

第 11 章 二值选择模型

11.1 二值选择模型

如果被解释变量 y 离散，称为“离散选择模型” (discrete choice model) 或“定性反应模型” (qualitative response model)。

最常见的离散选择模型是二值选择行为(binary choices)。

比如：考研或不考研；就业或待业；买房或不买房；买保险或不买保险；贷款申请被批准或拒绝；出国或不出国；回国或不回

国；战争或和平；生或死。

假设个体只有两种选择，比如 $y = 1$ (考研)或 $y = 0$ (不考研)。

最简单的建模方法为“线性概率模型”(Linear Probability Model, LPM):

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_K x_{iK} + \varepsilon_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i \quad (i = 1, \dots, n) \quad (11.1)$$

其中，解释变量 $\mathbf{x}_i \equiv (x_{i1} \ x_{i2} \ \cdots \ x_{iK})'$ ，而参数 $\boldsymbol{\beta} \equiv (\beta_1 \ \beta_2 \ \cdots \ \beta_K)'$ 。

LPM 的优点是，计算方便，容易得到边际效应(即回归系数)。

LPM 的缺点是，虽然 y 的取值非 0 即 1，但根据线性概率模型所作的预测值却可能出现 $\hat{y} > 1$ 或 $\hat{y} < 0$ 的不现实情形。

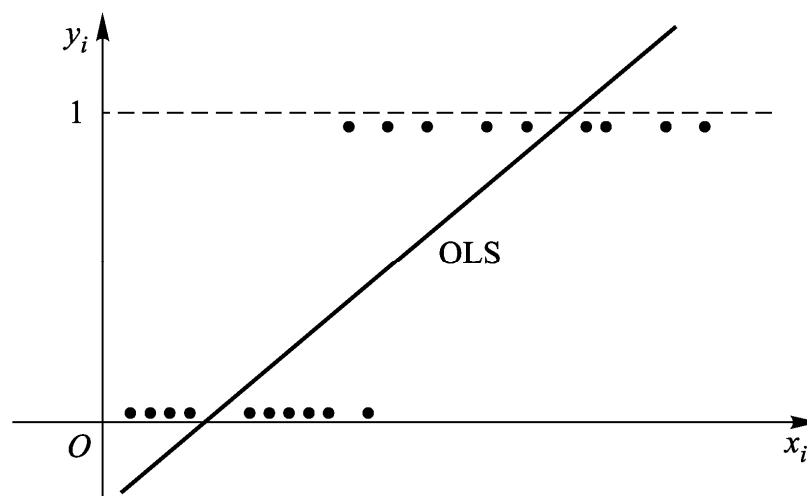


图 11.1 线性概率模型

为使 y 的预测值介于 $[0, 1]$ 之间，在给定 \mathbf{x} 的情况下，考虑 y 的两点分布概率：

$$\begin{cases} P(y = 1 | \mathbf{x}) = F(\mathbf{x}, \boldsymbol{\beta}) \\ P(y = 0 | \mathbf{x}) = 1 - F(\mathbf{x}, \boldsymbol{\beta}) \end{cases} \quad (11.2)$$

函数 $F(\mathbf{x}, \boldsymbol{\beta})$ 称为“连接函数” (link function)，因为它将 \mathbf{x} 与 y 连接起来。

y 的取值要么为 0，要么为 1，故 y 肯定服从两点分布。

连接函数的选择具有一定灵活性。

通过选择合适的连接函数 $F(\mathbf{x}, \boldsymbol{\beta})$ (比如，某随机变量的累积分布函数)，可保证 $0 \leq \hat{y} \leq 1$ ，并将 \hat{y} 理解为“ $y = 1$ ”发生的概率，因为

$$E(y | \mathbf{x}) = 1 \cdot P(y = 1 | \mathbf{x}) + 0 \cdot P(y = 0 | \mathbf{x}) = P(y = 1 | \mathbf{x}) \quad (11.3)$$

如果 $F(\mathbf{x}, \boldsymbol{\beta})$ 为标准正态的累积分布函数，则

$$P(y = 1 | \mathbf{x}) = F(\mathbf{x}, \boldsymbol{\beta}) = \Phi(\mathbf{x}'\boldsymbol{\beta}) \equiv \int_{-\infty}^{\mathbf{x}'\boldsymbol{\beta}} \phi(t) dt \quad (11.4)$$

$\phi(\cdot)$ 与 $\Phi(\cdot)$ 分别为标准正态的密度与累积分布函数；此模型称为“Probit”。

如果 $F(\mathbf{x}, \boldsymbol{\beta})$ 为“逻辑分布”(logistic distribution)的累积分布函数，则

$$P(y = 1 | \mathbf{x}) = F(\mathbf{x}, \boldsymbol{\beta}) = \Lambda(\mathbf{x}'\boldsymbol{\beta}) \equiv \frac{\exp(\mathbf{x}'\boldsymbol{\beta})}{1 + \exp(\mathbf{x}'\boldsymbol{\beta})} \quad (11.5)$$

其中，函数 $\Lambda(\cdot)$ 的定义为 $\Lambda(z) \equiv \frac{\exp(z)}{1 + \exp(z)}$ ；此模型称为“Logit”。

