



中山大學
SUN YAT-SEN UNIVERSITY

实验设计与数据分析

第十讲：线性回归

肖叶，副教授
环境科学与工程学院
中山大学

广州，2023

回归分析

✿ 基本概念

- 回归分析的起源
- 回归分析的研究内容和研究目的
- 变量之间的统计关系
- 回归模型与回归参数
- 总体回归函数和样本回归函数的关系



回归分析的起源

✿ 回归分析 (regression analysis)



A return to a former or less developed state

✿ 起源

- 1886年英国生物学家、统计学家F. Galton（高尔顿，1822-1911）指出：在同一种族里，父亲高的，儿子的平均高度比父亲矮，但高于种族平均高度；父亲矮的，儿子的平均高度比父亲高，但矮于种族平均高度。也就是说，儿子的高度有”回归“于种族平均高度的趋势。

✿ 发展

- 现代统计应用学的重要分支

回归的现代含义要比其原始含义广得多



回归分析的含义

✿ 回归分析的研究内容

- 一个变量或者一组变量（即自变量）之变化对另一个变量（即因变量）之变化的影响程度

✿ 回归分析的研究目的

- 根据已知的自变量的变异来估计或者预测因变量的变异情况
- 例如：探索在已知父亲们身高的条件下，儿子们的平均身高是怎样变动的



变量之间的关系

- ✿ 可以分为两大类
- ✿ 函数关系 (functional relation, deterministic relation)
- ✿ 统计关系 (statistic relation)



变量之间的函数关系

✿ 定义

- 变量之间数量上的确定性关系

✿ 数学表达

- 设 X 为自变量， Y 为因变量， X 和 Y 的函数关系可以表示为： $Y = f(X)$

✿ 特点

- 通过函数式，一个或一组变量在数量上的变化就可以**完全确定**另一个变量在数量上的变化。

✿ 举例

- 圆面积与半径之间的关系： $A = \pi r^2$



变量之间的统计关系

✿ 定义

- 在客观世界里，许多变量之间也存在相随变动，并具有某种规律性；
- 但是，数量关系往往**并非完全确定**

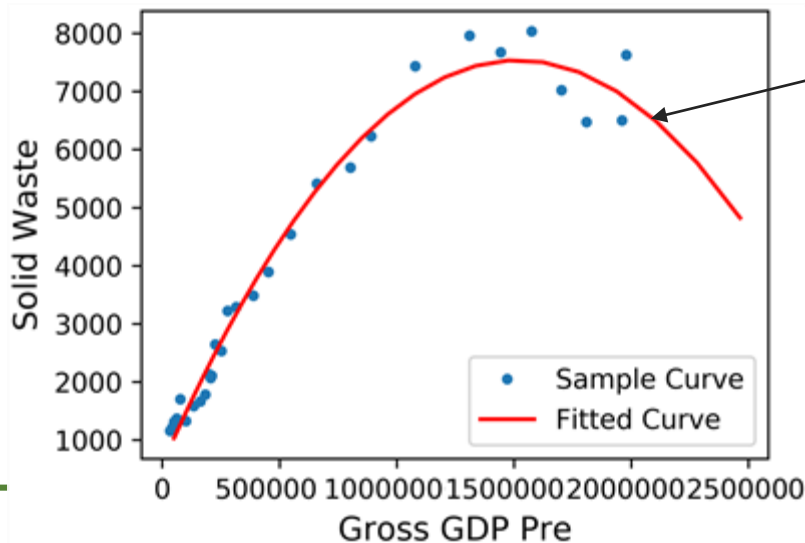
✿ 举例

- 固体废物产生量与经济发展水平（GDP）之间的关系

饮食文化习惯？

地区气候？

产业结构？



库兹涅茨曲线



变量之间的统计关系

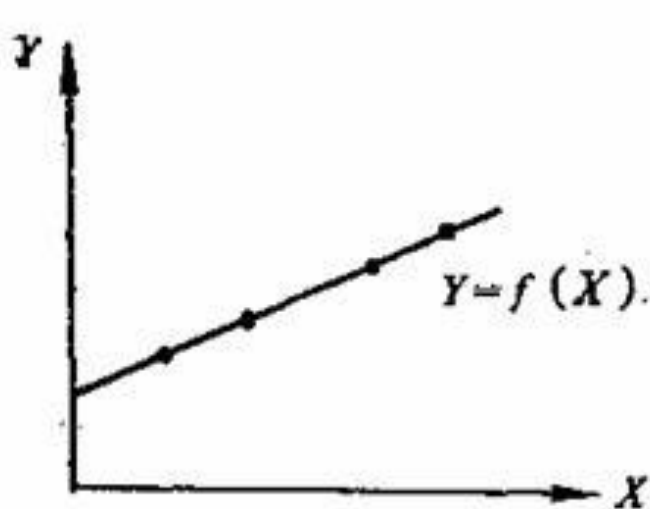
✿ 数学表达

- 如果把固体废物产生作为因变量 Y ,
- 把GDP、食物热量、经纬度和三产比例作为自变量, 并分别用 X_1 、 X_2 、 X_3 和 X_4 来表示,
- 那么它们之间的关系可以表达为:
- $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \varepsilon = f(X) + \varepsilon$
- ε 为随机扰动误差项

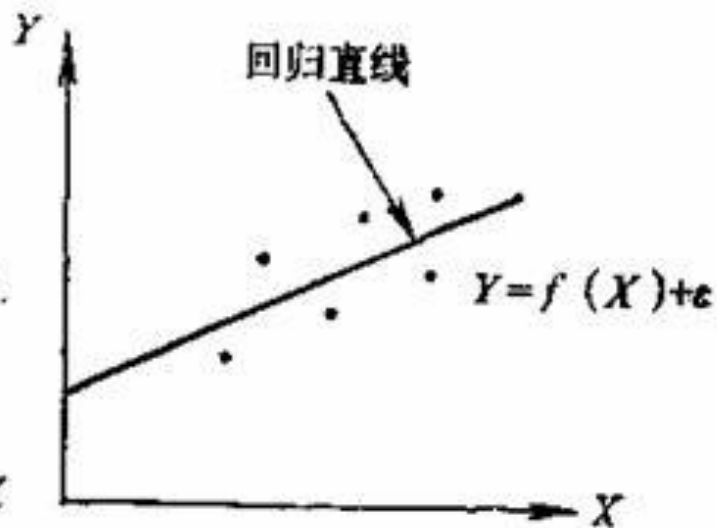


变量之间的两种关系

✿ X与Y变量之间的两种关系



(a) 函数关系



(b) 统计关系



回归模型

✿ 回归分析所研究的变量之间的关系：

- 是一种非确定性的统计关系

✿ 回归模型：

- 回归分析当中用到的对变量之间统计关系进行定量描述的数学模型

✿ 例如：

- $Y = f(X) + \varepsilon$
- $f(X)$ 称为回归函数
- ε 称为随机扰动误差项，或简称误差项



回归模型

✿ 随机扰动误差项 ε ，一般有四个来源：

- 被省略掉而没有纳入方程，但实际上又影响着因变量Y的种种因素；
 - 这些因素被省略有多种原因
 - 在回归分析时尚不知道，难于观测，难于量化，为了简化模型，等等
- 变量的观测误差；
- 模型的设定误差（结构误差，用非线性近似线性）；
- 随机误差

