

由此可见，曲线的拟合程度是很高的。可以用SPSS作多种函数的假设。



将变量X、Y进行对数变换，应用SPSS软件进行线性回归分析得：

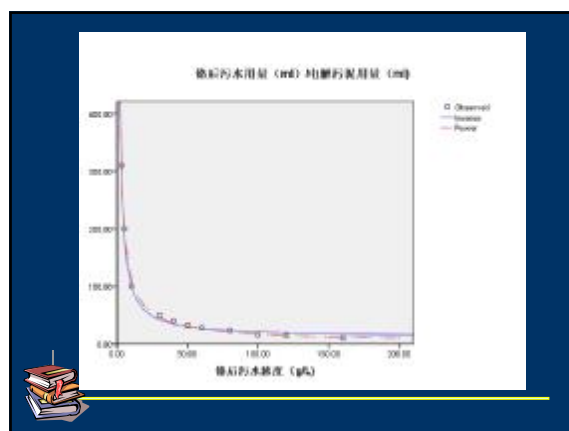
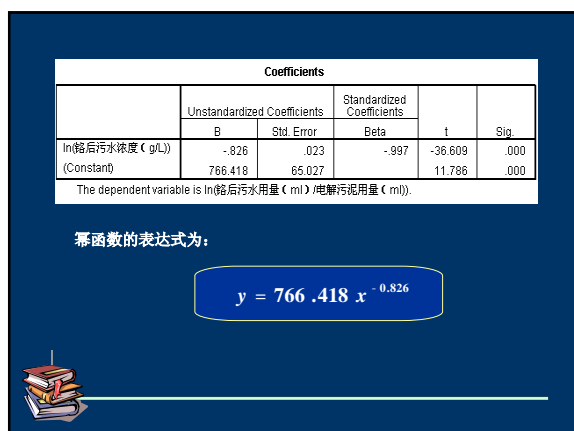
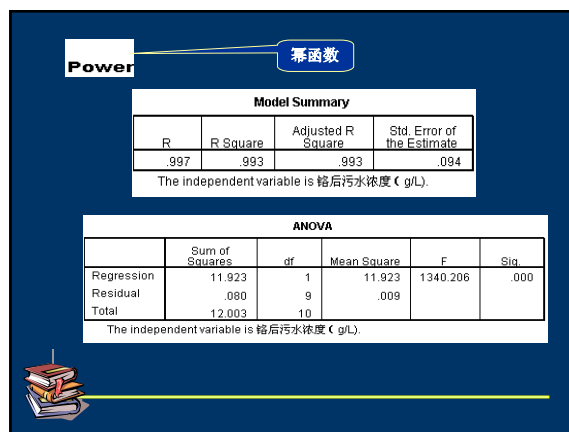
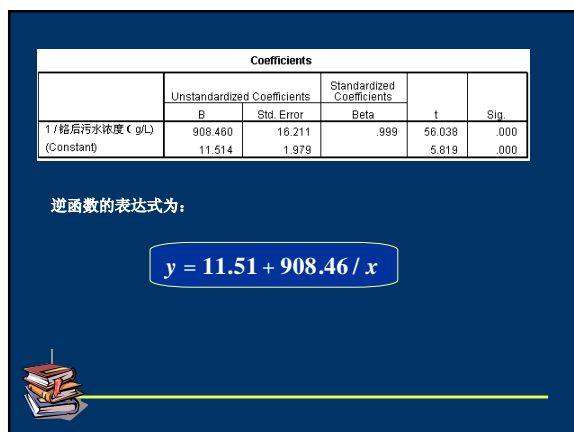
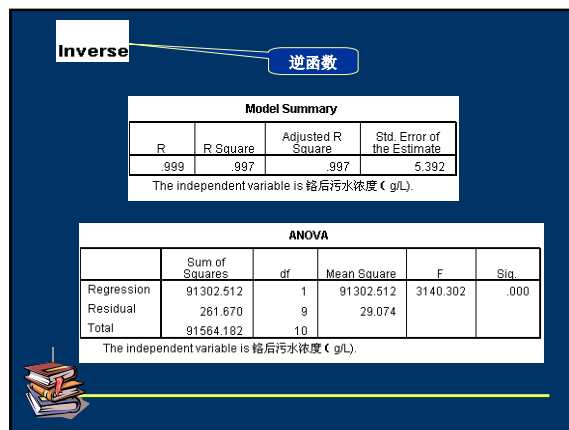
Model Summary ^a					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Durbin-Watson
1	.997 ^a	.993	.993	.09432	1.239

a. Predictors: (Constant), LNX
b. Dependent Variable: LNY

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	11.923	1	11.923	1340.206	.000 ^b
Residual	.080	9	.009		
Total	12.003	10			

a. Predictors: (Constant), LNX
b. Dependent Variable: LNY





一元线性回归小结:

1. 一元线性回归模型的数学公式(总体与样本), 要求掌握如何利用样本建立回归模型, 注意理解高斯假设的意义。

2. 理解线性回归的意义, 明确回归分析与相关分析的关系。

3. 掌握回归系数的求解方法——最小二乘法。

$$\hat{b}_1 = \frac{L_{xy}}{L_{xx}} = \frac{n \sum X_i Y_i - \sum X_i \sum Y_i}{n \sum X_i^2 - (\sum X_i)^2}, \quad \hat{b}_0 = \frac{\sum Y_i}{n} - \hat{b}_1 \frac{\sum X_i}{n} = \bar{Y} - \hat{b}_1 \bar{X}$$

4. 估计总体方差 σ^2 :

计算中常用 S^2 估计 σ^2 。

$$S^2 = \frac{\sum e_i^2}{n-2}$$



5. \hat{b}_0, \hat{b}_1 的分布

由假设得出Y服从正态分布, 且有: $y \sim N(b_0 + b_1 x, \sigma^2)$

可以证明, $\hat{b}_0 \sim N(b_0, S_{b_0}^2)$; $\hat{b}_1 \sim N(b_1, S_{b_1}^2)$

其中:

$$S_{b_0}^2 = S^2 \left\{ \frac{1}{n} + \frac{\bar{X}^2}{L_{xx}} \right\}, \quad S_{b_1}^2 = \frac{S^2}{L_{xx}}$$

由样本得出回归系数方差的无偏估计为:

$$\hat{S}_{b_0}^2 = S^2 \left\{ \frac{1}{n} + \frac{\bar{X}^2}{L_{xx}} \right\}, \quad \hat{S}_{b_1}^2 = \frac{S^2}{L_{xx}}$$



6. 模型检验

○ **拟合程度的检验:** 利用可决系数 r^2 对回归方程拟合程度的综合测量, 可决系数越大, 模型的拟合程度就越高, 反之, 可决系数越小, 模型的拟合程度就越差:

$$r^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

○ **回归系数的检验 (t 检验)**

提出假设: $H_0: b_1=0$; $H_1: b_1 \neq 0$

计算检验统计量值

$$t_{b_1} = \frac{\hat{b}_1}{S_{\hat{b}_1}} = \frac{\hat{b}_1}{S / \sqrt{L_{xx}}}$$

对于给定的显著性水平 α , 查 t 分布表确定临界值(或用 P 值)。要注意当假设是双边检验比较临界值与 t 值后作出判断。



○ 回归方程的检验 (F 检验)

提出假设: $H_0: b_1=0$; $H_1: b_1 \neq 0$

计算离差平方和, 列出方差分析表如下:

离差名称	离差平方和	自由度	均方差	F 值
回归平方和	$SSR = \hat{\hat{a}}(\sum Y_i - \bar{Y})^2$	1	$\frac{SSR}{SSE/(n-2)}$	$F = \frac{SSR}{SSE/(n-2)}$
残差平方和	$SSE = \sum e_i^2$	$n-2$		
总离差平方和	$SST = \hat{\hat{a}}(\sum Y_i - \bar{Y})^2$	$n-1$		