逻辑分布的密度函数关于原点对称，期望为 0，方差为 $\pi^2/3$ (大于标准正态的方差)，具有厚尾(fat tails)。

Probit 与 Logit 都很常用，二者的估计结果(比如边际效应)通常很接近。

Logit 模型的优势在于，逻辑分布的累积分布函数有解析表达式(标准正态没有)，故计算 Logit 更为方便；而且 Logit 的回归系数更易解释其经济意义。

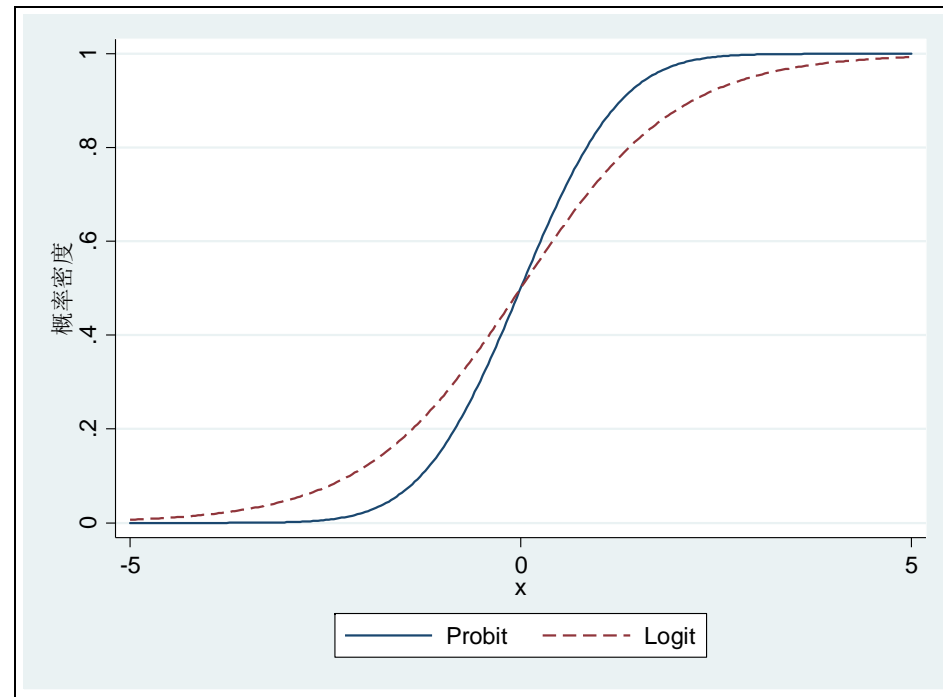


图 11.2 标准正态分布与逻辑分布的累积分布函数

11.2 最大似然估计的原理

Probit 与 Logit 模型本质上都是非线性模型，无法通过变量转换变为线性模型。

对于非线性模型，常使用最大似然估计法(Maximum Likelihood Estimation, MLE 或 ML)。

回顾概率统计中的最大似然估计法。

假设随机变量 y 的概率密度函数为 $f(y; \theta)$ ，其中 θ 为未知参数。

为估计 θ ，从 y 的总体中抽取样本容量为 n 的随机样本 $\{y_1, \dots, y_n\}$ 。

假设 $\{y_1, \dots, y_n\}$ 为 iid，样本数据的联合密度函数为

$$f(y_1; \theta) f(y_2; \theta) \cdots f(y_n; \theta) = \prod_{i=1}^n f(y_i; \theta) \quad (11.6)$$

其中， $\prod_{i=1}^n$ 表示连乘。

在抽样之前， $\{y_1, \dots, y_n\}$ 为随机向量。

抽样之后， $\{y_1, \dots, y_n\}$ 有了特定的样本值。

可将样本的联合密度函数视为在给定 $\{y_1, \dots, y_n\}$ 情况下，未知参数 θ 的函数。

定义似然函数(likelihood function)为

$$L(\theta; y_1, \dots, y_n) = \prod_{i=1}^n f(y_i; \theta) \quad (11.7)$$

似然函数与联合密度函数完全相等，只是 θ 与 $\{y_1, \dots, y_n\}$ 的角色互换，即把 θ 作为自变量，视 $\{y_1, \dots, y_n\}$ 为给定。

为运算方便，把似然函数取对数：

$$\ln L(\theta; y_1, \dots, y_n) = \sum_{i=1}^n \ln f(y_i; \theta) \quad (11.8)$$

MLE 的思想：给定样本取值后，该样本最可能来自参数 θ 为何值的总体。

寻找 $\hat{\theta}_{\text{ML}}$ ，使得观测到样本数据的可能性最大，即最大化对数似然函数(loglikelihood function)：

$$\max_{\theta} \ln L(\theta; y_1, \dots, y_n) \quad (11.9)$$

假设存在唯一内点解，一阶条件为

$$\frac{\partial \ln L(\theta; y_1, \dots, y_n)}{\partial \theta} = 0 \quad (11.10)$$

求解一阶条件，可得最大似然估计量 $\hat{\theta}_{\text{ML}}$ 。

例 假设 $y \sim N(\mu, \sigma^2)$ ，其中 σ^2 已知，得到样本容量为 1 的样本 $y_1 = 2$ ，求对 μ 的最大似然估计。根据正态分布的密度函数，此样本的似然函数为

$$L(\mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{\frac{-(2-\mu)^2}{2\sigma^2}\right\} \quad (11.11)$$

似然函数在 $\hat{\mu} = 2$ 处取最大值，参见图 11.3。

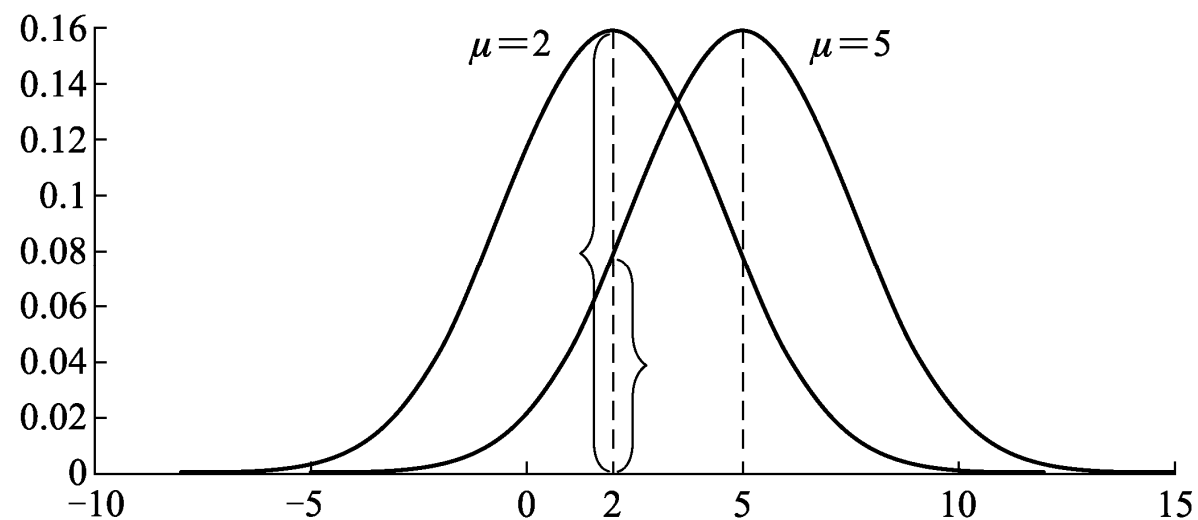


图 11.3 选择参数使观测到样本的可能性最大

例(非正式) 某人操一口浓重的四川口音，则判断他最有可能来自四川。

在一定的正则条件(regularity conditions)下, MLE 估计量具有良好的大样本性质, 可照常进行大样本统计推断。

(1) $\hat{\theta}_{\text{ML}}$ 为一致估计, 即 $\text{plim} \hat{\theta}_{\text{ML}} = \theta$ 。