✿ 回归函数：

- 总体回归函数
- 样本回归函数



总体回归函数

✿ 每当 X 取某一给定值时，都会有一个总体与之对应。该总体：

- 由所有 Y 的相应的可能值组成；
- 具有一定的概率分布（条件概率分布）

✿ 总体回归函数（PRF）：

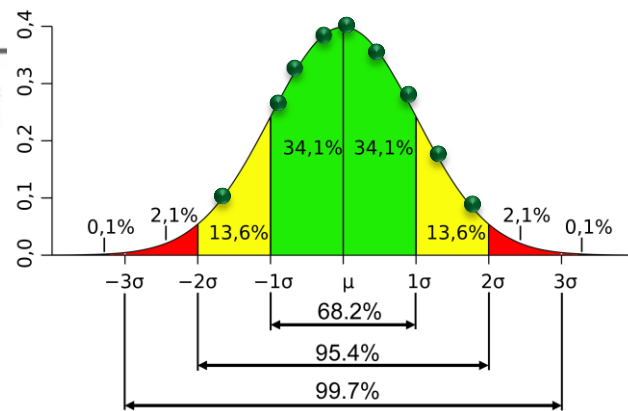
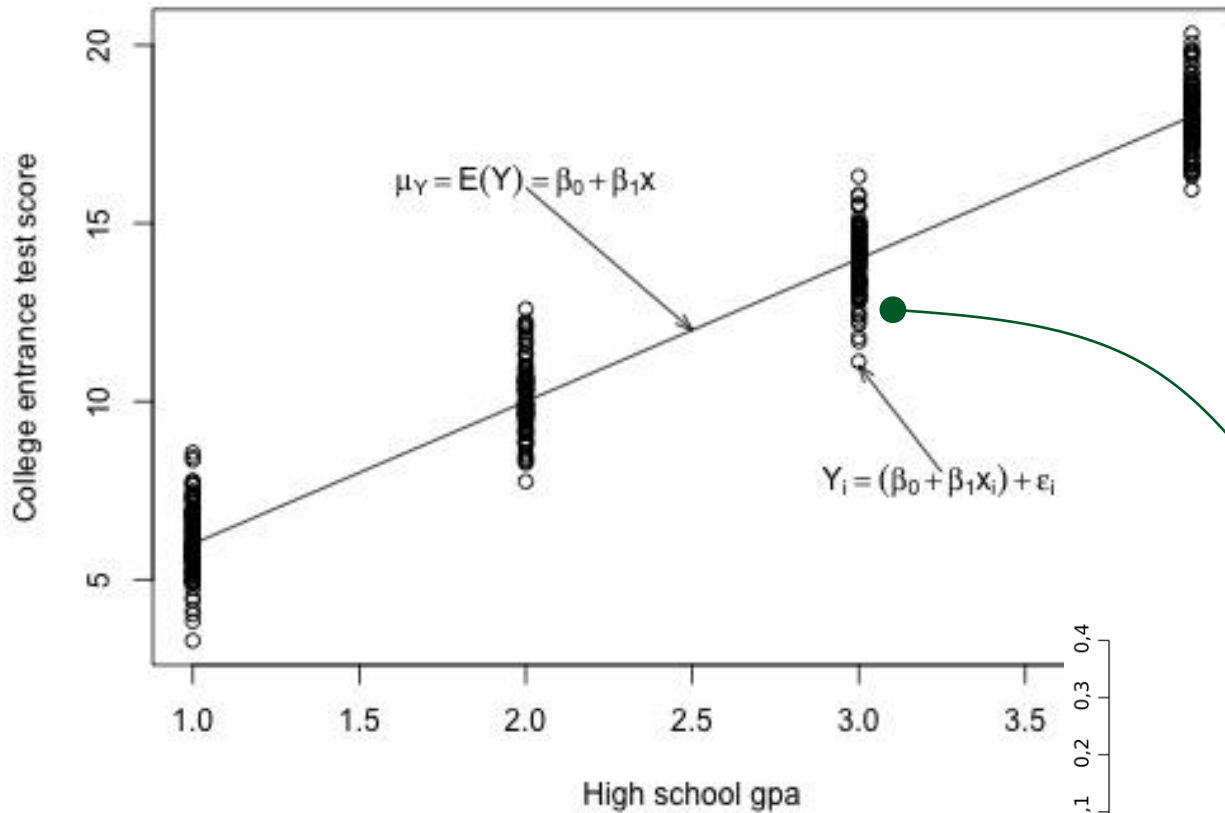
- 从 Y 的各个条件概率分布中求出的各个条件均值（条件数学期望值）所组成的一条 Y 对 X 的回归线（直线或者曲线）

✿ PRF的物理意义：

- PRF代表的是 Y 之条件均值，反映了在 X 取定值的条件下，因变量 Y 的平均变化状态。



总体回归函数



总体回归函数

✿ 理论上，总体回归函数的形式

- 根据定义，总体回归函数 PRF 可以记作：
- $f(X) = E(Y|X)$

✿ 实践中，

- 受各种因素的限制，在回归分析中，人们无法知道总体回归函数的确切形式
- 为了研究X与Y之间的数量变化规律，可以对总体回归函数的形式做必要的、合理的假设



总体回归函数

✿ 总体回归函数的形式：

- 假设 X 与 Y 之间是线性关系，那么，
 - $f(X) = \beta_0 + \beta_1 X_1$ ，该式又被称为一元线性回归方程
 - $Y = \beta_0 + \beta_1 X_1 + \varepsilon$ 被称为一元线性回归模型

✿ “元”：

- 是指回归模型中变量的个数，包括因变量 Y 在内
- 如果回归模型中含有**一个自变量**，则称为一元回归模型
- 如果回归模型中含有**多个自变量**，则称为多元回归模型



样本回归函数

- ✿ 实践中，在 X 取某定值的条件下，不可能掌握 Y 的所有可能值（总体），往往，
 - 只能掌握 Y 的一组样本值
 - 可以直接得到的只能是样本回归函数：



样本回归函数

✿ 回归分析的过程在本质上就是：

- 根据样本回归函数（Sample Regression Function, SRF）来估计未知的总体回归函数（Population Regression Function, PRF）
- 由于通常只知道总体中的一个样本，所以只能获得近似的估计结果。



PRF 与 SRF

✿ 举例

- 假设要研究某类企业的产量和污染物排放量之间的关系，以产量为自变量 X ，排污量为因变量 Y 。
- 一方面：
 - 随着产量的增加，污染物排放量必然也会相应增加
- 另一方面：
 - 除了产量之外，影响排污量的因素还有很多，如企业的不同技术水平、管理水平等等
- 因此，在研究该类企业的排污量与产量之间的数量变化规律时：
 - 在给定产量（ X ）的条件下所得到的排污量（ Y ）值将是围绕某些中心值上下波动的统计值
 - 与给定 X 值相对应的 Y 的所有可能



PRF 与 SRF

✿ 总体已知的情况:

- 全部61个企业的产量、排污量及其10个分组是已知的，并作为统计数据列于表中。根据这些数据就可以获得PRF。

	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	X ₉	X ₁₀
	30	20	60	80	40	50	60	30	70	60
Y ₁	73	79	139	159	102	129	130	84	142	133
Y ₂	91	46	128	173	80	97	135	88	140	136
Y ₃	70	71	127	155	87	96	137	82	146	132
Y ₄	75	50	138	169	98	108	122	96	137	120
Y ₅	85	84	118	146	90	140	126	90	148	143
Y ₆	98		115	170	110			69	154	116
Y ₇			145		119			65	155	
合计	492	330	910	972	686	570	650	574	1022	780
Y的条件概率	1/6	1/5	1/7	1/6	1/7	1/5	1/5	1/7	1/7	1/6
Y的条件均值	82	86	130	162	98	114	130	82	146	130

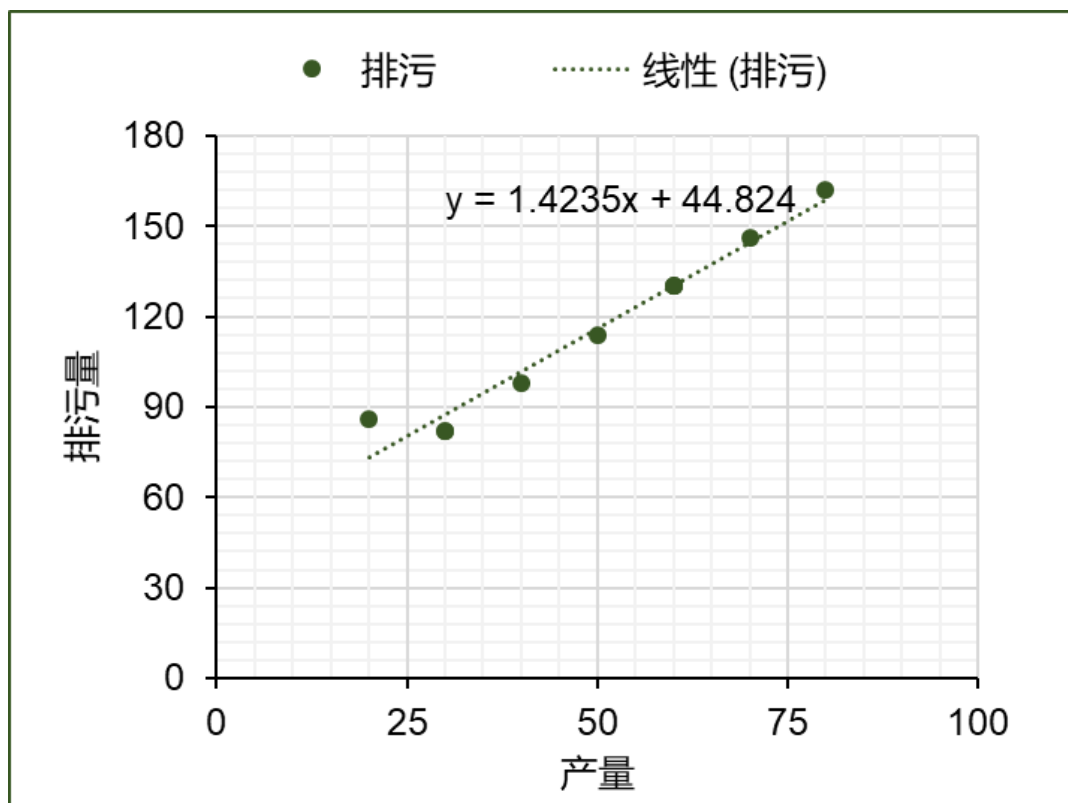


PRF 与 SRF

✿ 总体已知的情况:

- 全部61个企业的产量、排污量及其10个分组是已知的，并作为统计数据列于表中。根据这些数据就可以获得PRF。

	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	X ₉	X ₁₀
Y的条件均值	82	86	130	162	98	114	130	82	146	130



PRF 与 SRF

✿ 总体未知的情况:

- 从Y的总体里随机抽取两个样本，样本1和样本2。
- 用两组样本数据可以得到两条样本回归线： SRF_1 和 SRF_2 。
- 无论是 SRF_1 还是 SRF_2 ，都是对总体回归线PRF的近似。

样本1

编号	X_i	Y_i
1	30	73
2	20	50
3	60	128
4	80	170
5	40	87
6	50	108
7	60	135
8	30	69
9	70	148
10	60	132

样本2

编号	X_i	Y_i
1	30	75
2	20	71
3	60	118
4	80	173
5	40	90
6	50	96
7	60	126
8	30	96
9	70	140
10	60	136



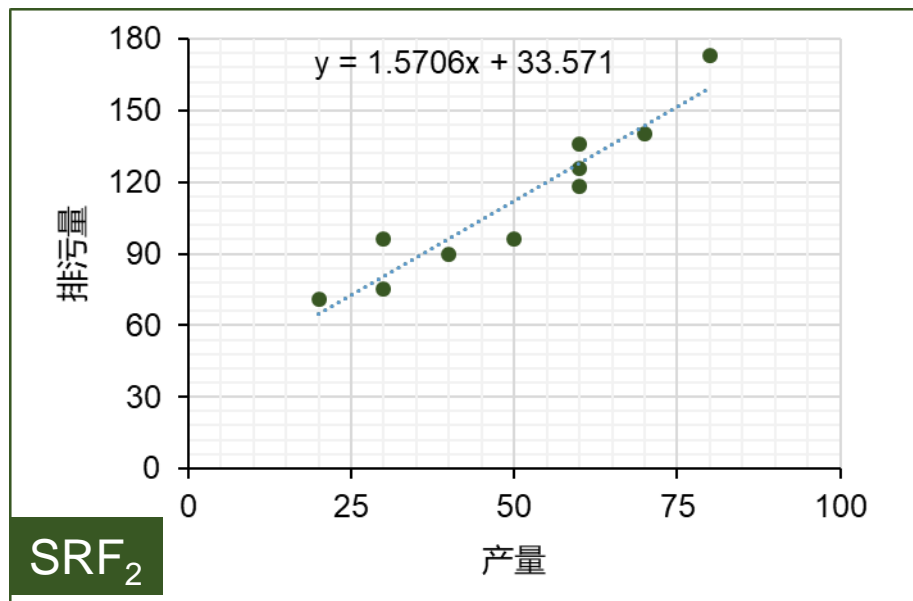
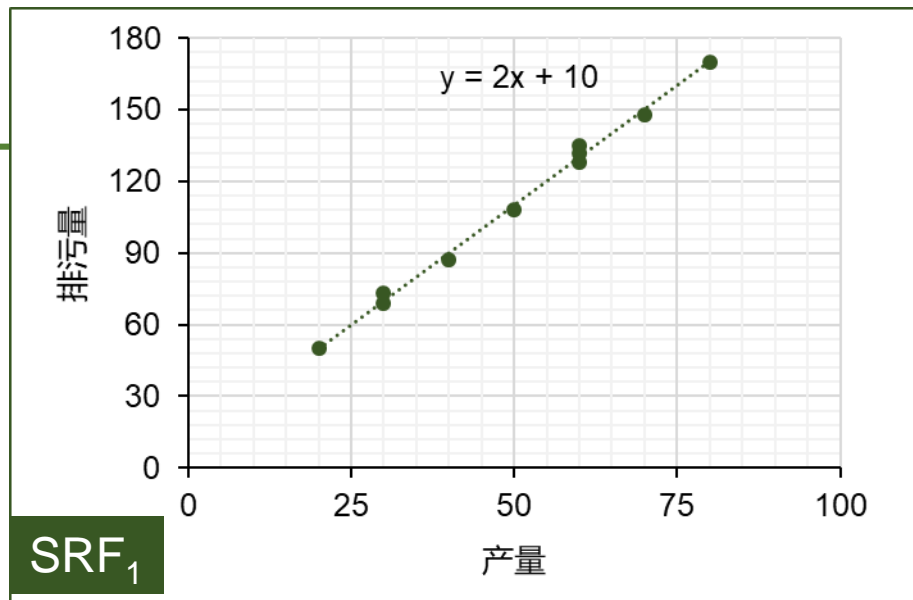
PRF 与 SRF

样本1

X_i	Y_i
30	73
20	50
60	128
80	170
40	87
50	108
60	135
30	69
70	148
60	132

样本2

X_i	Y_i
30	75
20	71
60	118
80	173
40	90
50	96
60	126
30	96
70	140
60	136



PRF 与 SRF

✿ 回归分析的重要任务:

- 根据样本回归函数 SRF 尽可能精确地估计总体回归函数 PRF
- 要达到这一目的, 就要使未知回归参数的估计量尽可能地接近其真值

✿ 如何做到这一点?



回归分析

✿ 基本概念

✿ 一元回归分析

- 回归参数的估计
- 回归模型的检验
- 回归预测

✿ 多元回归分析

✿ 回归函数的形式



回归参数的估计方法

✿ 回归参数的估计方法

- 最小平方法
- 高斯提出的
- 又称为最小二乘法

✿ 以一元线性回归模型为例讨论最小平方法



最小平方法

✿ 一元线性回归模型的一般形式:

- $Y = E(Y|X) + \varepsilon$

- $E(Y|X) = \beta_0 + \beta_1 X$

- Y 为因变量, X 为自变量, β_0 、 β_1 为回归参数, ε 为随机扰动误差项

✿ 对 Y 和 X 分别进行 n 次独立观测, 就可以取得 n 对观测值:

- $(X_i, Y_i), i = 1, 2, \dots, n$, 则有,

- $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad i = 1, 2, \dots, n$

- $E(Y_i|X_i) = \beta_0 + \beta_1 X_i \quad i = 1, 2, \dots, n$



最小平方法

✿ ε_i 实际上是无法观测的，不可能事先确定

- 必须要对它作出若干基本假设，才能做进一步分析

✿ 高斯对此作了以下4条假设：

- $E(\varepsilon_i) = 0$;
 - $Var(\varepsilon_i) = \sigma^2$;
 - $Cov(\varepsilon_i, \varepsilon_j) = 0, i \neq j$;
 - 自变量 X_i 为非随机变量。
-
- 满足以上4条假设的回归模型称为标准回归模型。



最小平方法

✿ 采用最小平方法估计未知回归参数 β_0 、 β_1

■ 估计准则：

■ 求 β_0 、 β_1 的估计量 $\widehat{\beta}_0$ 、 $\widehat{\beta}_1$ ，使得随机扰动误差项 ε_i 的平方和 S 最小。

$$S = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n [Y_i - (\beta_0 + \beta_1 X_i)]^2$$

■ 估计方法：

■ 对 S 分别求 β_0 、 β_1 的偏导数，并使之为零。



最小平方法

✿ 此时获得的估计量称为,

- 未知回归参数的最小平方估计量 (LSE)
- β_0 、 β_1 的LSE:

$$\widehat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{S_{XY}}{S_{XX}} = \frac{\sum_{i=1}^n X_i Y_i - \frac{1}{n} [(\sum_{i=1}^n X_i)(\sum_{i=1}^n Y_i)]}{\sum_{i=1}^n X_i^2 - \frac{1}{n} (\sum_{i=1}^n X_i)^2}$$

$$\widehat{\beta}_0 = \bar{Y} - \widehat{\beta}_1 \bar{X}$$



最小平方法

✿ 利用样本1的数据说明计算 $\hat{\beta}_1$ 和 $\hat{\beta}_0$ 的方法:

$$\hat{\beta}_1 = \frac{61800 - (500)(1100)/10}{28400 - (500)^2/10} = 2$$

$$\hat{\beta}_0 = \frac{1}{10}(1100 - 2 \times 500) = 10$$

✿ 用10个样本企业数据拟合所得的产量—排污量回归直线方程:

$$\hat{Y}_i = 10 + 2X_i$$

企业	产量 X_i	排污量 Y_i	$X_i Y_i$	X_i^2
1	30	73	2190	900
2	20	50	1000	400
3	60	128	7680	3600
4	80	170	13600	6400
5	40	87	3480	1600
6	50	108	5400	2500
7	60	135	8100	3600
8	30	69	2070	900
9	70	148	10360	4900
10	60	132	7920	3600
合计	500	1100	61800	28400