对于给定的显著性水平 α , 查 F 分布表确定临界值(或用 P 值), 当 $F > F_{\alpha}$ 时, 拒绝原假设。



○ 残差分析

用因变量或者自变量做横轴, 残差(标准化残差, 学生化残差等)为纵轴, 得到残差图, 通过残差图是否随机排列和残差值的范围考察回归模型是否存在异方差性, 是否存在异常值等。

7. 利用回归方程进行预测

○ **预测公式:**

由简单回归模型得出基本预测公式为:

$$\hat{y}_0 = \hat{b}_0 + \hat{b}_1 X_0$$

• **点预测:** 即给定 X_0 得到 Y_0 的预测值。

注意: 内插与外推预测的区别。



• Y_0 的区间预测:

在总体标准差 σ^2 未知的情况下, 用其无偏估计 S^2 来代替, 可以证明:

$$(y_0 - \hat{y}_0) / S_{e_0} \sim t(n-2)$$

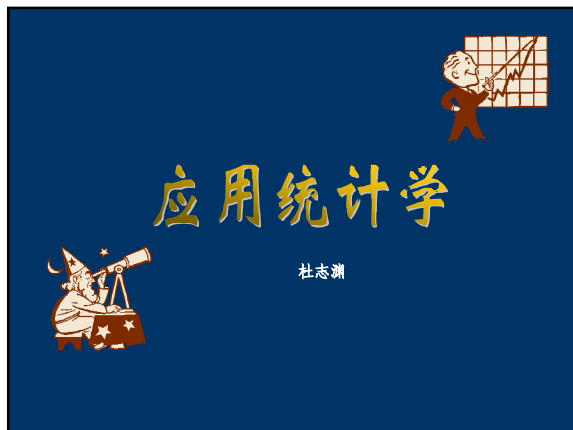
由置信区间的确定方法得 y_0 的置信度为 $(1-\alpha)$ 的置信区间为

$$(\hat{y}_0 \pm t_{\alpha/2}(n-2) \cdot S_{e_0})$$

• **预测的标准误 S_{e_0} :** 在标准假定下, 有

$$S_{e_0} = S \sqrt{\frac{1}{n} + \frac{(X_0 - \bar{X})^2}{L_{xx}}}$$





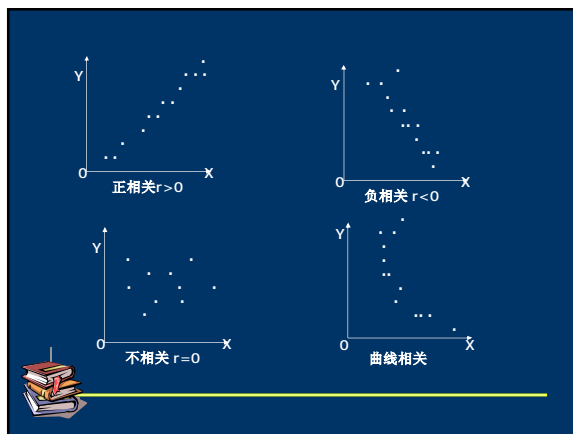
第四章 一元线性回归

§ 4.1 一元线性回归模型

一、变量之间的关系

函数关系 变量之间可以用数学公式表示。如圆的半径R与面积S等。

统计关系 两变量之间有一定依存关系，但没有严格的对应关系。如人的年龄和血压，身高和体重，储蓄额与居民收入等都是统计关系或相关关系。相关关系可以通过相关图表示出来。相关关系有线性相关和非线性相关（曲线相关）。



二、两变量之间的线性相关系数

1、定义：变量X与Y之间的线性相关程度可以用简单相关系数来度量。计算公式为：

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}} = \frac{L_{xy}}{\sqrt{L_{xx} L_{yy}}}$$

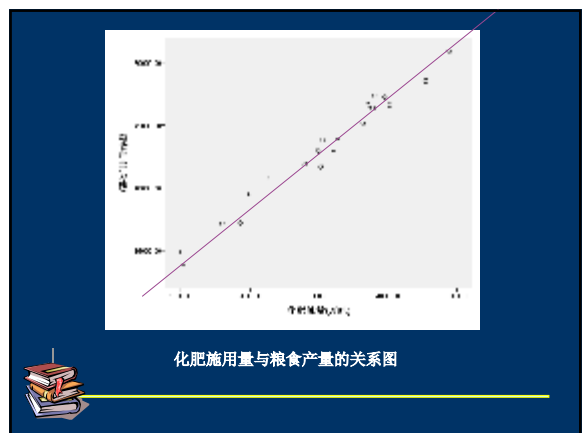
由样本数据得：

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

其中 $-1 \leq r \leq 1$ ，且 $r > 0$ 时称为正相关， $r < 0$ 时为负相关， $r = 0$ 时，两变量间不存在线性相关关系。

例 4.1 讨论化肥施用量与粮食产量之间的关系 数据如下表：

化肥施用量 x(万吨)	4541.05	3637.87	2287.49	3056.89	4883.7	3779.3	4021.09
粮食产量 y(万吨)	48526.69	45110.87	40753.79	43824.58	50890.11	46370.88	46577.91
化肥施用量 x(万吨)	2989.06	3021.9	3953.97	3212.13	3804.76	1598.28	1998.56
粮食产量 y(万吨)	42947.44	41673.21	47244.34	43061.53	47336.78	37127.89	39515.07
化肥施用量 x(万吨)	3710.56	3269.03	1017.12	1864.23	2797.24	1034.09	
粮食产量 y(万吨)	46598.04	44020.92	34866.91	37184.14	41864.77	33717.78	



2、相关系数的显著性检验

对样本相关系数 r 的相关性检验是确定两变量间的线性关系是否显著。

相关系数的检验过程：

- 1、假设 $H_0: r=0$; $H_1: r \neq 0$
- 2、根据给出的样本数据，计算样本相关系数 r 和统计量 t ；在原假设成立的条件下，检验统计量 t 服从 t 分布。

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \quad t \sim t(n-2)$$

- 3、根据给定的显著性水平 α ，查临界值 $t_{\alpha/2}(n-2)$ ；
- 4、若 $t \geq t_{\alpha/2}(n-2)$ ，说明 r 在统计上是显著的，即两个变量之间存在线性关系，若 $t < t_{\alpha/2}(n-2)$ ，则两个变量之间不存在线性相关。

r 为两变量的相关系数的真值



SPSS统计软件得到的化肥施用量与粮食产量的相关性输出表：

Correlations		
	化肥施用量(万吨)	粮食产量(万吨)
化肥施用量(万吨)	1	.989*
	Sig. (2-tailed)	.000
	N	20
粮食产量(万吨)	.989*	1
	Sig. (2-tailed)	.000
	N	20

**. Correlation is significant at the 0.01 level (2-tailed).

注意：在相关分析中两个变量的地位是平等的。



三、一元线性回归的数学模型

如果两个变量之间存在线性关系，设一个变量为自变量，另一个是因变量，则两个变量之间的关系可以用一元线性回归模型表达。

设自变量 X 为确定性变量（解释变量），因变量 Y （被解释变量）为随机变量，两者之间的数学结构式为

$$y = b_0 + b_1x + e$$

式中： b_0, b_1 是回归系数， e 是随机项，表示除了变量 X 之外其他因素对变量 Y 的影响。

线性方程由两部分组成，一部分由 X 的变化引起，另一部分是由随机因素引起。



基本假设

假设1：随机项 e 服从正态分布 $N(0, \sigma^2)$ ，即：

$$E(e) = 0, \quad Var(e) = \sigma^2$$

假设2：随机项 e_i 之间是相互独立的，并具有相同的方差，即：

$$Cov(e_i, e_j) = \begin{cases} \sigma^2 & i = j \\ 0 & i \neq j \end{cases} \quad i, j = 1, 2, \dots, n$$