(2) $\hat{\theta}_{\text{ML}}$ 服从渐近正态分布。

(3) 在大样本下, $\hat{\theta}_{\text{ML}}$ 是最有效率的估计(渐近方差最小)。

由于模型存在非线性, MLE 通常没有解析解, 只能寻找“数值解”(numerical solution)。

一般使用“迭代法”(iteration)进行数值求解。

常用的迭代法为“高斯-牛顿法”(Gauss-Newton method)。

MLE 的一阶条件可归结为求非线性方程 $f(x) = 0$ 的解。

假设 $f(x)$ 的导数 $f'(x)$ 处处存在，参见图 11.4。记该方程的解为 x^* ，满足 $f(x^*) = 0$ 。

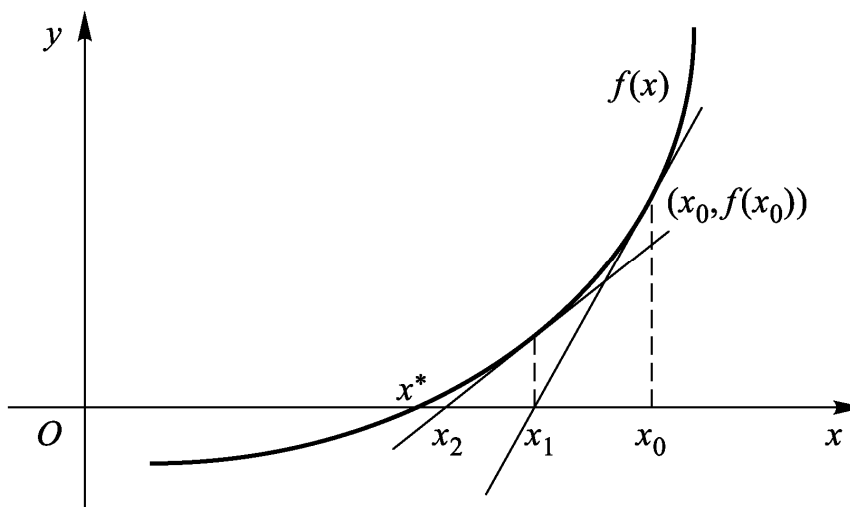


图 11.4 高斯-牛顿法

首先，猜初始值 x_0 ，在点 $(x_0, f(x_0))$ 处作曲线 $f(x)$ 的切线，记切线与横轴的交点为 x_1 。

然后，在点 $(x_1, f(x_1))$ 处再作切线，记切线与横轴的交点为 x_2 。

以此类推，不断迭代，可得序列 $\{x_0, x_1, x_2, x_3, \dots\}$ 。

一般情况下，该序列将收敛至 x^* (给定一个精确度，收敛到精确度范围内即停止)。

高斯-牛顿法的收敛速度很快，是二次的。如本次迭代的误差为 0.1，则下次迭代的误差约为 0.1^2 ，下下次迭代的误差约为 0.1^4 ，等等。

如果初始值 x_0 选择不当, 也可能出现迭代不收敛的情形。

使用牛顿法得到的可能只是“局部最大值”(local maximum), 而非“整体最大值”(global maximum)。

MLE 很容易应用于多参数的情形。

假设随机变量 y 的概率密度函数为 $f(y; \boldsymbol{\theta})$, 其中 $\boldsymbol{\theta} = (\theta_1 \ \theta_2)'$, 则对数似然函数为

$$\ln L(\boldsymbol{\theta}; y_1, \dots, y_n) = \sum_{i=1}^n \ln f(y_i; \boldsymbol{\theta}) \quad (11.12)$$

一阶条件为

$$\begin{cases} \frac{\partial \ln L(\boldsymbol{\theta}; y_1, \dots, y_n)}{\partial \theta_1} = 0 \\ \frac{\partial \ln L(\boldsymbol{\theta}; y_1, \dots, y_n)}{\partial \theta_2} = 0 \end{cases} \quad (11.13)$$

求解此联立方程组，可得最大似然估计量 $\hat{\theta}_{1, \text{ML}}$ 与 $\hat{\theta}_{2, \text{ML}}$ 。

高斯-牛顿法也适用于多元函数 $f(\mathbf{x})=0$ 的情形，只要在上述迭代过程中，将切线替换为(超)切平面即可。

11.3 二值选择模型的 MLE 估计

以 Logit 为例，将 MLE 应用于二值选择模型。

对于样本数据 $\{\mathbf{x}_i, y_i\}_{i=1}^n$ ，根据方程(11.5)，第 i 个观测数据的概率密度为

$$f(y_i | \mathbf{x}_i, \boldsymbol{\beta}) = \begin{cases} \Lambda(\mathbf{x}_i' \boldsymbol{\beta}), & \text{若 } y_i = 1 \\ 1 - \Lambda(\mathbf{x}_i' \boldsymbol{\beta}), & \text{若 } y_i = 0 \end{cases} \quad (11.14)$$

其中， $\Lambda(z) \equiv \frac{\exp(z)}{1 + \exp(z)}$ 为逻辑分布的累积分布函数。

上式可写为

$$f(y_i | \mathbf{x}_i, \boldsymbol{\beta}) = [\Lambda(\mathbf{x}_i' \boldsymbol{\beta})]^{y_i} [1 - \Lambda(\mathbf{x}_i' \boldsymbol{\beta})]^{1-y_i} \quad (11.15)$$

取对数可得

$$\ln f(y_i | \mathbf{x}_i, \boldsymbol{\beta}) = y_i \ln[\Lambda(\mathbf{x}_i' \boldsymbol{\beta})] + (1 - y_i) \ln[1 - \Lambda(\mathbf{x}_i' \boldsymbol{\beta})] \quad (11.16)$$

假设样本中的个体相互独立，整个样本的对数似然函数为

$$\ln L(\boldsymbol{\beta} | \mathbf{y}, \mathbf{x}) = \sum_{i=1}^n y_i \ln[\Lambda(\mathbf{x}_i' \boldsymbol{\beta})] + \sum_{i=1}^n (1 - y_i) \ln[1 - \Lambda(\mathbf{x}_i' \boldsymbol{\beta})] \quad (11.17)$$

把对数似然函数对 $\boldsymbol{\beta}$ 求偏导，可得一阶条件。

满足一阶条件的估计量即为 MLE 估计量，记为 $\hat{\boldsymbol{\beta}}_{\text{ML}}$ 。

11.4 边际效应

对于线性模型，回归系数 β_k 的经济意义就是变量 x_k 对 y 的边际效应(marginal effects)。

在非线性模型中，估计量 $\hat{\beta}_{ML}$ 一般并非边际效应。

以 Probit 为例，计算变量 x_k 的边际效应：

$$\frac{\partial P(y=1|\mathbf{x})}{\partial x_k} = \frac{\partial \Phi(\mathbf{x}'\boldsymbol{\beta})}{\partial x_k} = \frac{\partial \Phi(\mathbf{x}'\boldsymbol{\beta})}{\partial(\mathbf{x}'\boldsymbol{\beta})} \cdot \frac{\partial(\mathbf{x}'\boldsymbol{\beta})}{\partial x_k} = \phi(\mathbf{x}'\boldsymbol{\beta}) \cdot \beta_k \quad (11.18)$$

由于 Probit 与 Logit 所用分布函数不同，其参数估计值不直接可比。需分别计算二者的边际效应，然后比较。

对于非线性模型，边际效应通常不是常数，随着向量 \mathbf{x} 而变。

非线性模型常用的边际效应概念：

(1) 平均边际效应 (average marginal effect): 分别计算在每个样本观测值上的边际效应，然后进行简单算术平均。

(2) 样本均值处的边际效应 (marginal effect at mean): 计算在 $\mathbf{x} = \bar{\mathbf{x}}$ 处的边际效应。

(3) 在某代表值处的边际效应 (marginal effect at a representative value): 给定 \mathbf{x}^* ，计算在 $\mathbf{x} = \mathbf{x}^*$ 处的边际效应。

以上三种边际效应的计算结果可能有较大差异。

传统上，常计算样本均值处 $\mathbf{x} = \bar{\mathbf{x}}$ 的边际效应，因为计算方便。

但在非线性模型中，样本均值处的个体行为并不等于样本中个体的平均行为(average behavior of individuals differs from behavior of the average individual)。

对于政策分析而言，使用平均边际效应(Stata 的默认方法)，或在某代表值处的边际效应通常更有意义。

11.5 回归系数的经济意义

$\hat{\beta}_{\text{ML}}$ 并非边际效应，它究竟有什么含义？

对于 Logit 模型，记事件发生的概率为 $p \equiv P(y = 1 | \mathbf{x})$ ，则事件不发生的概率为 $1 - p = P(y = 0 | \mathbf{x})$ 。

由于 $p = \frac{\exp(\mathbf{x}'\boldsymbol{\beta})}{1 + \exp(\mathbf{x}'\boldsymbol{\beta})}$ ， $1 - p = \frac{1}{1 + \exp(\mathbf{x}'\boldsymbol{\beta})}$ ，故事件发生与不发生的几率比为

$$\frac{p}{1 - p} = \exp(\mathbf{x}'\boldsymbol{\beta}) \quad (11.19)$$

$\frac{p}{1 - p}$ 称为“几率比” (odds ratio) 或“相对风险” (relative risk)。

例 在检验药物疗效的随机实验中，“ $y=1$ ”表示“生”，“ $y=0$ ”表示“死”。如几率比为 2，意味着存活概率是死亡概率的两倍。

对方程(11.19)两边取对数：

$$\ln\left(\frac{p}{1-p}\right) = \mathbf{x}'\boldsymbol{\beta} = \beta_1 x_1 + \cdots + \beta_K x_K \quad (11.20)$$