回归分析

✿ 基本概念

✿ 一元回归分析

- 回归参数的估计
- 回归模型的检验
- 回归预测

✿ 多元回归分析

✿ 回归函数的形式



对回归模型的检验

- ✿ 在采用回归模型时，把自变量 X 与因变量 Y 之间的关系假设为：
 - 线性关系
- ✿ 这种假设是否适当？
 - 还需通过统计检验



对回归模型的检验

✿ 回归分析中的假设检验包含两个内容：

- 检验变量之间的总体线性关系是否显著，检验自变量与因变量之间的关系能否用一个适当的回归模型表示；
- 检验回归参数，检验回归模型中的每一个自变量对因变量的影响程度是否显著；

✿ 两种显著性统计检验在次序上不能颠倒：

- 只有当回归模型所代表的变量之间的线性关系通过检验后，进一步检验模型中的个别回归参数才有意义。



总体线性关系的检验

- ✿ 以一元线性回归模型为例，讨论变量之间总体线性关系的显著性检验。
- ✿ 两种方法：
 - 方差分析
 - 考察残差图



方差分析

✿ 方差分析:

- 对回归分析的初步计算结果进行总结的常用方法
- Analysis of Variance (ANOVA)
- 方差分析表的形式

Source of Variance 变差来源 SV	Degree of Freedom 自由度 df	Sum of Squares 平方和 SS	Mean Square 均方 MS	统计量 F
回归 Regression	1	$SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$	MSR= SSR/1	F=MSR/MSE
残差 Residual	n-2	$SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$	MSE= SSE/(n-2)	
总离差 Total	n-1	$SST = SSR + SSE = \sum_{i=1}^n (Y_i - \bar{Y})^2$	—	—



方差分析

✿ F 统计量的特征:

- F 值的高低直接反映因变量 Y 与自变量 X 之间线性关系的强弱;
- F 统计量服从自由度为 $(1, n-2)$ 的 F 分布。

✿ 建立 F 检验的决策规则:

- 设 $H_0: \beta_1 = 0$; $H_1: \beta_1 \neq 0$
- 若 $F \leq F_\alpha(1, n-2)$, 则接受 $\beta_1 = 0$, 即线性关系不显著
- 若 $F > F_\alpha(1, n-2)$, 则接受 $\beta_1 \neq 0$, 即线性关系显著
- $F_\alpha(1, n-2)$ 是 F 统计量的临界值, α 是事先确定的显著水平, 而 $(1-\alpha)$ 表示检验的可靠程度。例如取 $\alpha=0.01$, 由 F 分布表可以查到 $F_{0.01}(1, n-2)$ 的值, 而检验的可靠程度为99%。



方差分析

✿ 仍以样本1当中的产量、排污量数据为例。

企业	产量 X_i	排污量 Y_i			
1	30	73	70	1600	9
2	20	50	50	3600	0
3	60	128	130	400	4
4	80	170	170	3600	0
5	40	87	90	400	9
6	50	108	110	0	4
7	60	135	130	400	25
8	30	69	70	1600	1
9	70	148	150	1600	4
10	60	132	130	400	4
合计	500	1100	1100	13600	60

方差分析表

变差来源 SV	自由度 df	平方和 SS	均方 MS	统计量 F
回归	1	13600	13600	1813.3
残差	8	60	7.5	
总离差	9	13640	—	—

✿ 比较1813.3与 $F_{\alpha}(1, 8)$ 的大小就可以确定总体线性关系是否显著。

$$F_{0.1}(1, 8) = 3.46$$



考察残差图

✿ 残差

- $e_i = (Y_i - \hat{Y}_i) \sim N(0, \sigma^2)$

✿ 对 σ^2 的估计

- 最小平方估计量：对总离差进行分解

$$LSE(\hat{\sigma}^2) = \frac{SSE}{n-2} = \frac{1}{n-2} \sum_{i=1}^n e_i^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = MSE$$



考察残差图

✿ 残差的标准化:

- $e_i^* = (e_i/\sigma) \sim N(0, 1)$
- $P(|e_i^*| < 2) = 0.9545$, 即 e_i^* 的值落在区间 $(-2, +2)$ 上的概率应为95.45%

✿ 绘制标准化残差图:

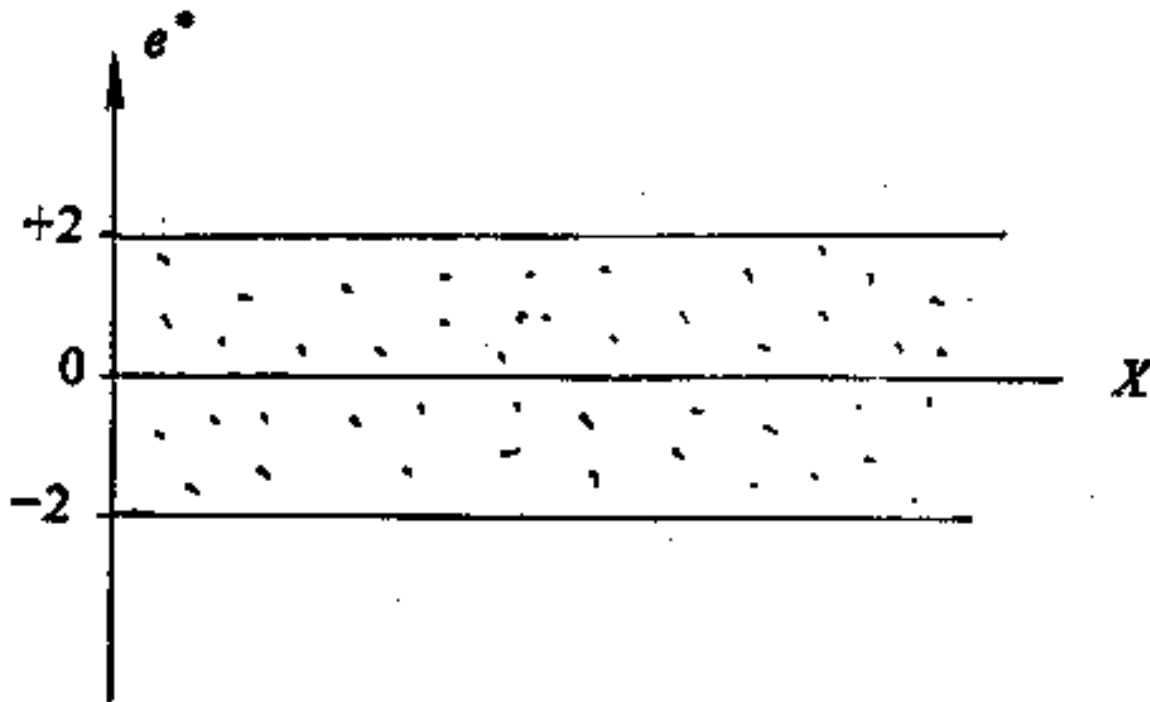
- X_i 为横坐标, e_i^* 为纵坐标
- 将数据点 (X_i, e_i^*) , $i = 1, \dots, n$, 画在平面图上



考察残差图

✿ 考察标准化残差图时，常遇到的4种情况：

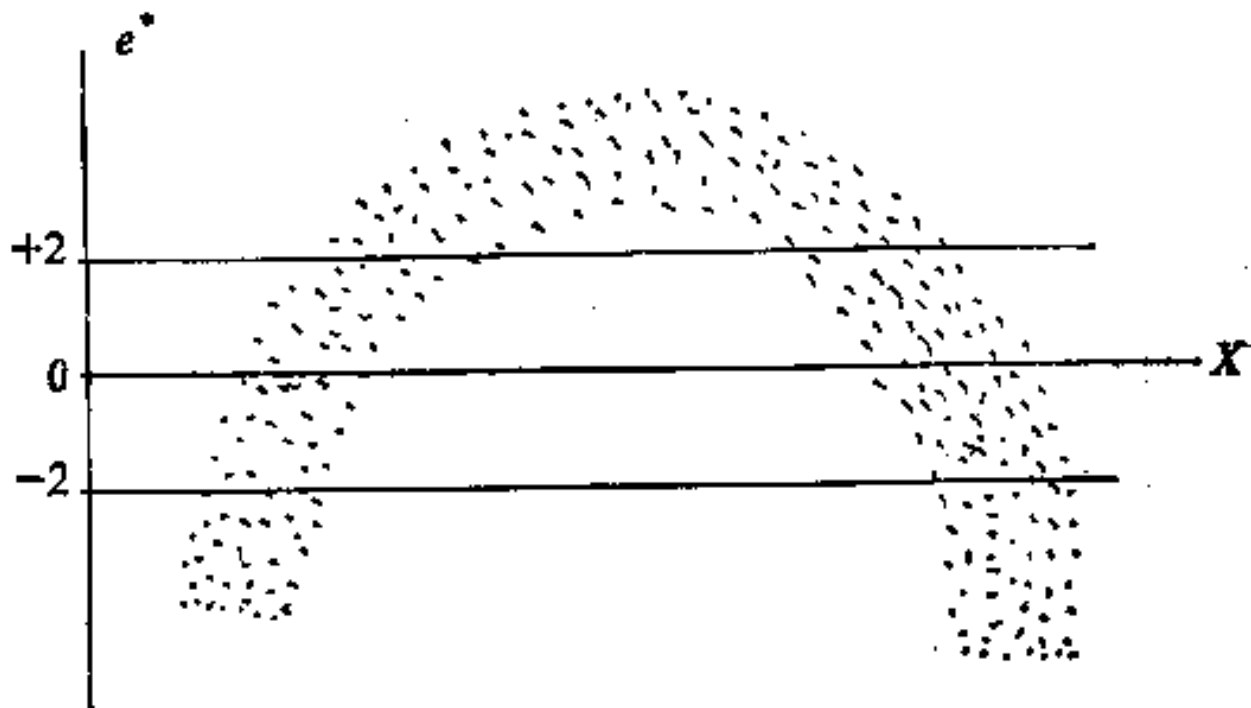
- (1) 回归方程拟合良好，绝大多数数据点都在 $(-2,+2)$ 水平带状区间内，且不帶有任何系统趋势，完全随机的散布。