假设3：样本数据是相互独立的。



由假设知，随机变量 y 也服从正态分布，由公式 $y = b_0 + b_1x + e$ 得出：

$$E(y) = b_0 + b_1x, \quad Var(y) = \sigma^2$$

在一般情况下，从研究的总体中抽取一个样本观察值 (x_i, y_i) ， $i=1, 2, \dots, n$ ，对于样本 X, Y 的每一组数，有

$$y_i = b_0 + b_1x_i + e_i \quad i = 1, 2, \dots, n$$

由假设条件知， $e_i \sim N(0, \sigma^2)$ ，且 $E(e_i) = 0, Var(e_i) = \sigma^2, i = 1, 2, \dots, n$

推导出观察值 $y_i (i=1, 2, \dots, n)$ 也是相互独立的正态随机变量，且

$$E(y_i) = b_0 + b_1x_i, \quad Var(y_i) = \sigma^2$$

线性模型在平均意义上表达了变量 Y 与 X 的统计规律性。



回归模型的矩阵表达式：

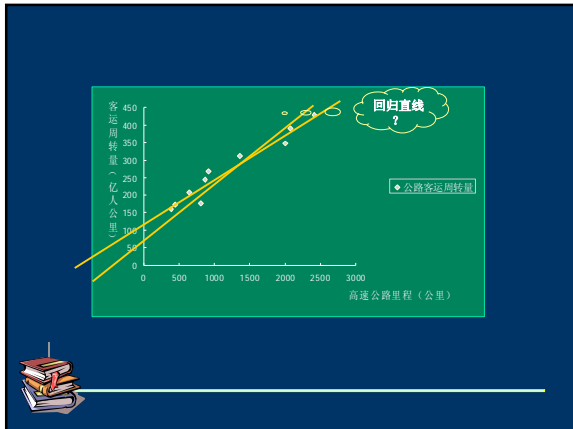
$$\begin{cases} Y = Xb + e \\ E(e) = 0 \\ Var(e) = \sigma^2 I_n \end{cases}$$

$$\text{其中: } Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad b = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}, \quad e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}, \quad X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

线性回归的任务是：通过样本观察值对回归系数进行估计，求出 b_0, b_1 的估计值 \hat{b}_0, \hat{b}_1 ，得出一元线性回归方程：

$$\hat{y} = \hat{b}_0 + \hat{b}_1x$$





§4.2 回归系数的最小二乘估计

一、普通最小二乘估计 (OLSE)

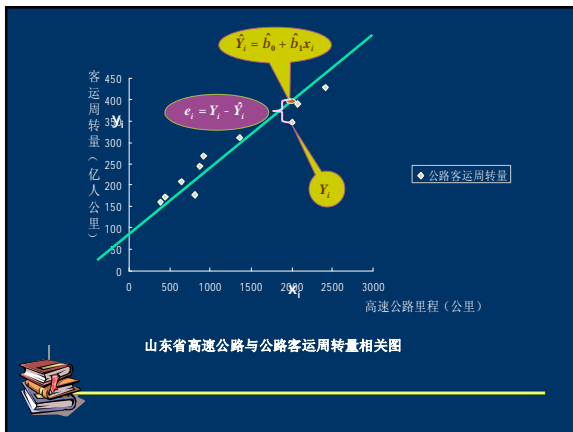
回归系数的估计是依赖于最小二乘法的基本思想, 考虑观察值 y 与回归值 \hat{y}_i 的差(称为残差, 记为 e_i). 通过使残差平方和为最小来估计回归系数。通过样本得到残差平方和为:

$$Q(b_0, b_1) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2 = \sum_{i=1}^n e_i^2$$

求极值得:

$$\begin{cases} \frac{\partial Q}{\partial b_0} = -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i) = 0 \\ \frac{\partial Q}{\partial b_1} = -2 \sum_{i=1}^n x_i (y_i - b_0 - b_1 x_i) = 0 \end{cases} \quad \leftarrow \begin{cases} \hat{a} e_i = 0 \\ \hat{a} e_i x_i = 0 \end{cases}$$

$\hat{y}_i = \hat{b}_0 + \hat{b}_1 x_i$



由上面方程组可得:

$$\begin{cases} \hat{a} y_i - n \hat{b}_0 - \hat{b}_1 \hat{a} x_i = 0 \\ \hat{a} x_i y_i - \hat{b}_0 \hat{a} x_i - \hat{b}_1 \hat{a} x_i^2 = 0 \end{cases}$$

$$\hat{b}_1 = \frac{L_{xy}}{L_{xx}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x} = \frac{1}{n} \sum_{i=1}^n y_i - \hat{b}_1 \frac{1}{n} \sum_{i=1}^n x_i$$

这种方法称普通最小二乘估计(OLSE), 是线性回归方程中回归系数求解的基本方法。这样由样本估计得到的回归方程称为一元线性经验回归方程 记为:

$$\hat{Y} = \hat{b}_0 + \hat{b}_1 X$$

例: 根据恩格尔定律得出食品支出 Y 与实际收入 X 的关系的一元线性回归模型来反映.用表的形式列出原始数据 X, Y 的值及相应的计算数据.见表

解: 根据计算表中的数据得出回归系数为:

$$\hat{b}_1 = \frac{15 \cdot 44632 - 1516 \cdot 423}{15 \cdot 163634 - 1516^2} = 0.1802$$

$$\hat{b}_0 = \frac{423}{15} - 0.1802 \cdot \frac{1516}{15} = 9.9872$$

所求的经验回归方程为:

$$\hat{y}_i = 9.99 + 0.1802 x_i$$

回归方程的**实际意义**是: 当收入每增加一个单位时, 食品支出会增加0.18单位, 即使在收入为0的情况下, 食品支出依然需要9.99单位。

续

编号	X	Y	XY	X ²	Y ²
1	102	27	2754	10404	729
2	96	26	2496	9216	676
3	97	25	2425	9409	625
4	102	28	2856	10404	784
5	91	27	2457	8281	729
6	158	36	5688	24964	1296
7	54	19	1026	2916	361
8	83	26	2158	6889	676
9	123	31	3813	15129	961
10	106	31	3286	11236	961
11	129	34	4386	16641	1156
12	138	38	5244	19044	1444
13	81	27	2187	6561	729
14	92	28	2576	8464	784
15	64	20	1280	4096	400
合计	1516	423	44632	163634	12311

返回

§4.3 回归系数的性质及统计意义

一、 \hat{b}_0, \hat{b}_1 是 y_i 的线性组合。

例如:
$$\hat{b}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{L_{xx}} = \frac{\sum (x_i - \bar{x})y_i}{L_{xx}} = \sum \frac{(x_i - \bar{x})}{L_{xx}} y_i$$

$$\hat{b}_1 = \sum C_i y_i \quad C_i = \frac{x_i - \bar{x}}{L_{xx}} \quad i = 1, 2, \dots, n$$

同理得:
$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x} = \sum \left(\frac{1}{n} - C_i \bar{x} \right) y_i$$



二、 \hat{b}_0, \hat{b}_1 服从正态分布, 且最小二乘估计的回归系数具有无偏性。

由于因变量 y_i 服从正态分布, \hat{b}_0 和 \hat{b}_1 是 y_i 的线性组合, 故也服从正态分布。

即有:
$$E(\hat{b}_0) = b_0 \quad S_{\hat{b}_0}^2 = S^2 \left\{ \frac{1}{n} + \frac{\bar{x}^2}{L_{xx}} \right\} = S^2 \left\{ \frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right\}$$

$$E(\hat{b}_1) = b_1 \quad S_{\hat{b}_1}^2 = \frac{S^2}{L_{xx}} = \frac{S^2}{\sum (x_i - \bar{x})^2}$$

$$\hat{b}_0 \sim N\left(b_0, \left(\frac{1}{n} + \frac{\bar{x}^2}{L_{xx}}\right) S^2\right); \quad \hat{b}_1 \sim N\left(b_1, \frac{S^2}{L_{xx}}\right)$$