$\ln\left(\frac{p}{1-p}\right)$ 称为“对数几率比” (log-odds ratio)。

回归系数 $\hat{\beta}_j$ 表示，变量 x_j 增加一个微小量引起对数几率比的边际变化。

取对数意味着百分比的变化，故可把 $\hat{\beta}_j$ 视为半弹性 (semi-elasticity)，即 x_j 增加一单位引起几率比 $\left(\frac{p}{1-p}\right)$ 的变化百分比。

例 $\hat{\beta}_j = 0.12$ ，意味着 x_j 增加一单位引起几率比增加 12%。

如 x_j 为离散变量(比如，性别、子女数)，可使用另一解释法。

假设 x_j 增加一单位，从 x_j 变为 x_j+1 ，记几率比 p 的新值为 p^* ，则新几率比与原几率比的比率可写为 (无法用微积分)

$$\frac{\frac{p^*}{1-p^*}}{\frac{p}{1-p}} = \frac{\exp[\beta_1 + \beta_2 x_2 + \cdots + \beta_j (x_j + 1) + \cdots + \beta_K x_K]}{\exp(\beta_1 + \beta_2 x_2 + \cdots + \beta_j x_j + \cdots + \beta_K x_K)} = \exp(\beta_j) \quad (11.21)$$

$\exp(\hat{\beta}_j)$ 表示变量 x_j 增加一单位引起几率比的变化倍数。

Stata 也称 $\exp(\hat{\beta}_j)$ 为几率比(odds ratio)。

例 $\hat{\beta}_j = 0.12$ ，则 $\exp(\hat{\beta}_j) = e^{0.12} = 1.13$ ，故当 x_j 增加一单位时，新几率比是原几率比的 1.13 倍，或增加 13%，因为 $\exp(\hat{\beta}_j) - 1 = 1.13 - 1 = 0.13$ 。

如果 $\hat{\beta}_j$ 较小, 则 $\exp(\hat{\beta}_j) - 1 \approx \hat{\beta}_j$ (将 $\exp(\hat{\beta}_j)$ 泰勒展开), 以上两种方法基本等价。

对于 Probit 模型, 无法对其系数 $\hat{\beta}_{ML}$ 进行类似解释。

11.6 拟合优度

不存在平方和分解公式, 无法计算 R^2 。

Stata 仍汇报“准 R^2 ”或“伪 R^2 ”(Pseudo R^2), 由 McFadden (1974) 提出, 定义为

$$\text{准}R^2 \equiv \frac{\ln L_0 - \ln L_1}{\ln L_0} \quad (11.22)$$

$\ln L_1$ 为原模型的对数似然函数之最大值, $\ln L_0$ 为以常数项为唯一解释变量的对数似然函数之最大值。

由于 y 为两点分布, 似然函数的最大可能值为 1 (取值概率为 1), 故对数似然函数的最大可能值为 0, 记为 $\ln L_{\max}$ 。

由于 $\ln L_{\max} = 0$, 可将 “准 R^2 ” 写为

$$\text{准}R^2 = \frac{\ln L_1 - \ln L_0}{\ln L_{\max} - \ln L_0} \quad (11.23)$$

显然, $0 \geq \ln L_1 \geq \ln L_0$, 而 $0 \leq \text{准}R^2 \leq 1$, 参见图 11.5。

分子为对数似然函数的实际增加值($\ln L_1 - \ln L_0$)；分母为对数似然函数的最大可能增加值($\ln L_{\max} - \ln L_0$)。

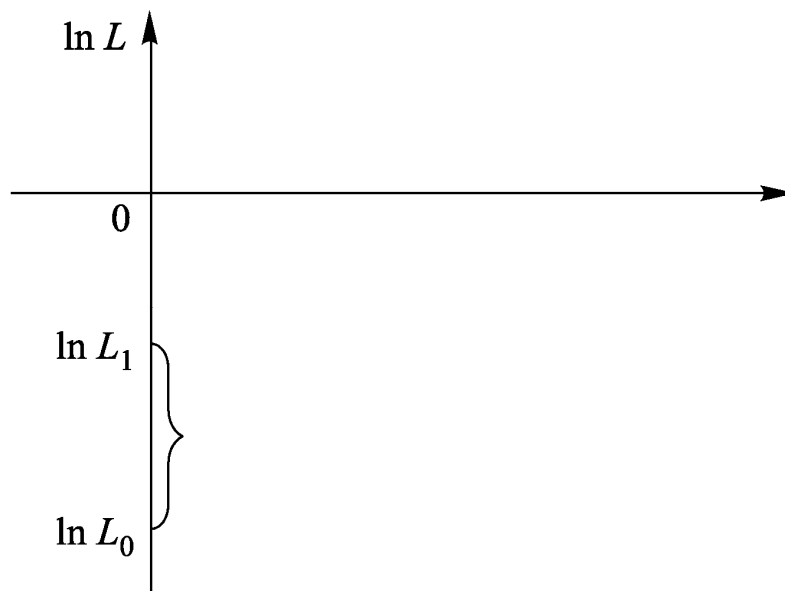


图 11.5 准 R^2 的计算

判断拟合优度的另一方法是计算“正确预测的百分比”(percent correctly predicted)。

如发生概率的预测值 $\hat{y} \geq 0.5$ ，则认为其预测 $y = 1$ ；

反之，则认为其预测 $y = 0$ 。

将预测值与实际值(样本数据)进行比较，可计算正确预测的百分比。

11.7 准最大似然估计

使用 MLE 的前提是对总体的分布函数作具体的假定。

Probit 与 Logit 模型分别假设 y 的两点分布概率为标准正态或逻辑分布的累积分布函数。

此分布函数的设定可能不正确,即存在“设定误差”(specification error)。

定义 使用不正确的分布函数所得到的最大似然估计量,称为“准最大似然估计”(Quasi MLE, 简记 QMLE)或“伪最大似然估计”(Pseudo MLE)。

准最大似然估计是否一定不一致?

