

定量研究方法期中考试材料

第三讲 概率与条件概率

**概率公理**

$P(A) \geq 0$   
 $P(\Omega = 1)$

**加法法则**  $P(A \text{或} B) = P(A) + P(B) - P(A \text{和} B)$

**全概率公式**  
 $P(A \text{或} B)$   
 $= P(A \text{和} B) + P(A \text{和} B^c) + P(B \text{和} A^c)P(A)$   
 $= P(A \text{和} B) + P(A \text{和} B^c)$

**条件概率**

$P(A|B) = \frac{P(A \text{和} B)}{P(B)}$  ← 联合概率  
← 边际概率

**乘法法则**  
 $P(A \text{和} B) = P(A|B)P(B) = P(B|A)P(A)$

**全概率公式与条件概率**  $P(A) = P(A|B)P(B) + P(A|B^c)P(B^c)$

**贝叶斯法则**

后验概率 →  $P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|A^c)P(A^c)}$  ← 先验概率

**独立性** 事件A、B独立的充要条件  $P(A \text{和} B) = P(A)P(B)$

**条件独立** 给定事件C, 事件A、B条件独立的充要条件  $P(A \text{和} B|C) = P(A|C)P(B|C)$  即  $P(A|B \text{和} C) = P(A|C)$

第四讲 随机变量与概率分布

**随机变量的种类**

离散随机变量: 数值为有限个, 如家庭成员个数、态度 (支持/反对) 等。  
→ 概率质量函数 (PMF)、累积分布函数 (CDF)

连续随机变量: 在某个实线区间内有无限个取值, 如长度、GDP等。  
→ 概率密度函数 (PDF)、累积分布函数 (CDF)

**分布**

**伯努利随机变量:** 一个二元随机变量, 该二元随机变量的数值为两个不同的数。

PMF:  $f(x) = \begin{cases} p & (x = 1) \\ 1 - p & (x = 0) \\ 0 & \text{其他} \end{cases}$  CDF:  $F(x) = \begin{cases} 0 & (x < 0) \\ 1 - p & (0 \leq x < 1) \\ 1 & (x \geq 1) \end{cases}$

**伯努利分布**

**均匀分布** 均匀随机变量: 均匀随机变量在给定区间[a, b]内取一值的可能性相同。

PDF: 在[a, b]内  $f(x) = \begin{cases} \frac{1}{b-a} & (a \leq x \leq b) \\ 0 & \text{其他} \end{cases}$  CDF: 在[a, b]内  $F(x) = \begin{cases} 0 & (x < a) \\ \frac{x-a}{b-a} & (a \leq x \leq b) \\ 1 & (x \geq b) \end{cases}$

$P(\mu - k\sigma \leq X \leq \mu + k\sigma) = P\left(-k \leq \frac{X - \mu}{\sigma} \leq k\right) = P(-k \leq Z \leq k)$

第五讲 点估计

**估计相关概念**

估计目标 (estimand) 是我们需要估计的总体的参数 (parameter)。通常用希腊字母表示 (如 $\mu$ 、 $\theta$ )。

估计量 (estimator) 是我们估算估计目标的方法。通常用希腊字母加hat表示 (如 $\hat{\mu}$ 、 $\hat{\theta}$ )。

估计值 (estimate) 是我们运用估计量从所得样本中进行估算所得的具体数值。

**期望和方差**

**期望**  
 $E(X) = \begin{cases} \sum_x x \cdot f(x) & \text{如果} X \text{是离散的} \\ \int x \cdot f(x)dx & \text{如果} X \text{是连续的} \end{cases}$  其中,  $f(x)$ 是离散变量X的概率质量函数 (PMF) 或连续变量X的概率密度函数 (PDF)

$E(a) = a$   $E[E(X)] = E(X)$   
 $E(aX) = aE(X)$  X与Y独立的必要非充分条件:  
 $E(aX + b) = aE(X) + b$   $E(XY) = E(X)E(Y)$   
 $E(aX + bY) = aE(X) + bE(Y)$

**方差**  $V(X) = E\{[X - E(X)]^2\} = E(X^2) - \{E(X)\}^2$   
 $V(a) = 0$   $V(aX) = a^2V(X)$  X与Y独立的必要非充分条件:  
 $V(X + b) = V(X)$   
 $V(aX + b) = a^2V(X)$   $V(X + Y) = V(X) + V(Y)$

