

# HomeWork

邹翔宇 | 2410833001

## 第一题

### 1. 点估计和区间估计的区别及应用场景

区别:

- 点估计是用单个数值(如样本均值、样本比例)直接估计总体参数(如总体均值、总体比例)。
- 区间估计是给出一个区间范围及对应的置信度(如 95% 置信区间),表示该区间以一定概率包含总体参数。

应用场景:

- 点估计:例如,工厂质检中估计某批次零件的平均长度为 10.2 毫米。
- 区间估计:例如,在药物临床试验中,估计某药的平均疗效提升为 5% 15% (置信度 95%),既给出范围又体现可信度。

### 2. 标准误的定义及正态分布、泊松分布的均数标准误

定义:标准误(Standard Error, SE)是统计量(如样本均值)的标准差,反映抽样误差的大小。公式:

- 正态分布:均数的标准误为  $SE = \frac{\sigma}{\sqrt{n}}$ ,其中 $\sigma$ 是总体标准差, $n$ 为样本量。
- 泊松分布:均数的标准误为  $SE = \sqrt{\frac{\lambda}{n}}$ ,其中 $\lambda$ 是总体均值(泊松分布的均值等于方差)。

### 3. 标准误的影响因素及样本量增加 4 倍时的变化

影响因素:

- 总体标准差( $\sigma$ 或 $\sqrt{\lambda}$ ):总体变异越大,标准误越大。
- 样本量( $n$ ):样本量越大,标准误越小。

样本量增加 4 倍时的变化:对于率的标准误(如二项分布),公式为  $SE = \sqrt{\frac{p(1-p)}{n}}$ 。当样本量 $n$ 增加 4 倍,标准误变为原来的 $\frac{1}{2}$ (因 $\sqrt{\frac{1}{4}} = \frac{1}{2}$ )。

### 4. 区间估计中使用 t 分布的情况

当同时满足以下条件时需使用 t 分布:

- 总体标准差未知,需用样本标准差 $s$ 代替。
- 样本量较小(通常 $n < 30$ )或总体不严格服从正态分布。

示例:研究 10 名学生的平均体重时,若总体方差未知,应使用 t 分布计算置信区间(而非正态分布)。

## 第二题

### 1. 总体均数的点估计值

点估计是用样本统计量直接估计总体参数。此处样本均值为 5.2 mmol/L,而样本均数是总体均数的无偏估计量。因此,总体均数的点估计值为 5.2 mmol/L。

### 2. 已知总体标准差( $\sigma=0.8$ )时,构建总体均数的 95% 置信区间

当总体标准差 $\sigma$ 已知时,使用正态分布(Z 分布)构建置信区间:

$$\text{置信区间} = \bar{X} \pm Z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$$

其中:

- $\bar{X} = 5.2$
- $\sigma = 0.8$
- $n = 25$
- $Z_{0.025} = 1.96$  (对应 95% 置信水平)

代入计算:

$$5.2 \pm 1.96 \cdot \frac{0.8}{\sqrt{25}} = 5.2 \pm 1.96 \cdot 0.16 = 5.2 \pm 0.3136$$

区间范围: 4.8864 ~ 5.5136 mmol/L

### 3. 已知样本标准差( $s=0.8$ )时,构建总体均数的 95% 置信区间

当总体标准差未知且样本量较小( $n = 25 < 30$ )时,需使用 t 分布:

$$\text{置信区间} = \bar{X} \pm t_{\frac{\alpha}{2}, \text{df}} \cdot \frac{s}{\sqrt{n}}$$

其中:

- $\bar{X} = 5.2$
- $s = 0.8$
- $n = 25$
- 自由度  $\text{df} = n - 1 = 24$
- $t_{0.025, 24} \approx 2.064$  (查 t 分布表)

代入计算:

$$5.2 \pm 2.064 \cdot \frac{0.8}{\sqrt{25}} = 5.2 \pm 2.064 \cdot 0.16 = 5.2 \pm 0.3302$$

区间范围: 4.8698 ~ 5.5302 mmol/L

依据: 小样本且总体方差未知时需用 t 分布。

#### 4. 置信区间比较及 t 分布与正态分布的区别与联系

##### (1) 比较两个置信区间大小

- 正态分布置信区间: 4.8864 ~ 5.5136 (宽度约 0.6272)
- t 分布置信区间: 4.8698 ~ 5.5302 (宽度约 0.6604)