不一定！

例 假设线性模型的扰动项服从正态分布，则 MLE 估计量与 OLS 估计量完全相同，而 OLS 估计量的一致性并不依赖于关于分布函数的具体假设。

关于 QMLE 估计量的标准误，可分两种情况考虑。

(1) 如果 QMLE 为一致估计量，由于可能存在对分布函数的设定误差，应使用稳健标准误(robust standard errors)，即相对于模型设定稳健的标准误。

此稳健标准误与异方差稳健的标准误是一致的，因为扰动项方差是否相同也是一种模型设定。

(2) 如 QMLE 估计量不一致, 即使采用稳健标准误也无济于事。

QMLE 估计量 $\hat{\boldsymbol{\beta}}_{\text{QML}} \xrightarrow{p} \boldsymbol{\beta}^* \neq \boldsymbol{\beta}$, 应首先担心估计量的一致性。

稳健标准误只是一致地估计了一个不一致估计量的方差(a consistent estimator of the variance of an inconsistent estimator)。

对于二值选择模型 (Probit 或 Logit), 只要条件期望函数 $E(y | \mathbf{x}) = F(\mathbf{x}, \boldsymbol{\beta})$ 设定正确, 则 MLE 估计就一致。

由于两点分布的特殊性, 在 iid 的情况下, 只要 $E(y | \mathbf{x}) = F(\mathbf{x}, \boldsymbol{\beta})$ 成立, 稳健标准误就等于普通标准误。

如果认为模型设定正确，就不必使用稳健标准误(使用稳健标准误也没错)。

如果模型设定不正确(即 $E(y | \mathbf{x}) \neq F(\mathbf{x}, \boldsymbol{\beta})$)，则 Probit 与 Logit 模型不能得到对系数 $\boldsymbol{\beta}$ 的一致估计，使用稳健标准误就没有太大意义；首先应解决参数估计的一致性问题。

对于二值选择模型，使用普通标准误或稳健标准误都可（文献中无定论）。

11.8 三类渐近等价的大样本检验

在计量中，常使用三类在大样本下渐近等价的统计检验。

考虑线性回归模型：

$$y_i = \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_K x_K + \varepsilon_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i \quad (i = 1, \dots, n) \quad (11.24)$$

其中，解释变量 $\mathbf{x} \equiv (x_1 \ x_2 \ \cdots \ x_K)'$ ，参数 $\boldsymbol{\beta} \equiv (\beta_1 \ \beta_2 \ \cdots \ \beta_K)'$ 。

检验以下原假设：

$$H_0 : \boldsymbol{\beta} = \boldsymbol{\beta}_0 \quad (11.25)$$

其中， $\boldsymbol{\beta}_0$ 已知，共有 K 个约束。

(1) 沃尔德检验(Wald Test)

沃尔德检验考察 $\boldsymbol{\beta}$ 的无约束估计量 $\hat{\boldsymbol{\beta}}$ 与 $\boldsymbol{\beta}_0$ 的距离。

基本思想：如果 H_0 正确，则 $(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$ 不应该很大。

由于 $(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$ 为多维向量，使用二次型：

$$W \equiv (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)' [\text{Var}(\hat{\boldsymbol{\beta}})]^{-1} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) \xrightarrow{d} \chi^2(K) \quad (11.26)$$

其中 K 为约束条件的个数。

第 5-6 章介绍的单一系数 t 检验、联合线性假设的 F 检验都是 Wald 检验。

(2) 似然比检验(Likelihood Ratio Test, 简记 LR)

似然比检验比较无约束估计量 $\hat{\beta}$ 与有约束估计量 $\hat{\beta}^*$ 的差别。

无约束的似然函数最大值 $\ln L(\hat{\beta})$ 比有约束的似然函数最大值 $\ln L(\hat{\beta}^*)$ 更大, 因为在无约束条件下的参数空间 Θ 比有约束条件下(即 H_0 成立时)参数的取值范围更大, 参见图 11.6。

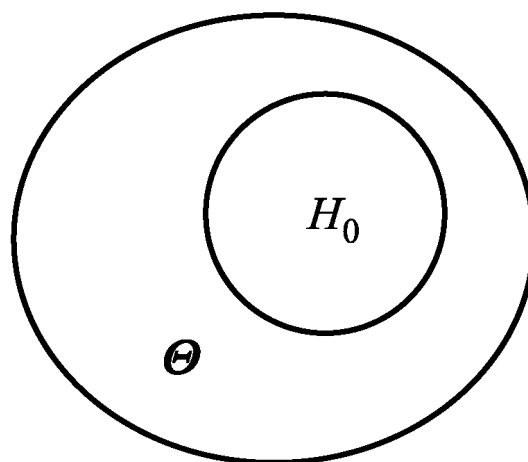


图 11.6 无约束与有约束的参数空间

基本思想：如果 H_0 正确，则 $[\ln L(\hat{\boldsymbol{\beta}}) - \ln L(\hat{\boldsymbol{\beta}}^*)]$ 不应该很大。

在此例中，有约束的估计量 $\hat{\boldsymbol{\beta}}^* = \boldsymbol{\beta}_0$ 。

LR 统计量为

$$LR \equiv -2\ln\left[\frac{L(\hat{\boldsymbol{\beta}}^*)}{L(\hat{\boldsymbol{\beta}})}\right] = 2\left[\ln L(\hat{\boldsymbol{\beta}}) - \ln L(\hat{\boldsymbol{\beta}}^*)\right] \xrightarrow{d} \chi^2(K) \quad (11.27)$$

在大样本下， LR 统计量也服从渐近 $\chi^2(K)$ 分布。

第 5 章介绍的 F 统计量的另一表达式 $F = \frac{(\text{SSR}^* - \text{SSR}) / (K - 1)}{\text{SSR} / (n - K)}$,

即依据似然比原理而设计。

在进行 Probit 或 Logit 回归时，Stata 会汇报一个似然比统计量，检验除常数项外所有参数的联合显著性。

(3) 拉格朗日乘子检验(Lagrange Multiplier Test, 简记 LM)

Wald 检验只考察无约束估计量 $\hat{\beta}$ 。

LR 检验同时考察无约束估计量 $\hat{\beta}$ 与有约束估计量 $\hat{\beta}^*$ 。

LM 检验则只考察有约束估计量 $\hat{\beta}^*$ 。

有约束条件的对数似然函数最大化问题：

$$\begin{aligned} \max_{\tilde{\beta}} \ln L(\tilde{\beta}) \\ s.t. \tilde{\beta} = \beta_0 \end{aligned} \quad (11.28)$$

$\tilde{\beta}$ 为在最大化过程中假想的参数 β 取值(hypothetical value)。

对于约束极值问题，引入拉格朗日乘子函数：

$$\max_{\tilde{\beta}, \lambda} \ln L(\tilde{\beta}) - \lambda'(\tilde{\beta} - \beta_0) \quad (11.29)$$

λ 为拉格朗日乘子向量(Lagrange Multiplier)，其经济含义为约束条件(比如资源约束)的影子价格(shadow price)。

如果 $\hat{\lambda} = \mathbf{0}$ ，约束条件完全不起作用(可无偿获取任意数量的资源)。

根据一阶条件(对 $\tilde{\beta}$ 求导)可知，

$$\hat{\lambda} = \frac{\partial \ln L(\hat{\beta}^*)}{\partial \tilde{\beta}} \equiv \begin{pmatrix} \frac{\partial \ln L(\hat{\beta}^*)}{\partial \tilde{\beta}_1} \\ \vdots \\ \frac{\partial \ln L(\hat{\beta}^*)}{\partial \tilde{\beta}_K} \end{pmatrix} \quad (11.30)$$

最优的拉格朗日乘子向量 $\hat{\lambda}$ 等于对数似然函数在约束估计量 $\hat{\beta}^*$ 处的一阶偏导数(切线的斜率)。

如 $\hat{\lambda} \approx \mathbf{0}$ ，说明约束条件不“紧”(tight)或不是“硬约束”(binding constraint)，加上约束条件不会使似然函数的最大值下降很多，即原假设 H_0 很可能成立。

如果原假设 H_0 成立，则 $(\hat{\lambda} - \theta)$ 的绝对值不应很大。

以二次型来度量此距离，可得 LM 统计量：

$$LM \equiv \hat{\lambda}' [\text{Var}(\hat{\lambda})]^{-1} \hat{\lambda} \xrightarrow{d} \chi^2(K) \quad (11.31)$$