**协方差** X与Y独立的必要非充分条件:  
 $Cov(X, Y) = E\{[X - E(X)][Y - E(Y)]\} = E(XY) - E(X)E(Y)$   $Cov(X, Y) = 0$

**适用于任何样本容量的有限样本特征**

**无偏性** 估计量无偏的充要条件:  
 $Bias(\hat{\theta}) = E(\hat{\theta} - \theta) = E(\hat{\theta}) - \theta$   $Bias(\hat{\theta}) = 0 \Leftrightarrow E(\hat{\theta}) = \theta$

**相对有效性** 比较 $\hat{\theta}_1$ 和 $\hat{\theta}_2$ , 如果的 $\hat{\theta}_1$ 方差小于 $\hat{\theta}_2$ 的方差, 那么 $\hat{\theta}_1$ 更有效  $V(\hat{\theta}_1) < V(\hat{\theta}_2)$

**平均数平方误差**  $MSE(\hat{\theta}) = E\{(\hat{\theta} - \theta)^2\} = V(\hat{\theta}) + Bias(\hat{\theta})^2$  (均方误, MSE)  
用于估计量都存在偏误时, 选取合适的估计方案。

**适用于大样本容量的特征**

**一致性** 如果一个估计量 $\hat{\theta}_n$ 的抽样分布 $\hat{\theta}_1, \dots, \hat{\theta}_n$ 随着样本容量n的增加越来越集中于估计目标 $\theta$ , 那么 $\hat{\theta}_n$ 就是一致的。  
 $\hat{\theta}_n \xrightarrow{p} \theta$  即  $p - \lim_{n \rightarrow \infty} \hat{\theta}_n = 0$

一致性体现在估计量的期望值越来越接近估计目标以及估计量的方差逐渐趋近于0。一致的估计量不一定无偏, 无偏误的估计量不一定具有一致性。

**渐进正态性**

**中心极限定理:** 样本平均数的分布随着样本数的增加接近正态分布。  
假设我们有独立并且相同分布 (i.i.d.) 的随机变量的样本 $X_1, X_2, \dots, X_n$ , 形成一个均值为 $\mu$ , 方差为 $\sigma^2$ 的概率分布。样本平均数用 $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ 表示, 那么中心极限定理为:  
 $Z_n = \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} \rightarrow N(0, 1)$

其中,  $\mu$ 可以用 $E(X)$ 来表示,  $\sigma^2$ 可以用 $V(X)$ 来表示。  
这一公式的含义是: 随着样本规模的增加, 分布样本平均数的z分数收敛成标准正态分布, 即 $N(0, 1)$ 。如果样本容量n足够大, 我们可以用样本标准差 $S_n$ 来代替 $\sigma$ 。

这一公式的含义是: 随着样本规模的增加, 样本平均数的分布收敛成正态分布, 即 $N(\mu, \frac{\sigma^2}{n})$

在样本容量足够大的情况下:  $\bar{X}_n \rightarrow N\left(\mu, \frac{\sigma^2}{n}\right)$

**第六讲 区间估计**

**标准误** 用于刻画随机抽样所得的样本平均值的离散程度。  
假设我们有独立且相同分布的随机变量样本 $X_1, X_2, \dots, X_n$ , 并且已知总体方差 $\sigma^2$ 的时候, 样本平均数的标准误为:  $\frac{\sigma}{\sqrt{n}}$   
可以通过样本方差来估计总体方差 $\sigma^2$   
 $S_n^2 = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{n - 1}$

**置信水平** 决定了在多大程度上确信区间内包括估计目标的真实值。研究者通常选取95%、90%和99%作为置信水平 (其中95%置信水平最为常见)。  
(置信度) 置信水平通常被写为 $(1-\alpha) \times 100\%$ 。