证明: $E(\hat{b}_1) = b_1$

因为
$$\hat{b}_1 = C_1 y_1 + C_2 y_2 + \dots + C_n y_n$$

由于
$$\sum C_i = \sum \frac{x_i - \bar{x}}{\sum (x_j - \bar{x})^2} = 0$$

$$\sum C_i x_i = \sum C_i (x_i - \bar{x}) = 1$$

$$\begin{aligned} E(\hat{b}_1) &= \sum C_i E(y_i) \\ &= \sum C_i (b_0 + b_1 x_i) \\ &= b_0 \sum C_i + b_1 \sum C_i x_i \\ &= b_1 \end{aligned}$$



同理得:
$$Var(\hat{b}_1) = \sum C_i^2 Var(y_i) = \sum C_i^2 S^2$$

由于
$$\sum C_i^2 = \frac{\sum (x_i - \bar{x})^2}{(\sum (x_j - \bar{x})^2)^2} = \frac{1}{\sum (x_j - \bar{x})^2} = \frac{1}{L_{xx}}$$

于是有
$$Var(\hat{b}_1) = S_{\hat{b}_1}^2 = \sum C_i^2 S^2 = \frac{S^2}{\sum (x_j - \bar{x})^2} = \frac{S^2}{L_{xx}}$$

同理, 可以证明关于 \hat{b}_0 的相应的数字特征的结论。

注意: 线性相关系数 r 与一元线性回归系数 \hat{b}_1 符号相同。因为

$$r = \frac{L_{xy}}{\sqrt{L_{xx} L_{yy}}} = \frac{L_{xy}}{L_{xx}} \sqrt{\frac{L_{xx}}{L_{yy}}} = \hat{b}_1 \sqrt{\frac{L_{xx}}{L_{yy}}}$$



在实际计算时, 由于总体的方差 σ^2 常常是未知的, 则用 S^2 近似地估计总体方差。故可得:

$$S_{\hat{b}_0}^2 = S_{\hat{b}_0}^2 = S^2 \left\{ \frac{1}{n} + \frac{\bar{x}^2}{L_{xx}} \right\} \quad S_{\hat{b}_1}^2 = S_{\hat{b}_1}^2 = \frac{S^2}{L_{xx}}$$

同理, 还可以得到回归系数 \hat{b}_0, \hat{b}_1 的协方差为:

$$Cov(\hat{b}_0, \hat{b}_1) = -\frac{\bar{x}}{L_{xx}} S^2$$

上式表明只有在 $\bar{x} = 0$ 时, 两回归系数才是相互独立的。

Ø 在基本假设条件下, 在各种线性无偏估计量中, 由最小二乘估计得到的回归系数具有最小的方差界。



三、方差 σ^2 的估计 S^2

在一元回归方程中, 可以证明, 未知参数 σ^2 的无偏估计 S^2 可以表示为:

$$S^2 = S^2 = \frac{\sum e_i^2}{n-2}$$

上式中, 分子是残差平方和,
分母是自由度,
 n 是样本容量。

由于在估计回归系数时用到 $\sum e_i = 0, \sum e_i x_i = 0$, 故可以得知自由度为 $n-2$ 。



§4.4 回归模型的检验

一、回归系数的显著性检验

首先，建立原假设： $H_0: b_1=0$

备选假设： $H_1: b_1 \neq 0$

原假设 $b_1=0$

其次，计算回归系数的检验统计量t值

$$t_{b_1} = \frac{\hat{b}_1 - b_1}{\hat{S}_{\hat{b}_1}} = \frac{\hat{b}_1}{\hat{S}_{\hat{b}_1}}$$

式中： $\hat{S}_{\hat{b}_1}$ 是由样本计算的回归系数的估计标准差。

$\hat{S}_{\hat{b}_1}$ 的计算公式：

$$\hat{S}_{\hat{b}_1} = \sqrt{\hat{S}^2 / L_{xx}} = \sqrt{S^2 / L_{xx}}$$

最后，根据给定的显著性水平 α ，查表得临界值 $t_{\alpha/2}(n-2)$ ，

如果 $|t| < t_{\alpha/2}(n-2)$ ，则接受原假设 H_0 ，即认为 $b_1=0$ ，说明回归系数与0的差异不显著，回归模型无效。



同理，可检验系数 b_0 。设原假设： $H_0: b_0=0$
备选假设： $H_1: b_0 \neq 0$

再计算检验统计量t值

$$t_{b_0} = \frac{\hat{b}_0 - b_0}{\hat{S}_{\hat{b}_0}} = \frac{\hat{b}_0}{\hat{S}_{\hat{b}_0}}$$

$\hat{S}_{\hat{b}_0}$ 的计算公式：

$$\hat{S}_{\hat{b}_0} = \sqrt{S^2 \left\{ \frac{1}{n} + \frac{\bar{x}^2}{L_{xx}} \right\}}$$

最后，根据给定的显著性水平 α ，查表得临界值 $t_{\alpha/2}(n-2)$ ，

如果 $|t| < t_{\alpha/2}(n-2)$ ，则接受原假设 H_0 ，即认为 $b_0=0$ ，说明回归系数 b_0 与0的差异不显著。



三、回归方程的显著性检验。

对于线性回归方程整体进行检验判定两变量之间是否存在线性关系。

建立原假设： $H_0: b_1=0$

对立假设： $H_1: b_1 \neq 0$

构造检验统计量：

考虑每个 y_i 与平均值 \bar{y} 之间的差异及总的偏差平方和。先将差异分解后再取平方得：

$$y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$$

$$\dot{a}(y_i - \bar{y})^2 = \dot{a}(\hat{y}_i - \bar{y})^2 + \dot{a}(y_i - \hat{y}_i)^2$$

交叉项为0



$$\dot{a}(y_i - \bar{y})^2 = \dot{a}(\hat{y}_i - \bar{y})^2 + \dot{a}(y_i - \hat{y}_i)^2$$

$\dot{a}(y_i - \bar{y})^2$ 称总偏差平方和，记为SST。

$\dot{a}(\hat{y}_i - \bar{y})^2$ 称回归平方和，记作SSR。反映由X变化引起Y的波动。

$\dot{a}(y_i - \hat{y}_i)^2$ 称残差平方和，记为SSE。反映了随机误差引起的波动。

$$SST = \dot{a}(y_i - \bar{y})^2 = L_{yy}$$

$$SSR = \dot{a}(\hat{y}_i - \bar{y})^2 = \hat{b}_1^2 L_{xx}$$

$$SSE = \dot{a}(y_i - \hat{y}_i)^2 = \dot{a}e_i^2$$



可以证明： $E(SSE) = (n-2)S^2$

$E(SSR) = S^2 + b_1^2 L_{xx}$

上式表明，当回归系数为0时，回归平方和只反映了随机误差引起的差异。故在原假设成立的条件下，有

$$E(SSE/(n-2)) = S^2; E(SSR) = S^2$$

在原假设成立的条件下还可以证明：

$$\frac{SSR}{S^2} \sim c^2(1); \quad \frac{SSE}{S^2} \sim c^2(n-2)$$



构造F统计量为：

$$F = \frac{SSR}{SSE / (n-2)}$$

根据给定的显著性水平 α 和两个自由度， $df_1=1$ ， $df_2=n-2$ ，查临界值 $F_{\alpha}(1, n-2)$ ，如果 $F > F_{\alpha}(1, n-2)$ ，则拒绝原假设 H_0 ，即回归效果是显著的；反之，回归效果是不显著的。

方差分析表

方差来源	平方和	自由度	均方	F值
回 归	SSR	1	$\frac{SSR}{1} = SSR$	$\frac{SSR}{SSE / (n-2)}$
残 差	SSE	n-2	$\frac{SSE}{n-2}$	
总 和	SST	n-1		