其中， $\text{Var}(\hat{\lambda})$ 为 $\hat{\lambda}$ 的协方差矩阵。

由于 $\hat{\lambda} = \frac{\partial \ln L(\tilde{\beta})}{\partial \tilde{\beta}}$ 称为“得分函数”(score function)或“得分向量”(score vector)，此检验也称“得分检验”(score test)。

另一直观理解是，由于在无约束估计量 $\hat{\beta}$ 处， $\frac{\partial \ln L(\hat{\beta})}{\partial \hat{\beta}} = \mathbf{0}$ (MLE

的一阶条件), 故如原假设 H_0 成立, 则在约束估计量 $\hat{\beta}^*$ 处,
 $\frac{\partial \ln L(\hat{\beta}^*)}{\partial \tilde{\beta}} \approx \mathbf{0}$, 而 LM 统计量反映的就是此接近程度。

在第 7-8 章, 对异方差与自相关所进行的 nR^2 形式的检验都来自 LM 检验的推导。

这三类检验在大样本下渐近等价, 从不同侧面考察同一事物, 参见图 11.7。

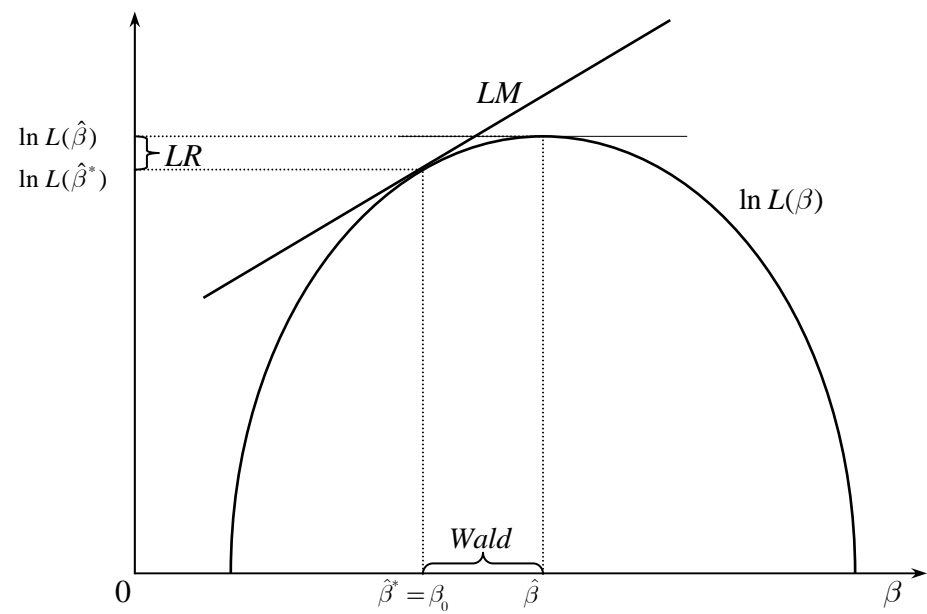


图 11.7 三类渐近等价的统计检验

究竟采取哪种检验常取决于“无约束估计”与“有约束估计”哪种更方便。

如果无约束估计更方便，常使用 **Wald** 检验(比如，对线性回归系数的显著性检验)；

如果有约束估计更方便，常使用 **LM** 检验(比如，对异方差、自相关的检验)；

如果二者都方便，可使用 **LR** 检验(比如，对非线性回归方程的显著性检验)。

11.9 二值选择模型的 Stata 命令与实例

二值模型的 Stata 命令为

`probit y x1 x2 x3,r` (probit 模型)

`logit y x1 x2 x3,r or` (logit 模型)

选择项 “r” 表示使用稳健标准误(默认为普通标准误);

选择项 “or” 表示显示几率比(odds ratio), 不显示回归系数。

完成 Probit 或 Logit 估计后, 可进行预测, 计算准确预测的百分比, 或计算边际效应:

`predict y1` (计算发生概率的预测值, 记为 `y1`)

`estat clas` (计算准确预测的百分比, `clas` 表示 classification)

`margins, dydx(*)` (计算所有解释变量的平均边际效应; “*” 代表所有解释变量)

`margins, dydx(*) atmeans` (计算所有解释变量在样本均值处的边际效应)

`margins, dydx(*) at(x1=0)` (计算所有解释变量在 `x1 = 0` 处的平均边际效应)

`margins, dydx(x1)` (计算解释变量 x_1 的平均边际效应)

`margins, eyex(*)` (计算平均弹性, 其中的两个“e”均指 elasticity)

`margins, eydx(*)` (计算平均半弹性, x 变化一单位引起 y 变化百分之几)

`margins, dyex(*)` (计算平均半弹性, x 变化 1% 引起 y 变化几个单位)

以数据集 `titanic.dta` 为例。

该数据集包括泰坦尼克号乘客的存活数据。

此数据集由 Dawson(1995)提供，原始数据来自英国贸易委员会 (British Board of Trade) 在沉船之后的调查。

该数据集的被解释变量为 `survive`(存活=1，死亡=0)；

解释变量包括 `child`(儿童=1，成年=0)，`female`(女性=1，男性=0)，`class1`(头等舱=1，其他=0)，`class2`(二等舱=1，其他=0)，`class3`(三等舱=1，其他=0)，`class4`(船员=1，其他=0)。

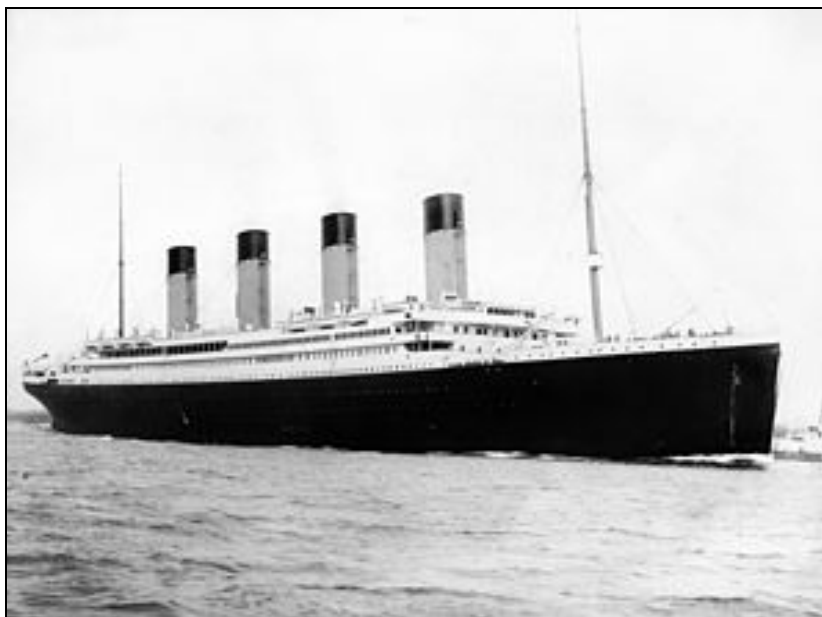


图 11.8 泰坦尼克号于 1914 年 4 月 10 日从英国南安普顿港出发

首先打开数据集，看原始数据。

```
. use titanic.dta, clear  
. list
```

	class1	class2	class3	class4	child	female	survive	freq
1.	0	0	1	0	1	0	0	35
2.	0	0	1	0	1	1	0	17
3.	1	0	0	0	0	0	0	118
4.	0	1	0	0	0	0	0	154
5.	0	0	1	0	0	0	0	387
6.	0	0	0	1	0	0	0	670
7.	1	0	0	0	0	1	0	4
8.	0	1	0	0	0	1	0	13
9.	0	0	1	0	0	1	0	89
10.	0	0	0	1	0	1	0	3
11.	1	0	0	0	1	0	1	5
12.	0	1	0	0	1	0	1	11
13.	0	0	1	0	1	0	1	13
14.	1	0	0	0	1	1	1	1
15.	0	1	0	0	1	1	1	13
16.	0	0	1	0	1	1	1	14
17.	1	0	0	0	0	0	1	57
18.	0	1	0	0	0	0	1	14
19.	0	0	1	0	0	0	1	75
20.	0	0	0	1	0	0	1	192
21.	1	0	0	0	0	1	1	140
22.	0	1	0	0	0	1	1	80
23.	0	0	1	0	0	1	1	76
24.	0	0	0	1	0	1	1	20