考察残差图

✿ 考察标准化残差图时，常遇到的4种情况：

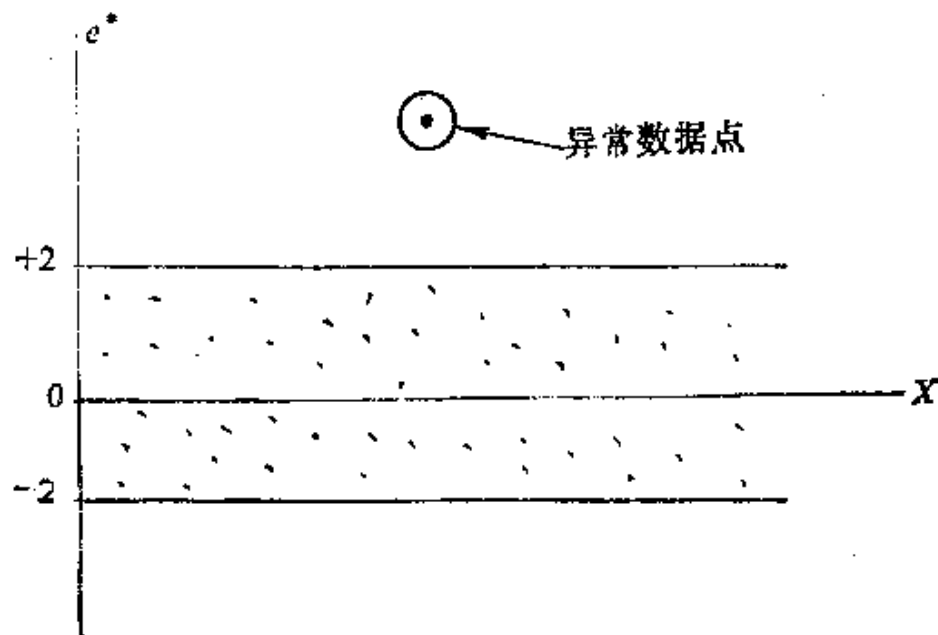
- (2) 回归函数具有曲线形式，如果回归函数的形式应为曲线，但却采用了直线回归方程，标准残差图就会出现类似下图的形式。



考察残差图

✿ 考察标准化残差图时，常遇到的4种情况：

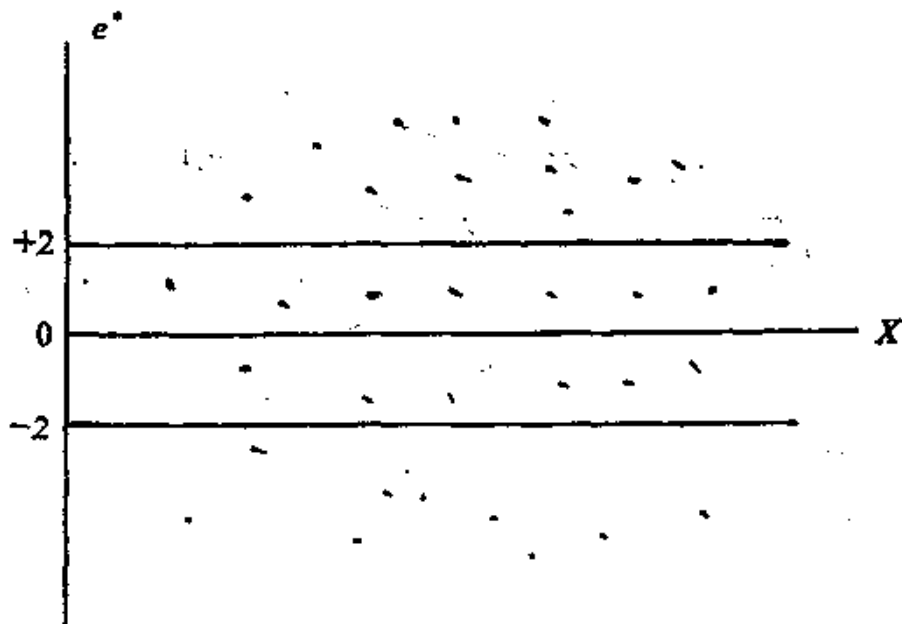
- (3) 样本数据中有一个或数个异常点 (outlier)，如果由于错误或其它特殊因素的影响，样本数据中出现了个别异常点，那么在标准残差图里就会相应出现远离大多数数据点的一个或数个异常数据点。



考察残差图

✿ 考察标准化残差图时，常遇到的4种情况：

- (4) 回归方程拟合不充分：如果有许多数据点落在区间 $(-2,+2)$ 的外面，就说明回归方程对数据的拟合是不充分的；原因可能是回归模型的函数形式选择不当，也可能是漏掉了一个（或数个）比较重要的自变量。



回归参数的检验

✿ 以一元回归模型为例，参数检验的任务

- 对回归参数 β_0 、 β_1 进行统计检验，并求出其置信区间。

✿ 对 β_1 的检验：

- β_1 代表了 X_i 的单位变动对 Y_i 的影响程度，所以对 β_1 的检验首先考察这种影响程度是否与零有显著差异
- 在一元回归模型中，对参数 β_1 的检验相当于检验总体线性关系



回归参数的检验

✿ 对 β_1 的检验过程:

■ 设 $H_0: \beta_1=0; H_1: \beta_1 \neq 0;$

■ 检验统计量 t :

$$t = \frac{\widehat{\beta}_1}{\sqrt{\widehat{\sigma}^2/S_{XX}}} \sim t_{n-2}$$

$$S_{XX} = \sum_{i=1}^n (X_i - \bar{X})^2$$
$$\widehat{\sigma}^2 = \frac{SSE}{n-2} = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \widehat{Y}_i)^2$$

• 服从自由度为 $(n-2)$ 的 t 分布

■ 决策规则:

- 若 $|t| \leq t_{\alpha/2, n-2}$, 则接受 $\beta_1 = 0$, 即X对Y的影响是不显著的;
- 若 $|t| > t_{\alpha/2, n-2}$, 则接受 $\beta_1 \neq 0$, 即X对Y的影响是显著的



回归参数的检验

✿ 对 β_0 的检验过程

- 设 $H_0: \beta_0 = \beta_0^*$; $H_1: \beta_0 \neq \beta_0^*$
 - β_0^* 是某个假设的特定值, 如果 $\beta_0^* = 0$, 相当于检验回归直线是否过原点
- 检验统计量 t :

$$t = \frac{\widehat{\beta}_0 - \beta_0^*}{\sqrt{\widehat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{S_{XX}} \right)}} \sim t_{n-2}$$

- 服从自由度为 $(n-2)$ 的 t 分布
- 决策规则:
 - 若 $|t| \leq t_{\alpha/2, n-2}$, 则接受 $\beta_0 = \beta_0^*$
 - 若 $|t| > t_{\alpha/2, n-2}$, 则接受 $\beta_0 \neq \beta_0^*$



回归参数的检验

✿ 点估计量

- 采用最小平方法获得的参数估计量是点估计量。
- 由于样本与样本之间数据差异产生的波动，单一的参数点估计量很可能与参数的真值不同。

✿ 如何对点估计量的可靠性进行度量：

- 通过点估计量围绕其真值变动的标准误差或者方差来度量
- 根据点估计量的已知概率分布来确定它落在某一区间（置信区间）内的概率
- 这一概率用 $(1 - \alpha)$ 来表示，又称置信水平