注意：一元线性回归方程的检验，回归系数的检验及相关系数的检验是等价的。



四、拟合优度检验（样本决定系数）

由于总离差平方和 $SST = SSR + SSE$ ，故在等式两边同除总离差平方和得：

$$1 = \frac{SSR}{SST} + \frac{SSE}{SST}$$

则样本决定系数 r^2 定义为：

$$r^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}; \text{ 或 } r^2 = \frac{\hat{\beta}^2 L_{xx}}{L_{yy}} = \frac{L_{xy}^2}{L_{xx} L_{yy}}$$

在一元线性回归模型中，线性相关系数 r 是决定系数 r^2 的平方根。

由定义知，样本决定系数是对回归方程拟合程度的综合测量， r^2 越接近1，拟合程度就越好。反之，样本决定系数越小，模型的拟合程度就越差。当 $r^2=0$ 时，有 $SSR=0$ ，变量 X 对 Y 没有关系。



五、估计标准误差（估计标准误差）

变量 Y 的观察值 y 与回归值 \hat{y} 的差异程度以下式表示：

$$s_y = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-2}} = \sqrt{\frac{\sum \hat{e}_i^2}{n-2}}$$

上式为最小二乘残差 \hat{e} 的标准差。被称估计标准误差。

估计标准误差是 y 值与回归直线变差的测度。可以用来判别回归方程的回归效果。

例：根据表提供的数据，建立某地区居民对某产品的需求量与居民收入的回归方程。

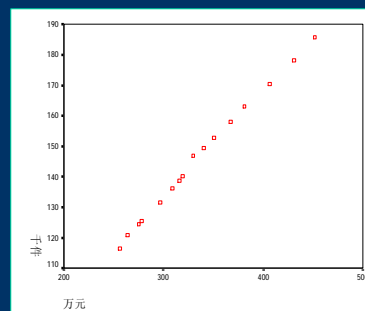
解：令居民收入为自变量 X ，需求量为因变量 Y ，根据表中数据绘制散点图，根据已经给出的数据计算：



序号	Y	X	Y ²	X ²	XY
1	116.5	255.7	13572.25	65382.49	29789.05
2	120.8	263.3	14592.64	69326.89	31806.64
3	124.4	275.4	15475.36	75845.16	34259.76
4	125.5	278.3	15750.25	77450.89	34926.65
5	131.7	296.7	17344.89	88030.89	39075.39
6	136.2	309.3	18550.44	95666.49	42126.66
7	138.7	315.8	19237.69	99729.64	43801.46
8	140.2	318.8	19656.04	101633.44	44695.76
9	146.8	330	21550.24	108900	48444
10	149.6	340.2	22380.16	115736.04	50893.92
11	153	350.7	23409	122990.49	53657.1
12	158.2	367.3	25027.24	134909.29	58106.86
13	163.2	381.3	26634.24	145389.69	62228.16
14	170.5	406.5	29070.25	165242.25	69308.25
15	178.2	430.8	31755.24	185588.64	76768.56
16	185.9	451.5	34558.81	203852.25	83933.85
合计	2339.4	5371.6	348564.74	1855674.54	803822.07



根据数据得到相关图：



解：令居民收入为自变量 X ，需求量为因变量 y ，根据表中数据绘制散点图，根据已经给出的数据采用最小二乘法计算回归系数：

$$\hat{\beta}_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} = 0.3524$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 27.9123$$

得一元线性回归方程：

$$\hat{y} = 27.9123 + 0.3524 X$$

并计算出各离差平方和为：

$$SST = \sum (y_i - \bar{y})^2 = 5191.35 \quad SSE = \sum (y_i - \hat{y}_i)^2 = 17.57$$

$$SSR = \sum (\hat{y}_i - \bar{y})^2 = \hat{\beta}_1^2 L_{xx} = 5173.78$$



回归系数检验

建立假设： $H_0: \beta_1=0; H_1: \beta_1 \neq 0$

$$\text{计算检验统计量 } t_{\beta_1} = \frac{\hat{\beta}_1}{S_{\hat{\beta}_1}}$$

$$S^2 = s^2 = \frac{\sum \hat{e}_i^2}{n-2} = \frac{17.57}{14} = 1.225$$

$$S_{\hat{\beta}_1} = \sqrt{\frac{S^2}{\sum (x_i - \bar{x})^2}} = \sqrt{\frac{1.225}{40495.87}} = 0.0055$$

$$t_{\beta_1} = \frac{\hat{\beta}_1}{S_{\hat{\beta}_1}} = \frac{0.3524}{0.0055} = 64.2069$$

根据显著性水平 $\alpha=0.05$ ，查 $t_{0.025}(14)=2.1448$ 。显然有 $t_{\beta_1} > t_{0.025}(14)=2.1448$ ，表明回归系数是显著不等于0。



Ø 回归方程的显著性检验

建立假设: $H_0: \beta_1=0$; $H_1: \beta_1 \neq 0$

$$F = \frac{\hat{\alpha} (\hat{y}_i - \bar{y})^2}{\hat{\alpha} (y_i - \hat{y}_i)^2 / n - 2} = \frac{SSR / 1}{SSE / n - 2}$$

$$Q SSR = 5173.78$$

$$SSE = 17.57$$

$$F = \frac{5173.78 / 1}{17.57 / 14} = 4122.53$$

根据给定的显著性水平 $\alpha=0.05$, 两个自由度, $df_1=1, df_2=14$. 查临界值 $F_{0.05}(1, 14)=4.60$ 显然有

$$F=4122.53 > F_{0.05}(1, 14)=4.60$$

F检验通过, 可以认为回归方程的回归效果是显著的;



Ø 拟合程度测定

$$\text{计算样本可决系数: } r^2 = \frac{SSR}{SST} = \frac{5173.78}{5191.35} = 0.9966$$

r^2 接近于1, 表明回归直线与样本点的拟合程度很高。

六、回归系数的置信区间

根据区间估计的计算方法, 给定置信水平 $1-\alpha$, 则可以求出回归系数 b_0, b_1 的置信区间分别为:

$$(\hat{b}_0 - \hat{s}_{b_0} t_{\alpha/2}(n-2), \hat{b}_0 + \hat{s}_{b_0} t_{\alpha/2}(n-2))$$

$$(\hat{b}_1 - \hat{s}_{b_1} t_{\alpha/2}(n-2), \hat{b}_1 + \hat{s}_{b_1} t_{\alpha/2}(n-2))$$

由此得出例4.1中回归系数 b_1 的95%的置信区间为:

$$(0.3524 - 2.1448 \cdot 0.0055, 0.3524 + 2.1448 \cdot 0.0055)$$

$$= (0.3406, 0.3642)$$



§4.5 残差分析

一、残差的概念与残差图

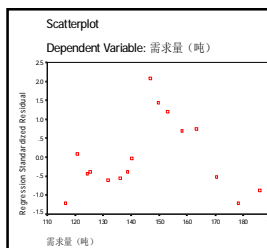
残差指实际观察值与回归值之差 e_i , 是随机项 ε_i 的估计。

$$e_i = y_i - \hat{y}_i = y_i - \hat{b}_0 - \hat{b}_1 x_i$$

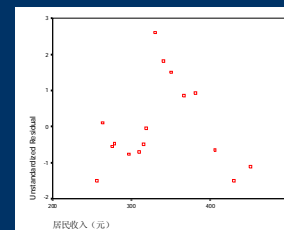
$$e_i = y_i - b_0 - b_1 x_i$$

以自变量X或因变量Y为横轴, 以残差为纵轴可以得到残差图。

如果一个回归模型满足给定的基本假设条件, 则残差应当在 $e=0$ 附近的带状区域内随机排列。否则, 则表明模型存在一定的问題。书上给出了几种典型的残差图的实例。



例题的残差图



二、残差的性质

性质1 $E(e_i) = 0$

性质2 $Var(e_i) = \frac{\hat{\sigma}^2}{n} \cdot \frac{1}{n} - \frac{(x_i - \bar{x})^2}{L_{xx}} \cdot \frac{\hat{\sigma}^2}{n} = (1 - h_{ii}) \hat{\sigma}^2$

式中 $h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{L_{xx}}$ 称为杠杆值 (Leverage Value),