结论: t 分布的置信区间更宽, 说明在小样本且总体方差未知时, t 分布考虑了额外的不确定性。

##### (2) t 分布与正态分布的区别与联系 区别:

1. 尾部厚度: t 分布尾部更厚, 对极端值更敏感, 适用于小样本。
2. 自由度影响: t 分布形状由自由度决定, 自由度越大越接近正态分布; 而正态分布形态固定。
3. 应用条件:
  - 正态分布: 总体方差已知或大样本 ( $n \geq 30$ )。
  - t 分布: 总体方差未知且小样本。

联系:

1. 对称性: 两者均为对称分布, 以均值为中心。
2. 趋近性: 当自由度  $\text{df} \rightarrow \infty$  时, t 分布收敛于标准正态分布。
3. 参数估计: 两者均可用于总体均数的区间估计, 但适用条件不同。

总结: t 分布是小样本统计推断的核心工具, 其灵活性弥补了正态分布在小样本场景中的不足。

### 第三题

#### 1. 高血压患病率的点估计值

$$\hat{p} = \frac{48}{400} = 0.12$$

解释: 点估计值为样本患病率, 即  $\frac{48}{400} = 0.12$ 。

#### 2. 正态近似条件验证

条件:  $n\hat{p} \geq 5$  且  $n(1 - \hat{p}) \geq 5$ 。 计算:

$$n\hat{p} = 400 \times 0.12 = 48 \geq 5, \quad n(1 - \hat{p}) = 400 \times 0.88 = 352 \geq 5$$

结论: 正态近似条件成立。

#### 3. 患病率的 95% 置信区间

公式:

$$\hat{p} \pm Z_{\frac{\alpha}{2}} \cdot \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

计算步骤:

- 标准误 (SE):

$$\sqrt{\frac{0.12 \times 0.88}{400}} = \sqrt{0.000264} \approx 0.01625$$

- 95% 置信水平的 Z 值:  $Z_{0.975} = 1.96$
- 误差范围:

$$1.96 \times 0.01625 \approx 0.0318$$

- 置信区间:

$$0.12 \pm 0.0318 \Rightarrow [0.088, 0.152]$$

答案: 95% 置信区间为 (0.088, 0.152)。

#### 4. 样本量增至 1600 人时置信区间的宽度变化

结论: 置信区间的宽度会 缩小为原来的一半。

原因:

- 标准误公式为  $\sqrt{\frac{p(1-p)}{n}}$ , 当样本量  $n$  增至 4 倍时, 标准误变为原来的  $\frac{1}{\sqrt{4}} = \frac{1}{2}$ 。
- 置信区间宽度为  $2 \times Z_{\frac{\alpha}{2}} \times SE$ , 因此宽度也会缩小为原来的 1/2。

### 第四题

#### 1. 总体均数之差的点估计

$$\hat{\mu}_1 - \hat{\mu}_2 = 120 - 115 = 5$$

#### 2. 两样本标准差均为 15 时的 95% 置信区间

条件: 假设总体方差相等, 使用合并方差计算。公式:

$$(\bar{x}_1 - \bar{x}_2) \pm t_{\frac{\alpha}{2}, \text{df}} \cdot s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

计算步骤:

##### 1. 合并方差:

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} = \frac{14 \times 225 + 9 \times 225}{23} = 225 \Rightarrow s_p = 15$$

##### 2. 标准误 (SE):

$$s_p \sqrt{\frac{1}{15} + \frac{1}{10}} = 15 \times \sqrt{\frac{1}{6}} \approx 6.124$$

3. 自由度 (df):  $n_1 + n_2 - 2 = 23$ , 查 t 表得  $t_{0.975, 23} = 2.069$

##### 4. 误差范围:

$$2.069 \times 6.124 \approx 12.67$$

##### 5. 置信区间:

$$5 \pm 12.67 \Rightarrow [-7.67, 17.67]$$

答案: 95% 置信区间为 (-7.67, 17.67)。

#### 3. 样本方差不同 ( $s_1=15$ , $s_2=12$ ) 时的 95% 置信区间

条件: 总体方差不等, 使用 Welch-Satterthwaite 方法。公式:

$$(\bar{x}_1 - \bar{x}_2) \pm t_{\frac{\alpha}{2}, \text{df}} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