回归参数的检验

✿ 参数 β_1 的置信区间

- 把参数 β_1 的标准误差估计量定义为：

$$s(\widehat{\beta}_1) = \sqrt{\widehat{\sigma}^2 / S_{XX}}$$

- 参数 β_1 的 $100(1 - \alpha)\%$ 置信区间为：

$$\widehat{\beta}_1 \pm t_{\alpha/2, n-2} \cdot s(\widehat{\beta}_1)$$



回归参数的检验

✿ 参数 β_0 的置信区间

- 把参数 β_0 的标准误差估计量定义为：

$$S(\widehat{\beta}_0) = \sqrt{\widehat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{S_{XX}} \right)}$$

- 参数 β_0 的 $100(1 - \alpha)\%$ 置信区间为：

$$\widehat{\beta}_0 \pm t_{\alpha/2, n-2} \cdot S(\widehat{\beta}_0)$$



回归预测问题

✿ 在一元回归模型中，预测是指：

- 根据自变量 X 的某一已知值 X_0 求因变量 Y 的相应值 Y_0 的过程
- X_0 可以是样本数据中的某个数据；
- X_0 也可以是样本数据值域范围内的某一个取值；
- 可以预测 Y_0 的平均值， $E(Y_0|X_0) = \beta_0 + \beta_1 X_0$ ，简称为均值预测；
- 可以预测与 X_0 相对应的总体个别值（或称为特定值） Y_0 ， $Y_0 = \beta_0 + \beta_1 X_0 + \varepsilon_0$ ，简称为个别值预测



均值预测

✱ $E(Y_0|X_0)$ 的点估计量:

$$\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 X_0$$

✱ $E(Y_0|X_0)$ 的点估计量与真值之差是一种预测误差

✱ 把 Y_0 的估计标准误差定义为:

$$S(\hat{Y}_0) = \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{(X_0 - \bar{X})^2}{S_{XX}} \right]}$$

✱ $E(Y_0|X_0)$ 的 $100(1 - \alpha)\%$ 置信区间为:

$$\hat{Y}_0 - t_{\frac{\alpha}{2}, n-2} \cdot S(\hat{Y}_0) \leq Y_0 \leq \hat{Y}_0 + t_{\frac{\alpha}{2}, n-2} \cdot S(\hat{Y}_0)$$



个别值预测

✿ Y 变量的总体个别值 Y_0^* 与 Y_0 的平均值 $E(Y_0|X_0)$ 有:

■ 相同的预测值:

$$\widehat{Y}_0^* = \widehat{\beta}_0 + \widehat{\beta}_1 X_0$$

■ 不同的预测误差:

- Y_0^* 的预测区间比 $E(Y_0|X_0)$ 的置信区间要宽

✿ 把 Y_0^* 的估计标准误差定义为

$$S(\widehat{Y}_0^*) = \sqrt{\widehat{\sigma}^2 \left[1 + \frac{1}{n} + \frac{(X_0 - \bar{X})^2}{S_{XX}} \right]}$$

✿ Y_0^* 的 $100(1 - \alpha)\%$ 预测区间为:

$$\widehat{Y}_0^* - t_{\frac{\alpha}{2}, n-2} \cdot S(\widehat{Y}_0^*) \leq Y_0^* \leq \widehat{Y}_0^* + t_{\frac{\alpha}{2}, n-2} \cdot S(\widehat{Y}_0^*)$$



回归分析

✿ 基本概念

✿ 一元回归分析

- 回归参数的估计
- 回归模型的检验
- 回归预测

✿ 多元回归分析

✿ 回归函数的形式



多元回归分析

✿ 多元回归分析：

■ 二元回归分析的延伸

✿ 多元线性回归模型：

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots \cdots + \beta_p X_p + \varepsilon$$

✿ 多元线性回归方程：

$$E(Y|X = X_1, X = X_2, \dots \dots, X = X_p) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots \cdots + \beta_p X_p$$

其中， $\beta_0, \beta_1, \beta_2, \dots, \beta_p$ 是 $(p + 1)$ 个未知的回归参数；

X_1, X_2, \dots, X_p 为 p 个自变量；

Y 为因变量；

ε 为未知的随机扰动误差项。



多元回归分析

✿ 通过对 X_1, X_2, \dots, X_p, Y 进行 n 项独立的观测, 并取得

■ 样本: $(X_{i1}, X_{i2}, \dots, X_{ip}, Y_i) \quad i = 1, \dots, n$

■ 模型: $Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} + \varepsilon_i$

■ 相应的回归方程: $E(Y|X = X_1, X = X_2, \dots, X = X_p) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$

✿ 多元正态线性回归模型:

■ 假设 ε_i 是独立、正态的随机变量, $\varepsilon_i \sim N(0, \sigma^2)$;

■ 假设 p 个自变量为非随机的变量, 并且自变量之间不存在完全的或者接近完全的线性相关性;

■ Y_i 也是一个正态随机变量, 其分布为:

$$Y_i \sim N(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip}, \sigma^2) \quad i = 1, \dots, n$$



多元回归分析

✿ 用矩阵形式讨论多元正态线性回归模型

✿ 设 $k = p + 1$, 模型可表示为: $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{array}{cccc} \mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} & \mathbf{X} = \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1p} \\ 1 & X_{21} & X_{22} & \cdots & X_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{np} \end{bmatrix} & \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} & \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix} \\ n \times 1 & n \times k & k \times 1 & n \times 1 \end{array}$$

因变量向量

设计矩阵

待估参数向量

误差向量



多元回归模型的参数估计

✿ 待估参数的最小平方估计量:

$$(LSE) \hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$



多元回归模型的检验

✿ 变量之间线性关系的检验:

■ 方法1: 利用方差分析表进行F检验

变差来源	自由度df	平方和SS	均方MS	统计量F
回归	p	$SSR = \hat{\beta}'\mathbf{X}'\mathbf{Y} - n\bar{Y}^2$	$MSR = SSR/p$	$F = MSR/MSE$
残差	$n - p - 1$	$SSE = \mathbf{Y}'\mathbf{Y} - \hat{\beta}'\mathbf{X}'\mathbf{Y}$	$MSE = SSE/(n - p - 1)$	
总离差	$n - 1$	$SST = \mathbf{Y}'\mathbf{Y} - n\bar{Y}^2$		

- 假设: $H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$; $H_1: \beta_1\beta_2 \dots \beta_p$ 中至少一个不等于零;
- 若统计量 $F \leq F_\alpha(p, n - p - 1)$, 接受 H_0 , 说明变量之间的线性关系是不显著的;
- 若统计量 $F > F_\alpha(p, n - p - 1)$, 接受 H_1 , 说明变量之间的线性关系是显著的。



多元回归模型的检验

✿ 变量之间线性关系的检验:

■ 方法2: 利用多元测定系数

■ 定义:

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

■ 值域: $0 \leq R^2 \leq 1$

■ R^2 是评价多元回归模型对变量之间线性关系代表性的一个指标, 它的值越大, 说明拟合优度越好

■ 为了去除自变量个数对 R^2 的影响, 对它进行调整:

$$R_a^2 = 1 - \frac{SSE/(n-p-1)}{SST/(n-1)} = 1 - (1 - R^2) \frac{n-1}{n-p-1}$$



多元回归模型的检验

✿ 回归参数的检验:

- 假设: $H_0: \beta_j = \beta_{j0}$ (β_{j0} 是 β_j 的某个假设值, 常取 $\beta_{j0} = 0$) ; $H_1: \beta_j \neq \beta_{j0}$;

- 记

$$S^2(\hat{\boldsymbol{\beta}}) = \widehat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1} = (MSE) \begin{bmatrix} a_{00} & a_{01} & \cdots & a_{0p} \\ a_{10} & a_{11} & \cdots & a_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ a_{p0} & a_{p1} & \cdots & a_{pp} \end{bmatrix}$$

- 检验统计量

$$t = \frac{\hat{\beta}_j - \beta_{j0}}{S(\hat{\beta}_j)} = \frac{\hat{\beta}_j - \beta_{j0}}{\hat{\sigma}\sqrt{a_{jj}}} = \frac{\hat{\beta}_j - \beta_{j0}}{\sqrt{(MSE)a_{jj}}}, \quad j = 0, 1, 2, \dots, p$$

- 若 $|t| \leq t_{\frac{\alpha}{2}, n-p-1}$, 则接受 H_0 ; $|t| > t_{\frac{\alpha}{2}, n-p-1}$, 则接受 H_1 。



多元回归模型的检验

✿ 回归参数的检验:

- β_j 的 $100(1 - \alpha)\%$ 置信区间为:

$$\hat{\beta}_j - t_{\frac{\alpha}{2}, n-p-1} \sqrt{(MSE)a_{jj}} \leq \beta_j \leq \hat{\beta}_j + t_{\frac{\alpha}{2}, n-p-1} \sqrt{(MSE)a_{jj}}$$
$$j = 0, 1, 2, \dots, p$$