有 $0 < h_{ii} < 1$. 说明 x_i 远离平均值时杠杆值比较大, 则方差较小。

性质3 残差满足约束条件: $\sum \hat{a} e_i = 0; \sum \hat{a} e_i x_i = 0$

即残差项不是独立的。



三、标准化残差与学生化残差

残差分析中的异常值是指超过 $\pm 3s$ 的残差。由于普通的残差的方差不等, 故引入标准化残差和学生化残差。

Ø 标准化残差

$$ZRE_i = \frac{e_i}{\hat{s}}$$

Ø 学生化残差

$$SRE_i = \frac{e_i}{\hat{s} \sqrt{1 - h_{ii}}}$$

标准化残差使残差具有可比性, 一般 $|ZRE| > 3$ 的相应观测值判为异常值。但没有解决方差不等的问题。而学生化残差则进一步解决了标准化残差没有解决的方差不等的问题, 在寻找异常值 $|SRE| > 3$ 方面有优越性。

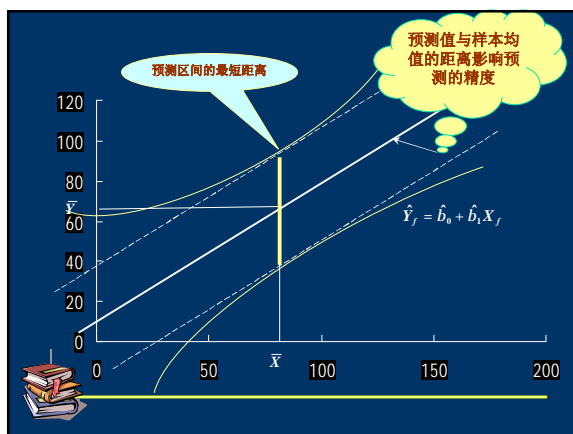


一、预测值及预测区间

点预测

$$\hat{y}_0 = \hat{b}_0 + \hat{b}_1 X_0$$

区间预测

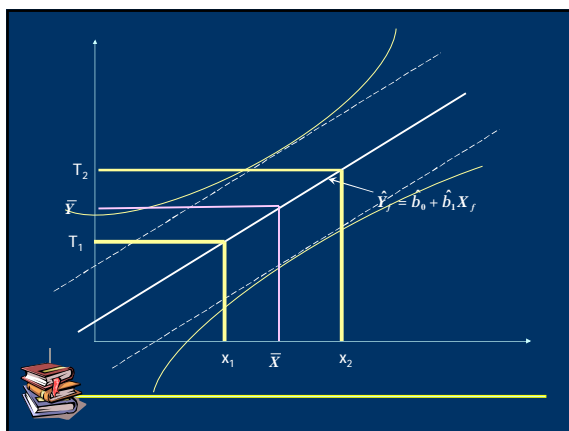
$$\hat{y}_0 - y_0 \sim N[0, s^2(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2})]$$
$$[\hat{y}_0 \pm t_{\alpha/2}(n-2) \times \hat{S} \sqrt{1 + h_{ii}}]$$
$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_j - \bar{x})^2}$$
$$\hat{y}_0 \pm Z_{\alpha/2} \hat{S}$$

$$\hat{y}_0 = 9.99 + 0.1802 \cdot 200 = 46.03$$
$$\begin{aligned} \hat{y}_0 \pm t_{\alpha/2}(n-2) \times \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\hat{\sigma}^2(x_i - \bar{x})^2}} \\ = 46.03 \pm t_{0.025}(13) \cdot 3.334 \cdot \sqrt{1 + \frac{1}{15} + \frac{200 - 101.07}{20224.74}} \\ = (42.069, 49.988) \end{aligned}$$

上式的实际意义：当收入达到200单位时，以95%的把握程度认为：食品支出大约在42.069至49.988之间。

二、控制问题

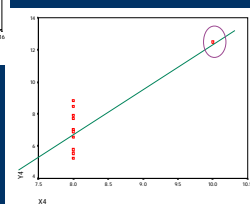
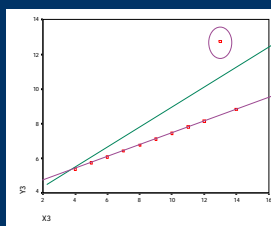
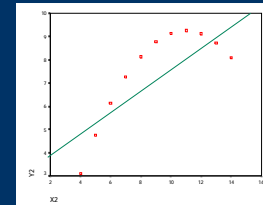
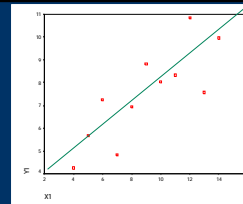
$$P\{T_1 < y < T_2\} = 1 - a$$
$$\begin{aligned} \hat{\mathbf{y}}(x) - 2\hat{\mathbf{S}} &= \hat{\mathbf{y}}(x) - 2\mathbf{S} > T_1 \\ \hat{\mathbf{y}}(x) - 2\hat{\mathbf{S}} &= \hat{\mathbf{y}}(x) + 2\mathbf{S} < T_2 \end{aligned}$$

2近似等于标准正态分布95%的置信概率。



回归模型的进一步讨论：见例题 P 101.

第一组		第二组		第三组		第四组	
X	Y	X	Y	X	Y	X	Y
4	4.26	4	3.1	4	5.39	8	6.58
5	5.68	5	4.74	5	5.73	8	5.76
6	7.24	6	6.13	6	6.08	8	7.71
7	4.82	7	7.26	7	6.44	8	8.84
8	6.95	8	8.14	8	6.77	8	8.47
9	8.81	9	8.77	9	7.11	8	7.04
10	8.04	10	9.14	10	7.46	8	5.25
11	8.33	11	9.26	11	7.81	8	5.56
12	10.84	12	9.13	12	8.15	8	7.91
13	7.58	13	8.74	13	12.74	8	6.89
14	9.96	14	8.1	14	8.84	10	12.5



Model Summary ^b					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Durbin-Watson
1	.816 ^a	.667	.629	1.23660	2.788
a. Predictors: (Constant), X1					
Model Summary ^b					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Durbin-Watson
1	.816 ^a	.666	.629	1.23721	.385
a. Predictors: (Constant), X2					
Model Summary ^b					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Durbin-Watson
1	.816 ^a	.666	.629	1.23618	2.638
a. Predictors: (Constant), X3					
Model Summary ^b					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Durbin-Watson
1	.817 ^a	.667	.630	1.23570	1.296
a. Predictors: (Constant), X4					
b. Dependent Variable: Y4					

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	27.510	1	27.510	17.990	.002 ^a
Residual	13.763	9	1.529		
Total	41.273	10			

a. Predictors: (Constant), X1
b. Dependent Variable: Y1

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	27.500	1	27.500	17.966	.002 ^a
Residual	13.776	9	1.531		
Total	41.276	10			

a. Predictors: (Constant), X2
b. Dependent Variable: Y2

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	27.430	1	27.430	17.950	.002 ^a
Residual	13.753	9	1.528		
Total	41.183	10			

a. Predictors: (Constant), X3
b. Dependent Variable: Y3

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	27.490	1	27.490	18.003	.002 ^a
Residual	13.742	9	1.527		
Total	41.232	10			

a. Predictors: (Constant), X4
b. Dependent Variable: Y4

Coefficients ^a					
Model		Unstandardized Coefficients		Standardized Coefficients	Sig.
		B	Std. Error	Beta	
1	(Constant)	3.000	1.125		.026
	X1	.500	.118	.816	.002

a. Dependent Variable: Y1

Coefficients ^a					
Model		Unstandardized Coefficients		Standardized Coefficients	Sig.
		B	Std. Error	Beta	
1	(Constant)	3.001	1.125		.026
	X2	.500	.118	.816	.002