**$\sigma^2$ 已知的置信区间**

假设我们有独立并且相同分布的随机变量的样本 $X_1, X_2, \dots, X_n$ , 并且来自于总体平均值为 $\mu$ , 方差为 $\sigma^2$ 的正态分布中, 总体方差已知的前提下, 对于样本平均数 $\bar{X}$ 的 $(1-\alpha) \times 100\%$ 置信区间为:  
 $CI(\alpha) = \left[ \bar{X} - z_{\frac{\alpha}{2}} \times \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\frac{\alpha}{2}} \times \frac{\sigma}{\sqrt{n}} \right]$  → 临界值 $z_{\frac{\alpha}{2}}$ 等于标准正态分布 $N(0, 1)$ 的 $(1 - \frac{\alpha}{2})$ 分位数。最大临界值 $z_{\frac{\alpha}{2}}$ 和最小临界值 $-z_{\frac{\alpha}{2}}$ 围成的面积为置信水平 $(1-\alpha) \times 100\%$ 。置信水平越高, 最大临界值越高, 在其他条件不变的情况下, 所得置信区间的范围越大。

其中 $\frac{\sigma}{\sqrt{n}}$ 为标准误,  $z_{\frac{\alpha}{2}}$ 为临界值

n%置信区间的含义是: 在重复数据产生的过程中, 有n%的数据的区间内包含估计目标。

在总体方差未知的情况下，对于样本平均数 $\bar{x}$ 的 $(1-\alpha) \times 100\%$ 的大样本置信区间为：

$$CI(\alpha) = \left[ \bar{x} - t_{\frac{\alpha}{2}, n-1} \times \frac{S_n}{\sqrt{n}}, \bar{x} + t_{\frac{\alpha}{2}, n-1} \times \frac{S_n}{\sqrt{n}} \right]$$

其中 $\frac{S_n}{\sqrt{n}}$ 为标准误， $t_{\frac{\alpha}{2}, n-1}$ 为临界值  
在相同置信水平的情况下，t分布的临界值要稍大于正态分布的临界值

---

随着自由度（df）的增加，t分布接近正态分布。  
随着 $n \rightarrow \infty$ ，则（见右式）：

$$T_n = \frac{\bar{X}_n - \mu}{\frac{S_n}{\sqrt{n}}} \rightarrow N(0, 1)$$

这意味着当样本容量足够大的时候我们不需要假设总体服从正态分布。

**p值检验** p值指的是在零假设正确的情况下，至少观测到一次检验统计量的概率；p值越小，拒绝零假设的证据越强

单边： $p\text{值} = P(Z_n \geq Z_n | H_0)$

双边： $p\text{值} = P(Z_n \geq Z_n \text{ 或 } Z_n \leq Z_n | H_0)$

**检验统计量**  $Z_n$   $Z_n = \frac{\bar{X}_n - \mu_0}{\frac{\hat{\sigma}}{\sqrt{n}}}$

样本容量足够大时,  
可以用样本方差  $S_n^2$  来估计  $\sigma$   
且  $Z_n \sim N(0,1)$

\*临界值  $z$  会在卷面上给出。

$$Z_n = \frac{\bar{X}_a - \bar{X}_b}{\sqrt{\frac{\hat{\sigma}_a^2}{n_a} + \frac{\hat{\sigma}_b^2}{n_b}}}$$

**样本协方差**

$$Cov(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

当Cov(X, Y)为正值时，整体数据云的趋势向上；为负值时，则向下。

---

**相关系数**

$$Cor(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}} = \frac{Cov(X, Y)}{S_x S_y}$$

其中， $S_x$ 、 $S_y$ 分别是变量X和Y的标准差

$$\hat{\beta} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

根据协方差和样本方差的公式有：

$$\hat{\beta} = \frac{Cov(X, Y)}{S_X^2}$$

根据相关系数的公式有：

$$\hat{\beta} = \rho_{XY} \times \frac{S_Y}{S_X}$$

其中， $\rho_{XY}$ 是X与Y的相关系数

在模型中，解释平方和是可被回归模型解释的部分，残差平方和则是没有被回归模型预测到的部分。  

$$TSS = ESS + SSR$$

$$S_y^2 = S_y^2 + S_e^2$$

---

**决定系数  $R^2$**  为预测值  $\hat{X}$  可以解释的总平方和 (TSS) 的比例，表明线性模型与数据的拟合程度

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{SSR}{TSS}$$

**性质：**

- $0 \leq R^2 \leq 1$ ;
- 当  $R^2 = 1$  时，所有的数据点都在回归线上；
- 当  $R^2 = 0$  时，X 和 Y 之间没有相关性。

应注意， $R^2$  越大，数据的拟合程度不一定越好

[illegible]