多元回归的预测

✿ 均值预测

- $E(Y_0|\mathbf{X}_0)$ 的估计量为:

$$\hat{Y}_0 = \mathbf{x}'_0 \hat{\boldsymbol{\beta}} \quad \mathbf{X}_0 = \begin{bmatrix} 1 \\ X_{01} \\ X_{02} \\ \vdots \\ X_{0p} \end{bmatrix}$$

- $E(Y_0|\mathbf{X}_0)$ 的100(1 - α)%置信区间为:

$$\hat{Y}_0 - t_{\frac{\alpha}{2}, n-p-1} \sqrt{(MSE) \mathbf{X}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}_0} \leq E(Y_0|\mathbf{X}_0) \leq \hat{Y}_0 + t_{\frac{\alpha}{2}, n-p-1} \sqrt{(MSE) \mathbf{X}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}_0}$$



多元回归的预测

✿ 个别值预测:

- 个别值 Y_0 的估计量为:

$$\hat{Y}_0 = \mathbf{x}'_0 \hat{\boldsymbol{\beta}} \quad \mathbf{x}_0 = \begin{bmatrix} 1 \\ X_{01} \\ X_{02} \\ \vdots \\ X_{0p} \end{bmatrix}$$

- Y_0 的 $100(1 - \alpha)\%$ 置信区间为:

$$\begin{aligned} & \hat{Y}_0 - t_{\frac{\alpha}{2}, n-p-1} \sqrt{(MSE)(1 + \mathbf{x}'_0(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_0)} \leq Y_0 \\ & \leq \hat{Y}_0 + t_{\frac{\alpha}{2}, n-p-1} \sqrt{(MSE)(1 + \mathbf{x}'_0(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_0)} \end{aligned}$$



多元回归模型计算举例

✿ 例1: 某企业开发了新型环保产品, 目前在国内有15个销售点。为了分析该产品的市场情况, 决定以销售点所在地区的人口 (万人) 为自变量 X_1 , 以市场地区的人均年收入 (元) 为自变量 X_2 , 以环保产品的年销售量 (套) 为因变量 Y , 并已取得样本数据的初步计算结果。设 Y 与 X_1 、 X_2 的关系为线性且满足多元线性回归模型的各种假定, 取 $\alpha = 0.05$, 试进行多元回归分析。

$$\begin{aligned}n &= 15, \sum X_{i1} = 3626, \sum X_{i2} = 44428, \sum Y_i = 2259, \sum X_{i1}^2 = 1067614, \\ \sum X_{i2}^2 &= 139063428, \sum Y_i^2 = 394107, \sum X_{i1} X_{i2} = 11419181, \\ \sum X_{i1} Y_i &= 647107, \sum X_{i2} Y_i = 7096619\end{aligned}$$



多元回归模型计算举例

✿ 解：根据题意， Y 和 X_1 、 X_2 之间的样本模型为

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i \quad i = 1, 2, 3, \dots, 15$$

■ 其矩阵形式为

$$\mathbf{Y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

■ (1) 估计回归参数

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ X_{11} & X_{21} & \cdots & X_{15,1} \\ X_{12} & X_{22} & \cdots & X_{15,2} \end{bmatrix} \begin{bmatrix} 1 & X_{11} & X_{12} \\ 1 & X_{21} & X_{22} \\ \vdots & \vdots & \vdots \\ 1 & X_{15,1} & X_{15,2} \end{bmatrix} = \begin{bmatrix} n & \sum X_{i1} & \sum X_{i2} \\ \sum X_{i1} & \sum X_{i1}^2 & \sum X_{i1} X_{i2} \\ \sum X_{i2} & \sum X_{i1} X_{i2} & \sum X_{i2}^2 \end{bmatrix}$$

$$= \begin{bmatrix} 15 & 3626 & 44428 \\ 3626 & 1067614 & 11419181 \\ 44428 & 11419181 & 139063428 \end{bmatrix}$$

$$\mathbf{X}'\mathbf{Y} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ X_{11} & X_{21} & \cdots & X_{15,1} \\ X_{12} & X_{22} & \cdots & X_{15,2} \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_{15} \end{bmatrix} = \begin{bmatrix} \sum Y_i \\ \sum X_{i1} Y_i \\ \sum X_{i2} Y_i \end{bmatrix} = \begin{bmatrix} 2259 \\ 647107 \\ 7096619 \end{bmatrix}$$



多元回归模型计算举例

✿ 估计回归参数

- 待估参数向量的最小平方估计量为

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Y}) \\ &= \begin{bmatrix} 1.2463484 & 2.1296642 \times 10^{-4} & -4.1567125 \times 10^{-4} \\ 2.1296642 \times 10^{-4} & 7.7329030 \times 10^{-6} & -7.0302518 \times 10^{-7} \\ -4.1567125 \times 10^{-4} & -7.0302518 \times 10^{-7} & 1.9771851 \times 10^{-7} \end{bmatrix} \times \begin{bmatrix} 2259 \\ 647107 \\ 7096619 \end{bmatrix} \\ &= \begin{bmatrix} 3.4526 \\ 0.4960 \\ 0.0092 \end{bmatrix} = \begin{bmatrix} \widehat{\beta}_0 \\ \widehat{\beta}_1 \\ \widehat{\beta}_2 \end{bmatrix}\end{aligned}$$

- 因此，估计的回归方程为：

$$\hat{Y} = 3.4526 + 0.4960X_1 + 0.0092X_2$$



多元回归模型计算举例

■ (2) 建立方差分析表 (ANOVA)

$$Y'Y = \sum Y_i^2 = 394107 \quad n\bar{Y}^2 = 15 \times \left[\frac{1}{15} \sum Y_i \right]^2 = 15 \times \left[\frac{1}{15} \times 2259 \right]^2 = 340205.4$$

$$\hat{\beta}'X'Y = [3.4526 \quad 0.4960 \quad 0.0092] \begin{bmatrix} 2259 \\ 647107 \\ 7096619 \end{bmatrix} = 394050.116$$

$$SST = Y'Y - n\bar{Y}^2 = 394107 - 340205.4 = 53901.6$$

$$SSE = Y'Y - \hat{\beta}'X'Y = 394107 - 394050.116 = 56.884$$

$$SSR = SST - SSE = 53901.6 - 56.884 = 53844.716$$

变差来源	自由度	平方和	均方	统计量F
回归	2	SSR=53844.716	MSR=26922.358	F=5680
残差	12	SSE=56.884	MSE=4.740	
总离差	14	SST=53901.600	——	——



多元回归模型计算举例

- (3) 检验总体线性关系的显著性
- 设 $H_0: \beta_1 = \beta_2 = 0$; $H_1: \beta_1$ 和 β_2 中至少有一个不等于零
- 因为检验统计量 $F = 5680 \geq F_{0.05}(2,12) = 3.89$, 所以在 $\alpha = 0.05$ 的水平上接受 H_1 , 即该环保产品销售量 Y 跟销售地区人口 X_1 和销售地区人均年收入 X_2 之间存在的线性关系是显著的。
- (4) 考察回归方程的拟合优度
- 计算多元测定系数和调整后的多元测定系数

$$R^2 = \frac{SSR}{SST} = \frac{53844.716}{53901.6} = 0.9989 \quad R_a^2 = 1 - (1 - 0.9989) \left(\frac{15 - 1}{15 - 2 - 1} \right) = 0.9987$$

- 多元测定系数的值较高, 说明该回归方程拟合良好。



多元回归模型计算举例

- (5) 检验回归参数的显著性
- 检验 $H_0: \beta_1=0$, $H_1: \beta_1 \neq 0$
- 因为:

$$a_{11} = 7.732903 \times 10^{-6} \quad \widehat{\sigma}^2 = MSE = 4.74 \quad \widehat{\beta}_1 = 0.496$$

$$t = \frac{\widehat{\beta}_1 - 0}{\sqrt{\widehat{\sigma}^2 a_{11}}} = \frac{0.496}{\sqrt{4.74 \times 7.732903 \times 10^{-6}}} = 81.926 > t_{0.025,12} = 2.179$$

- 所以在 $\alpha = 0.05$ 的水平上接受 H_1 , 即该环保产品销售量 Y 跟销售地区人口数 X_1 之间存在显著的线性关系。 β_1 的95%置信区间为:

$$0.496 - 2.179 \times 0.006054 \leq \beta_1 \leq 0.496 + 2.179 \times 0.006054$$

$$0.4828 \leq \beta_1 \leq 0.5092$$



多元回归模型计算举例

- (5) 检验回归参数的显著性
- 检验 $H_0: \beta_2=0$, $H_1: \beta_2 \neq 0$
- 因为:

$$a_{22} = 1.9771851 \times 10^{-7} \quad \widehat{\sigma}^2 = MSE = 4.74 \quad \widehat{\beta}_1 = 0.0092$$

$$t = \frac{\widehat{\beta}_2 - 0}{\sqrt{\widehat{\sigma}^2 a_{22}}} = \frac{0.0092}{\sqrt{4.74 \times 1.9771851 \times 10^{-7}}} = 9.508 > t_{0.025,12} = 2.179$$