a. Dependent Variable: Y2

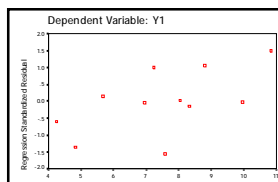


Coefficients ^a					
Model		Unstandardized Coefficients		Standardized Coefficients	Sig.
		B	Std. Error	Beta	
1	(Constant)	3.008	1.124		.025
	X3	.499	.118	.816	.002

a. Dependent Variable: Y3

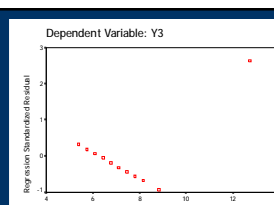
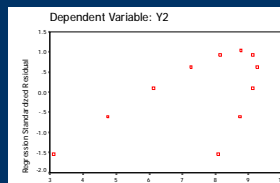
Coefficients ^a					
Model		Unstandardized Coefficients		Standardized Coefficients	Sig.
		B	Std. Error	Beta	
1	(Constant)	-14.995	5.315		.020
	X4	2.750	.648	.817	.002

a. Dependent Variable: Y4



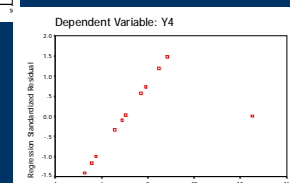
第一组数据的残差图有序列自相关的嫌疑。

第二组数据的残差图有曲线形式



第三组和第四组数据中都明显存在一个异常值。

这四个模型都需要进一步的分析。



一元线性回归应用步骤:

1. 首先根据研究目的确定因变量Y和自变量X;
2. 收集数据后, 首先进行相关分析, 大致确定回归方程(线性或通过变换得到线性模型);
3. 建立一元线性回归模型, 根据样本建立数据表;
4. 用SPSS统计软件对本数据进行分析, 得到相应的输出结果。
5. 对照输出结果对回归模型进行检验: 包括拟合优度的检验(R平方)、回归系数的检验P值、回归方程的检验(方差分析表)、残差分析、寻找异常值等。
6. 通过检验的回归模型称为有效方程, 可以在实际中应用, 也可以用于预测和控制问题。



§4.7 可化为线性的曲线函数

一、几种可化为线性的曲线函数

1. 抛物线函数 (Quadratic 二次函数) $y = a + bX + cX^2$ 经变换得:

$$y = a + bX_1 + cX_2$$

其中: $X_1 = X$, $X_2 = X^2$, 故方程变换后成为一个线性方程。

2. 双曲线函数 (Inverse逆函数) $y = a + b/X$ 做变换得:

$$y = a + bX_1 \quad X_1 = 1/X$$

3. 幂函数 (Power) $y = aX^b$

$$\ln y = \ln a + b \ln X$$

即: $y^c = a^c + bX^c$ 其中 $y^c = \ln y$, $a^c = \ln a$, $X^c = \ln X$

参数b是y对于自变量X的弹性, 是X变动1%时, y变动的百分比。



4. 指数函数 (Exponential)

$$y = ae^{bx} \quad \ln y = \ln a + bx$$

5. 逆函数 (Inverse)

$$y = b_0 + (b_1/x) \quad y = b_0 + b_1 x^c \quad x^c = 1/x$$

6. 对数函数 (Logarithmic)

$$y = a + b \log x$$

7. 逻辑函数 (Logistic)

$$y = \frac{L}{1 + ae^{-bx}} \quad (L, a, b > 0)$$

$$ae^{-bx} = \frac{L}{y} - 1 \quad \ln\left(\frac{L}{y} - 1\right) = \ln a - bx$$

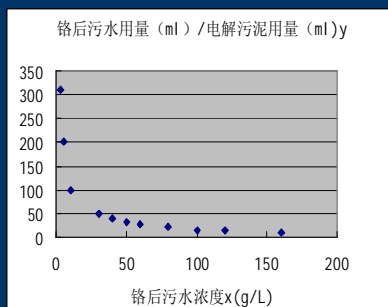


实例分析：某厂表面处理车间试验将铬后污水同电解污泥混合生成无毒溶液。但实际排污的浓度不完全相同，且一定浓度的定量铬后污水只有同定量的电解污泥混合后，才能反应完全。现通过试验，找出铬后污水用量与电解污泥用量之比对于铬后污水浓度之间的关系，试验数据如下：

序号	铬后污水浓度x (g/L)	铬后污水用量 (ml) / 电解污泥用量 (ml) y
1	3	310
2	5	200
3	10	100
4	30	49
5	40	40
6	50	32
7	60	28
8	80	23
9	100	16
10	120	14
11	160	10



解：首先根据数据得到散点图如下：



根据散点图的形状，试用幂函数进行曲线拟合。

假设曲线为幂函数，

$$y = aX^b \quad \text{对数变换} \quad \ln y = \ln a + b \ln X$$

$$\text{取} y^c = \ln y \quad x^c = \ln x \quad a^c = \ln a$$

以变换后计算得：

$$L_{x^c y^c} = 17.462, \quad L_{y^c y^c} = 12.003, \quad L_{x^c y^c} = -14.429$$

由此得回归系数为：

$$\hat{b} = \frac{L_{x^c y^c}}{L_{x^c x^c}} = -0.8263, \quad a^c = 6.6417$$

$$\hat{a} = e^{a^c} = 766.4$$

$$\text{回归方程为: } Y = 766.4x^{-0.8263}$$



对于线性函数应用相关系数表达其线性相关程度，对于非线性函数也可以用相关指数来反映一元非线性回归方程的优劣。

相关指数的定义为：

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{a}_i^2 e_i^2}{L_{yy^c}}$$

由例中的计算得相关指数为：

$$\hat{a}_i e_i^2 = \hat{a} (y_i - \hat{y}_i)^2$$

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{a}_i^2 e_i^2}{L_{yy^c}} = 1 - \frac{0.080}{12.003} = 0.993$$