计算步骤:

##### 1. 标准误 (SE):

$$\sqrt{\frac{15}{15} + \frac{12}{10}} = \sqrt{1 + 1.2} \approx 1.483$$

2. 自由度近似 (Welch 公式):

$$df \approx \frac{\left(\frac{15^2}{15} + \frac{12^2}{10}\right)^2}{\frac{\left(\frac{15^2}{15}\right)^2}{14} + \frac{\left(\frac{12^2}{10}\right)^2}{9}} \approx 22$$

3. 查 t 表得  $t_{0.975,22} = 2.074$

4. 置信区间:

$$5 \pm 2.074 \times \sqrt{\frac{115^2}{15} + \frac{12^2}{12}} \Rightarrow [-6.24, 16.24]$$

答案: 95% 置信区间为  $[-6.24, 16.24]$ 。

#### 4. 总体方差已知 ( $\sigma_1=15, \sigma_2=12$ ) 时的 95% 置信区间

条件: 总体方差已知, 使用正态分布 (Z 分布)。公式:

$$(\bar{x}_1 - \bar{x}_2) \pm Z_{\frac{\alpha}{2}} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

计算步骤:

1. 标准误 (SE):

$$\sqrt{\frac{15}{15} + \frac{12}{10}} = \sqrt{1 + 1.2} \approx 1.483$$

2. 95% 置信水平的 Z 值:  $Z_{0.975} = 1.96$

3. 置信区间:

$$5 \pm 1.96 \times \sqrt{\frac{115^2}{15} + \frac{12^2}{12}} \Rightarrow [-5.6273, 15.6273]$$

答案: 95% 置信区间为  $(-5.6273, 15.6273)$ 。

### 第五题

#### 1. 疫苗 A 和疫苗 B 未感染率的点估计值

· 疫苗 A:  $\hat{p}_A = \frac{240}{300} = 0.80$

· 疫苗 B:  $\hat{p}_B = \frac{180}{300} = 0.60$

疫苗 A 未感染率点估计值为 **0.80**, 疫苗 B 为 **0.60**。

#### 2. 两独立样本率之差的标准误公式及计算

公式:

$$SE(\hat{p}_A - \hat{p}_B) = \sqrt{\frac{\hat{p}_A(1-\hat{p}_A)}{n_A} + \frac{\hat{p}_B(1-\hat{p}_B)}{n_B}}$$

· 疫苗 A:  $\frac{0.80 \times 0.20}{300} = 0.000533$

· 疫苗 B:  $\frac{0.60 \times 0.40}{300} = 0.000800$

· 标准误:  $\sqrt{0.000533 + 0.000800} = \sqrt{0.001333} \approx 0.0365$

#### 3. 两疫苗未感染率之差的 95% 置信区间及结论

公式:

$$(\hat{p}_A - \hat{p}_B) \pm Z_{\frac{\alpha}{2}} \cdot SE$$

计算过程:

1. 率之差:  $0.80 - 0.60 = 0.20$

2. 标准误: 0.0365

3. 95% 置信水平的 Z 值:  $Z_{0.975} = 1.96$

4. 误差范围:  $1.96 \times 0.0365 \approx 0.0715$

5. 置信区间:

$$0.20 \pm 0.0715 \Rightarrow [0.1285, 0.2715]$$

结论：由于置信区间  $[0.13, 0.27]$  不包含 0，可以认为两种疫苗的有效性存在 **显著差异** ( $p < 0.05$ )。

#### 4. 样本量扩大至每组 1000 人时置信区间的宽度变化

原因：

- 标准误公式中，样本量  $n$  位于分母，且与标准误成反比。
- 当样本量从 300 增至 1000（扩大约 3.33 倍），标准误变为原来的  $\frac{1}{\sqrt{3.33}} \approx 0.55$  倍。
- 置信区间宽度与标准误成正比，因此宽度会 **缩小至原来的约 55%**。

置信区间宽度会 **显著变窄**，具体缩小至原宽度的约 55%。