原始数据只有 24 个观测值，但每个观测值可能重复多次；其重复次数以最后一列变量 `freq` 表示。

第一行数据显示，乘坐三等舱的男孩死亡者有 35 人；第二行数据显示，乘坐三等舱的女孩死亡者有 17 人；以此类推。

对于观测值重复的数据，在估计时，须以重复次数(`freq`)作为权重才能得到正确结果。

其效果相当于在数据文件中，将第一行数据重复 35 次，第二行数据重复 17 次，以此类推（不同于以方差倒数为权重的 WLS）。

假设观测值的重复次数记录于变量 `freq`，在 Stata 中，可通过在命令的最后加上“`[fweight=freq]`”来实现加权计算或估计；其中“`fweight`”指“frequency weight” (频数权重)。

首先看各变量的统计特征。

```
. sum [fweight=freq]
```

Variable	Obs	Mean	Std. Dev.	Min	Max
survive	2201	.323035	.4677422	0	1
child	2201	.0495229	.2170065	0	1
female	2201	.2135393	.4098983	0	1
class1	2201	.1476602	.3548434	0	1
class2	2201	.1294866	.335814	0	1
class3	2201	.3207633	.466876	0	1

样本容量为 2201(旅客与船员总人数)，而非 24。从变量 `survive` 的平均值可知，平均存活率为 0.32。

分别计算小孩、女士以及各等舱旅客的存活率。

```
. sum survive if child [fweight=freq]
```

Variable	Obs	Mean	Std. Dev.	Min	Max
survive	109	.5229358	.5017807	0	1

```
. sum survive if female [fweight=freq]
```

Variable	Obs	Mean	Std. Dev.	Min	Max
survive	470	.7319149	.4434342	0	1

```
. sum survive if class1 [fweight=freq]
```

Variable	Obs	Mean	Std. Dev.	Min	Max
survive	325	.6246154	.4849687	0	1

```
. sum survive if class2 [fweight=freq]
```

Variable	Obs	Mean	Std. Dev.	Min	Max
survive	285	.4140351	.493421	0	1


```
. sum survive if class3 [fweight=freq]
```

Variable	Obs	Mean	Std. Dev.	Min	Max
survive	706	.2521246	.4345403	0	1

```
. sum survive if class4 [fweight=freq]
```

Variable	Obs	Mean	Std. Dev.	Min	Max
survive	885	.239548	.427049	0	1

小孩、女士、一等舱、二等舱的存活率分别为 0.52、0.73、0.62、0.41，高于平均存活率；三等舱、船员的存活率分别为 0.25、0.24，低于平均存活率。

下面进行回归分析。首先使用 OLS 估计线性概率模型。

```
. reg survive child female class1 class2 class3  
[fweight=freq],r
```

Linear regression				Number of obs = 2201		
				F(5, 2195) = 221.66		
				Prob > F = 0.0000		
				R-squared = 0.2529		
				Root MSE = .40474		
survive	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]	
child	.1812957	.0479499	3.78	0.000	.0872639	.2753275
female	.4906798	.0239292	20.51	0.000	.4437535	.5376061
class1	.1755538	.0291386	6.02	0.000	.1184117	.232696
class2	-.0105263	.0258402	-0.41	0.684	-.0612	.0401475
class3	-.1311806	.0212996	-6.16	0.000	-.17295	-.0894112
_cons	.2267959	.0139872	16.21	0.000	.1993664	.2542254

将虚拟变量 class4(船员)作为参照类别，不放入回归方程。

儿童(child)、妇女(female)与头等舱旅客(class1)的存活概率均显著更高，三等舱旅客(class3)的存活概率显著更低，二等舱旅客(class2)的存活概率与船员无显著差异。

其次，进行 Logit 估计：

```
. logit survive child female class1 class2 class3  
[fweight=freq],nolog
```

选择项 “nolog” 表示不显示 MLE 数值计算的迭代过程。

Logistic regression				Number of obs	=	2201
				LR chi2(5)	=	559.40
				Prob > chi2	=	0.0000
Log likelihood = -1105.0306				Pseudo R2	=	0.2020
survive	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
child	1.061542	.2440257	4.35	0.000	.5832608	1.539824
female	2.42006	.1404101	17.24	0.000	2.144862	2.695259
class1	.8576762	.1573389	5.45	0.000	.5492976	1.166055
class2	-.1604188	.1737865	-0.92	0.356	-.5010342	.1801966
class3	-.9200861	.1485865	-6.19	0.000	-1.21131	-.6288619
_cons	-1.233899	.0804946	-15.33	0.000	-1.391666	-1.076133

Logit 估计结果在显著性方面与 OLS 完全一致。

准 R^2 为 0.20。检验整个方程显著性的 LR 统计量(LR chi2(5))为 559.40, p 值为 0.000, 整个方程高度显著。

使用稳健标准误进行 Logit 估计。

```
. logit survive child female class1 class2 class3
[fweight=freq],nolog r
```

Logistic regression

Number of obs = 2201

Wald chi2(5) = 467.05

Prob > chi2 = 0.0000

Log pseudolikelihood = -1105.0306

Pseudo R2 = 0.2020

survive	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
child	1.061542	.2767452	3.84	0.000	.5191318	1.603953
female	2.42006	.1363096	17.75	0.000	2.152898	2.687222
class1	.8576762	.1475218	5.81	0.000	.5685387	1.146814
class2	-.1604188	.1502193	-1.07	0.286	-.4548432	.1340056
class3	-.9200861	.1621035	-5.68	0.000	-1.237803	-.602369
_cons	-1.233899	.0798876	-15.45	0.000	-1.390476	-1.077322

稳健标准误与普通标准误比较接近。

由于此回归中的解释变量均为虚拟变量，只能变化一个单位(从 0 变为 1)，让 Stata 汇报几率比而非系数。

```
. logit survive child female class1 class2 class3
[fweight=freq],or nolog
```

Logistic regression		Number of obs = 2201				
		LR chi2(5) = 559.40				
		Prob > chi2 = 0.0000				
Log likelihood = -1105.0306		Pseudo R2 = 0.2020				
survive	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
child	2.890826	.7054359	4.35	0.000	1.791872	4.663769
female	11.24654	1.579128	17.24	0.000	8.540859	14.80936
class1	2.357675	.3709541	5.45	0.000	1.732036	3.209306
class2	.851787	.1480291	-0.92	0.356	.6059037	1.197453
class3	.3984847	.0592095	-6.19	0.000	.2978068	.5331983
_cons	.2911551	.0234364	-15.33	0.000	.2486608	.3409114