$$L_{yy} = \hat{a} (y_i - \bar{y})^2$$

由此可见，曲线的拟合程度是很高的。可以用SPSS作多种函数的假设。



将变量X、Y进行对数变换，应用SPSS软件进行线性回归分析得：

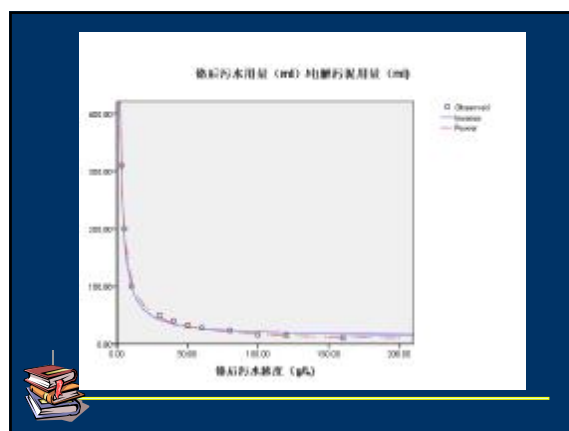
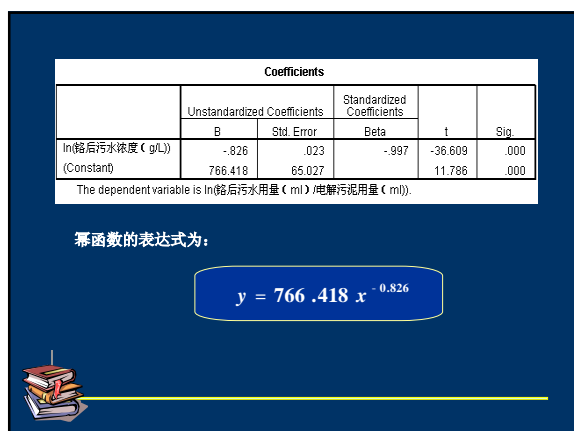
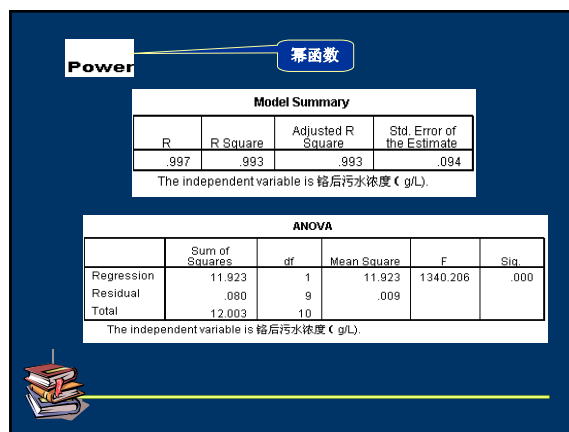
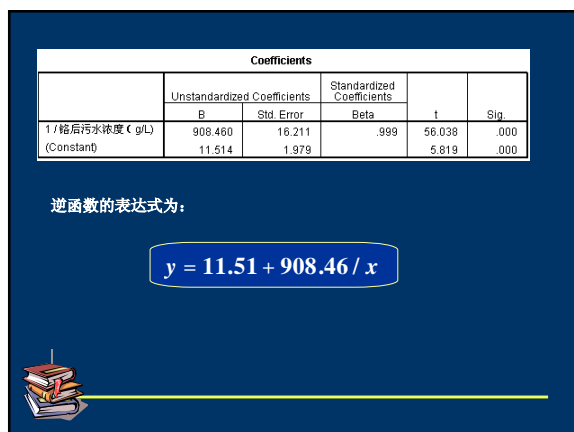
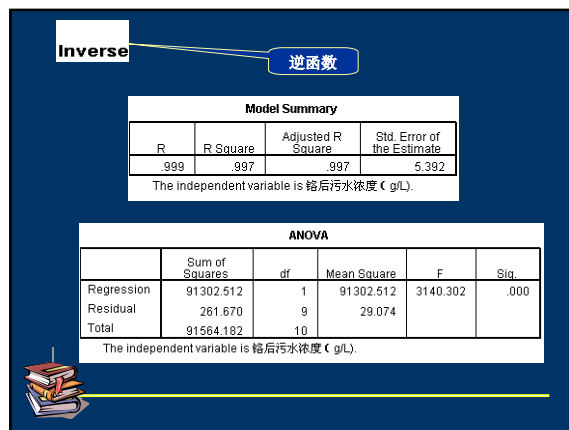
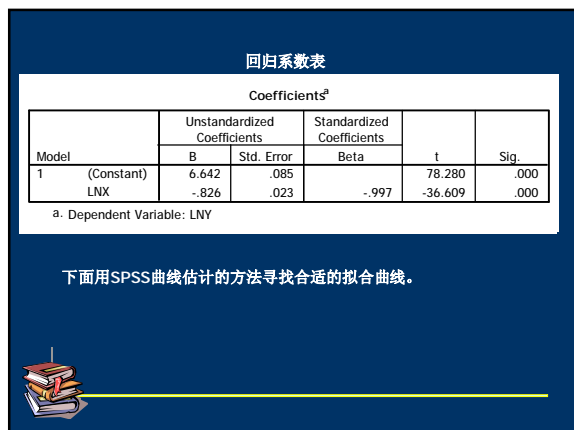
Model Summary ^a					
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Durbin-Watson
1	.997 ^a	.993	.993	.09432	1.239

a. Predictors: (Constant), LNX
b. Dependent Variable: LNY

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	11.923	1	11.923	1340.206	.000 ^b
Residual	.080	9	.009		
Total	12.003	10			

a. Predictors: (Constant), LNX
b. Dependent Variable: LNY





一元线性回归小结:

1. 一元线性回归模型的数学公式(总体与样本), 要求掌握如何利用样本建立回归模型, 注意理解高斯假设的意义。

2. 理解线性回归的意义, 明确回归分析与相关分析的关系。

3. 掌握回归系数的求解方法——最小二乘法。

$$\hat{b}_1 = \frac{L_{xy}}{L_{xx}} = \frac{n \sum X_i Y_i - \sum X_i \sum Y_i}{n \sum X_i^2 - (\sum X_i)^2}, \quad \hat{b}_0 = \frac{\sum Y_i}{n} - \hat{b}_1 \frac{\sum X_i}{n} = \bar{Y} - \hat{b}_1 \bar{X}$$

4. 估计总体方差 σ^2 :

计算中常用 S^2 估计 σ^2 。

$$S^2 = \frac{\sum e_i^2}{n-2}$$



5. \hat{b}_0, \hat{b}_1 的分布

由假设得出Y服从正态分布, 且有: $y \sim N(b_0 + b_1 x, \sigma^2)$

可以证明, $\hat{b}_0 \sim N(b_0, S_{b_0}^2)$; $\hat{b}_1 \sim N(b_1, S_{b_1}^2)$

其中:

$$S_{b_0}^2 = S^2 \left\{ \frac{1}{n} + \frac{\bar{X}^2}{L_{xx}} \right\}, \quad S_{b_1}^2 = \frac{S^2}{L_{xx}}$$

由样本得出回归系数方差的无偏估计为:

$$\hat{S}_{b_0}^2 = S^2 \left\{ \frac{1}{n} + \frac{\bar{X}^2}{L_{xx}} \right\}, \quad \hat{S}_{b_1}^2 = \frac{S^2}{L_{xx}}$$



6. 模型检验

○ **拟合程度的检验:** 利用可决系数 r^2 对回归方程拟合程度的综合测量, 可决系数越大, 模型的拟合程度就越高, 反之, 可决系数越小, 模型的拟合程度就越差:

$$r^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

○ 回归系数的检验 (t 检验)

提出假设: $H_0: b_1=0$; $H_1: b_1 \neq 0$

计算检验统计量值

$$t_{b_1} = \frac{\hat{b}_1}{S_{\hat{b}_1}} = \frac{\hat{b}_1}{S / \sqrt{L_{xx}}}$$

对于给定的显著性水平 α , 查 t 分布表确定临界值(或用 P 值)。要注意当假设是双边检验比较临界值与 t 值后作出判断。



○ 回归方程的检验 (F 检验)

提出假设: $H_0: b_1=0$; $H_1: b_1 \neq 0$

计算离差平方和, 列出方差分析表如下:

离差名称	离差平方和	自由度	均方差	F 值
回归平方和	$SSR = \hat{\hat{a}}(\bar{Y} - \bar{Y})^2$	1	$\frac{SSR}{SSE/(n-2)}$	$F = \frac{SSR}{SSE/(n-2)}$
残差平方和	$SSE = \sum e_i^2$	$n-2$		
总离差平方和	$SST = \sum (Y_i - \bar{Y})^2$	$n-1$		

对于给定的显著性水平 α , 查 F 分布表确定临界值(或用 P 值), 当 $F > F_{\alpha}$ 时, 拒绝原假设。



○ 残差分析

用因变量或者自变量做横轴, 残差(标准化残差, 学生化残差等)为纵轴, 得到残差图, 通过残差图是否随机排列和残差值的范围考察回归模型是否存在异方差性, 是否存在异常值等。

7. 利用回归方程进行预测

○ 预测公式:

由简单回归模型得出基本预测公式为:

$$\hat{y}_0 = \hat{b}_0 + \hat{b}_1 X_0$$

• 点预测: 即给定 X_0 得到 Y_0 的预测值。

注意: 内插与外推预测的区别。



• Y_0 的区间预测:

在总体标准差 σ^2 未知的情况下, 用其无偏估计 S^2 来代替, 可以证明:

$$(y_0 - \hat{y}_0) / S_{e_0} \sim t(n-2)$$

由置信区间的确定方法得 y_0 的置信度为 $(1-\alpha)$ 的置信区间为

$$(\hat{y}_0 \pm t_{\alpha/2}(n-2) \cdot S_{e_0})$$

• 预测的标准误 S_{e_0} : 在标准假定下, 有

$$S_{e_0} = S \sqrt{\frac{1}{n} + \frac{(X_0 - \bar{X})^2}{L_{xx}}}$$