- 所以在 $\alpha = 0.05$ 的水平上接受 H_1 , 即该环保产品销售量 Y 跟销售地区人均年收入 X_2 之间存在显著的线性关系。 β_2 的95%置信区间为:

$$0.0092 - 2.179 \times 0.000968 \leq \beta_2 \leq 0.0092 + 2.179 \times 0.000968$$

$$0.0071 \leq \beta_2 \leq 0.00113$$



多元回归模型计算举例

- (6) 预测
- 该企业准备为它的产品在一个新地区开设销售点。已知新地区的人口为220（万人），人均年收入2500（元），试预测产品在新地区的年销售量 Y 的平均值 $E(Y_0|\mathbf{X}_0)$ 和个别值 Y_0 。
- 预测产品在新地区的年销售量 Y 的平均值 $E(Y_0|\mathbf{X}_0)$ 。
- 已知

$$\mathbf{X}_0 = \begin{bmatrix} 1 \\ 220 \\ 2500 \end{bmatrix}$$

- $E(Y_0|\mathbf{X}_0)$ 的估计值为

$$\hat{Y}_0 = X_0' \hat{\beta} = [1 \quad 220 \quad 2500] \begin{bmatrix} 3.4526 \\ 0.4960 \\ 0.0092 \end{bmatrix} = 135.5726$$



多元回归模型计算举例

- (6) 预测

- 预测产品在新地区的年销售量 Y 的平均值 $E(Y_0|\mathbf{X}_0)$ 。

- 因为

$$\mathbf{X}'_0(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}_0 = [1 \quad 220 \quad 2500] (\mathbf{X}'\mathbf{X})^{-1} \begin{bmatrix} 1 \\ 220 \\ 2500 \end{bmatrix} = 0.09835$$

$$s(\hat{Y}_0) = \sqrt{(MSE)\mathbf{X}'_0(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}_0} = \sqrt{4.74 \times 0.09835} = 0.6828$$

$$t_{0.025,12} = 2.179$$

- 所以， $E(Y_0|\mathbf{X}_0)$ 的95%置信区间为

$$135.5726 - 2.179 \times 0.6828 \leq E(Y_0|\mathbf{X}_0) \leq 135.5726 + 2.179 \times 0.6828$$

$$134.09 \leq E(Y_0|\mathbf{X}_0) \leq 137.06$$



多元回归模型计算举例

- (6) 预测

- 预测产品在新地区的年销售量 Y 的个别值 Y_0 。

- Y_0 的点估计值仍然为

$$\hat{Y}_0 = X_0' \hat{\beta} = [1 \quad 220 \quad 2500] \begin{bmatrix} 3.4526 \\ 0.4960 \\ 0.0092 \end{bmatrix} = 135.5726$$

- 因为

$$s(\hat{Y}_0^*) = \sqrt{(MSE)[1 + \mathbf{X}_0'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}_0]} = \sqrt{4.74 \times 1.09835} = 2.2817$$

- 所以, Y_0 的95%置信区间为

$$135.5726 - 2.179 \times 2.2817 \leq Y_0 \leq 135.5726 + 2.179 \times 2.2817$$

$$130.60 \leq Y_0 \leq 140.54$$



多元回归模型计算举例

- ✿ 例2：已知历年实测的湖水中化学耗氧量COD的浓度与相应的环境自然经济资料。根据已知数据分析COD与环境、自然和经济状况之间的关系。

年份	项目 COD浓度 y (mg/l)	农业产量 x_1 (亿斤)	工业总产值 x_2 (亿元)	湖泊水位 x_3 (米)
1960	2.50	0.25	4.00	3.17
1975	3.63	0.92	21.10	3.24
1976	3.15	0.87	29.10	3.02
1977	2.52	0.60	33.00	3.24
1978	4.06	0.63	37.50	2.63
1979	3.72	0.65	42.40	2.80
1980	2.82	0.42	49.25	3.65
1981	3.31	0.40	50.00	2.97



多元回归模型计算举例

✿ 解:

- (1) 建立回归模型:
- 设: $y = b_0 + b_1x_1 + b_2x_2 + b_3x_3 + \varepsilon$
- 其中, b_0 、 b_1 、 b_2 、 b_3 是待估参数, ε 是误差项。
- (2) 求回归方程:
- 经计算获得各参数的估计量和回归方程 (计算过程略)
$$\widehat{b}_0 = 5.6609, \widehat{b}_1 = 0.9613, \widehat{b}_2 = 0.011, \widehat{b}_3 = -1.0948$$
$$\widehat{y} = 5.6609 + 0.9613x_1 + 0.0011x_2 - 1.0948x_3$$
- (3) 显著性检验 (给定 $\alpha = 0.25$)
- 方差分析



多元回归模型计算举例

■ 方差分析表

方差来源	自由度	平方和	均方	F值
回归	3	1.6210	0.5373	3.1022
残差	4	0.6928	0.1732	
总计	7	2.3048		

- 因为 $F = 3.1022 > F_{0.25}(3,4) = 2.05$ ，所以回归方程总体线性显著。
- 回归系数的显著性检验
- （计算过程略）在给定的置信水平下， x_1 、 x_3 对 y 的影响显著， x_2 对 y 的影响不显著。
- 将 x_2 从模型中剔除，在 y 与 x_1 、 x_3 之间重新进行回归分析。



多元回归模型计算举例

- (4) 在 y 与 x_1 、 x_3 之间建立回归方程
- (计算过程略) $\hat{y} = 6.1297 + 0.9411x_1 - 1.1241x_3$
- (5) 再次进行显著性检验
- 方差分析

方差来源	自由度	平方和	均方	F值
回归	2	1.4099	0.7050	3.9387
残差	5	0.8948	0.1790	
总计	7	2.3048		

- 因为 $F = 3.9387 > F_{0.25}(2,5) = 1.85$ ，所以新的回归方程显著。



多元回归模型计算举例

- (5) 再次进行显著性检验
- 回归系数的显著性检验
- (计算过程略) 在给定的置信水平下, x_1 、 x_3 对 y 的影响显著, 新的回归方程有实际意义。
- (6) 课后练习:
- 基于该课程案例, 给出完整的多元回归过程。



回归分析

✿ 基本概念

✿ 一元回归分析

- 回归参数的估计
- 回归模型的检验
- 回归预测

✿ 多元回归分析

✿ 回归函数的形式



回归函数的形式

✿ 如果自变量与因变量之间的关系为某种曲线关系，但是回归分析中采用的回归函数是直线形式，那么

- 模型对样本数据的代表性就很差。这一点会反映在：
 - 统计检验
 - 残差图

✿ 纠正办法之一：

- 对样本数据进行某种变换，以使回归模型对于变换后的样本数据，其直线函数关系变得适宜；
- transformation on the data



回归函数的形式

✿ 常用的变换方法:

■ 对数变换:

- 半对数: 只对变量 X 和 Y 中的一个进行对数变换。

$$Y = a_0 a_1^X \varepsilon \Rightarrow \lg Y = \lg a_0 + X \lg a_1 + \lg \varepsilon \Rightarrow Y' = \beta_0 + \beta_1 X + \varepsilon'$$

- 双对数: 对变量 X 和 Y 都进行对数变换。

$$Y = \beta_0 X^{\beta_1} \varepsilon \Rightarrow \lg Y = \lg \beta_0 + \beta_1 \lg X + \lg \varepsilon \Rightarrow Y' = \alpha + \beta_1 X' + \varepsilon'$$

■ 倒数变换

$$Y = \beta_0 + \frac{\beta_1}{X} + \varepsilon \Rightarrow Y = \beta_0 + \beta_1 X' + \varepsilon$$



回归函数的形式

✿ 常用的变换方法:

■ 对数和倒数联合变换:

$$Y = \frac{1}{a + b \exp(\beta_0 + \beta_1 X + \varepsilon)} \Rightarrow Y' = \ln \left[\left(\frac{1}{Y} - a \right) / b \right] \Rightarrow Y' = \beta_0 + \beta_1 X + \varepsilon$$

■ 多项式回归

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \cdots + \beta_m X^m + \varepsilon \Rightarrow Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_m X_m + \varepsilon$$



回归分析总结

✿ 基本概念

- 起源、研究任务、变量之间的统计关系、回归模型、回归函数、PRF、SRF

✿ 回归参数估计

- 最小平方法

✿ 回归模型检验

- 总体线性检验：方差分析、残差图、多元测定系数
- 参数检验：统计假设检验、置信区间

✿ 回归预测

- 均值预测、个别值预测





中山大學
SUN YAT-SEN UNIVERSITY

谢谢!