儿童的生存几率比是成年人的近 3 倍(几率比 2.89), 妇女的存活几率比是男人的 11 倍多(几率比 11.25), 头等舱旅客的存活几率比是船员的 2.36 倍, 三等舱旅客的存活几率比只是船员的 39.8%; 二等舱旅客的存活几率比也略低于船员(几率比 0.85), 但此差别不显著。

计算 Logit 模型的平均边际效应:

```
. margins, dydx( *)
```

Average marginal effects			Number of obs = 2201			
Model VCE : OIM						
Expression : Pr(survive), predict()						
dy/dx w.r.t. : child female class1 class2 class3						
	Delta-method					
	dy/dx	Std. Err.	z	P> z	[95% Conf. Interval]	
child	.1732315	.0393799	4.40	0.000	.0960484	.2504147
female	.394926	.0171966	22.97	0.000	.3612214	.4286307
class1	.1399629	.0250922	5.58	0.000	.0907831	.1891427
class2	-.0261785	.0283616	-0.92	0.356	-.0817663	.0294093
class3	-.1501475	.0238334	-6.30	0.000	-.1968602	-.1034348

Logit 模型的平均边际效应与 OLS 回归系数相差不大。

作为演示，计算在样本均值处的边际效应。

```
. margins, dydx(*) atmeans
```

Conditional marginal effects

Number of obs = 2201

Model VCE : OIM

Expression : Pr(survive), predict()

dy/dx w.r.t. : child female class1 class2 class3

at : child = .0495229 (mean)

female = .2135393 (mean)

class1 = .1476602 (mean)

class2 = .1294866 (mean)

class3 = .3207633 (mean)

	Delta-method					
	dy/dx	Std. Err.	z	P> z	[95% Conf. Interval]	
child	.2223422	.0510772	4.35	0.000	.1222328	.3224516
female	.5068865	.0303542	16.70	0.000	.4473934	.5663797
class1	.179642	.0332374	5.40	0.000	.1144979	.2447861
class2	-.0336	.0363774	-0.92	0.356	-.1048983	.0376983
class3	-.1927139	.0308186	-6.25	0.000	-.2531173	-.1323105

在样本均值处的边际效应与平均边际效应有所不同。

计算 Logit 模型准确预测的比率：

```
. estat clas
```


Logistic model for survive			
Classified	True		Total
	D	~D	
+	349	126	475
-	362	1364	1726
Total	711	1490	2201
Classified + if predicted $\Pr(D) \geq .5$			
True D defined as survive != 0			
Sensitivity	$\Pr(+ D)$		49.09%
Specificity	$\Pr(- \sim D)$		91.54%
Positive predictive value	$\Pr(D +)$		73.47%
Negative predictive value	$\Pr(\sim D -)$		79.03%
False + rate for true ~D	$\Pr(+ \sim D)$		8.46%
False - rate for true D	$\Pr(- D)$		50.91%
False + rate for classified +	$\Pr(\sim D +)$		26.53%
False - rate for classified -	$\Pr(D -)$		20.97%
Correctly classified			77.83%

正确预测的比率为 $(349 + 1364)/2201 = 77.83\%$ 。

根据 Logit 回归结果, 预测每位乘客的存活概率, 记为变量 prob。

```
. predict prob  
(option pr assumed; Pr(survive))
```

考察给定某种特征旅客的生存概率。

计算 Ms. Rose (头等舱、成年、女性)的存活概率：

```
. list prob survive freq if class1==1 & child==0  
& female==1
```

	prob	survive	freq
7.	.8853235	0	4
21.	.8853235	1	140

Ms. Rose 的存活概率高达 88.5%。从频率上看，在所有头等舱的 144 位成年女性中，只有 4 位死亡。

计算 Mr. Jack (三等舱、成年、男性)的存活概率:

```
. list prob survive freq if class3==1 & child==0  
& female==0
```

	prob	survive	freq
5.	.1039594	0	387
19.	.1039594	1	75

Mr. Jack 的存活概率仅有 10.4%。从频率上看, 在所有三等舱的 462 位成年男性中, 只有 75 位生还。

类似地, 可对此数据集进行 Probit 估计。

```
. probit survive child female class1 class2  
class3 [fweight=freq],nolog
```

Probit regression				Number of obs	=	2201
				LR chi2(5)	=	556.83
				Prob > chi2	=	0.0000
Log likelihood = -1106.3142				Pseudo R2	=	0.2011
survive	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
child	.5803382	.1377535	4.21	0.000	.3103463	.85033
female	1.44973	.0808635	17.93	0.000	1.29124	1.608219
class1	.5399101	.0951552	5.67	0.000	.3534092	.7264109
class2	-.0898158	.1028857	-0.87	0.383	-.2914681	.1118364
class3	-.4875252	.0800342	-6.09	0.000	-.6443893	-.3306611
_cons	-.7530486	.0468804	-16.06	0.000	-.8449325	-.6611648

Probit 与 Logit 的回归系数不可比。考察 Probit 模型的平均边际效应及预测准确度。

```
. margins, dydx(*)
```


Probit model for survive			
Classified	True		Total
	D	~D	
+	349	126	475
-	362	1364	1726
Total	711	1490	2201
Classified + if predicted $\Pr(D) \geq .5$			
True D defined as survive != 0			
Sensitivity	$\Pr(+ D)$	49.09%	
Specificity	$\Pr(- \sim D)$	91.54%	
Positive predictive value	$\Pr(D +)$	73.47%	
Negative predictive value	$\Pr(\sim D -)$	79.03%	
False + rate for true ~D	$\Pr(+ \sim D)$	8.46%	
False - rate for true D	$\Pr(- D)$	50.91%	
False + rate for classified +	$\Pr(\sim D +)$	26.53%	
False - rate for classified -	$\Pr(D -)$	20.97%	
Correctly classified		77.83%	

Probit 的平均边际效应、准 R^2 与正确预测比率与 Logit 十分接近，基本等价。

使用 Probit 预测每位个体的存活概率，记为变量 prob1，并考察 prob1 与 prob(Logit 预测结果)的相关性。

```
. predict prob1  
(option pr assumed; Pr(survive))  
  
. corr prob prob1 [fweight=freq]  
(obs=2201)
```

	prob	prob1
prob	1.0000	
prob1	0.9997	1.0000

Probit 与 Logit 对个体存活概率的预测相关系数高达 0.9997，基本无差异。

11.10 其他离散选择模型

(1) 多值选择(multiple choices): 比如, 对交通方式的选择(步行、骑车、自驾车、打的、地铁), 对不同职业的选择, 对手机品牌的选择。

(2) 计数数据(count data): 有时被解释变量只能取非负整数。比如, 企业在某段时间内获得的专利数; 某人在一定时间内去医院看病的次数; 某省在一年内发生煤矿事故的次数。

(3) 排序数据(ordered data): 有些离散数据有着天然的排序。比如, 公司债券的评级(AAA, AA, A, B, C 级), 对“春节联欢晚会”的满意度(很满意、满意、不满意、很不满意)。

对于以上离散数据，一般也不宜直接进行 OLS 回归，主要估计方法仍为 MLE。

由于离散选择模型主要用于微观经济学的实证研究中，故是“微观计量经济学”(Microeconometrics)的重要组成部分。

除了离散数据外，微观计量经济学还关注的另一类数据类型为“受限被解释变量”(limited dependent variable)，即被解释变量的取值范围受到限制(包括断尾回归、归并回归与样本选择模型等)。