

Inferência Estatística

Docente: Manuel G. Scotto

Departamento de Matemática
IST, ULisboa

População e amostra

- **População:** conjunto de indivíduos ou objetos que apresentam uma ou mais características em comum, podendo ser finita ou infinita, dependendo do número de elementos que a compõem.

População e amostra

- **População:** conjunto de indivíduos ou objetos que apresentam uma ou mais características em comum, podendo ser finita ou infinita, dependendo do número de elementos que a compõem.
- O objetivo usual é **inferir** sobre os parâmetros (ou a forma) da distribuição de X .

População e amostra

- **População:** conjunto de indivíduos ou objetos que apresentam uma ou mais características em comum, podendo ser finita ou infinita, dependendo do número de elementos que a compõem.
- O objetivo usual é **inferir** sobre os parâmetros (ou a forma) da distribuição de X .
- **Amostra:** subconjunto de uma população da qual são estudadas as características. Tem que ser **representativa** relativamente à população de onde foi retirada.

Amostra representativa



Amostra representativa



Amostras NÃO representativas

- **Amostras auto selecionadas:** apresentar uma questão e solicitar os espectadores que telefonem para um número se a sua opinião é “sim” e para outro número se a sua opinião é “não”.

Amostras NÃO representativas

- **Amostras auto selecionadas:** apresentar uma questão e solicitar os espectadores que telefonem para um número se a sua opinião é “sim” e para outro número se a sua opinião é “não”.

- 1 Deve haver pagamento de propinas no ensino superior público?

Amostras NÃO representativas

- **Amostras auto selecionadas:** apresentar uma questão e solicitar os espectadores que telefonem para um número se a sua opinião é “sim” e para outro número se a sua opinião é “não”.
- 1 Deve haver pagamento de propinas no ensino superior público?
 - 2 Os hipermercados devem fechar aos domingos?

Amostras NÃO representativas

- **Amostras auto selecionadas:** apresentar uma questão e solicitar os espectadores que telefonem para um número se a sua opinião é “sim” e para outro número se a sua opinião é “não”.

- 1 Deve haver pagamento de propinas no ensino superior público?
- 2 Os hipermercados devem fechar aos domingos?
- 3 Há hipótese de Portugal voltar a ter monarquia?

Amostras NÃO representativas

- **Amostras inadequadas:** em geral não acontece inadvertidamente. É intencional a generalização abusiva do que foi observado numa amostra muito diminuta para uma população de dimensão considerável.





Amostra aleatória

- Seja X uma variável aleatória de interesse e X_1, \dots, X_n variáveis aleatórias independentes e identicamente distribuídas (i.i.d.) de X . Então, o vetor aleatório $\underline{X} = (X_1, \dots, X_n)$ diz-se uma amostra aleatória (a.a.) respeitante à variável aleatória X .
- À observação particular da a.a. dá-se o nome de amostra e representa-se por $\underline{x} = (x_1, \dots, x_n)$.

Caraterização da amostra aleatória

Caso discreto: Função de probabilidade conjunta de \underline{X}

$$\begin{aligned}P(\underline{X} = \underline{x}) &= P(X_1 = x_1, \dots, X_n = x_n) \\&= \prod_{i=1}^n P(X_i = x_i) \quad X_i \text{ ind} \\&= \prod_{i=1}^n P(X = x_i) \quad X_i \sim X\end{aligned}$$

Caraterização da amostra aleatória

Caso contínuo: Função de densidade conjunta de \underline{X}

$$\begin{aligned} f_{\underline{X}}(\underline{x}) &= f_{X_1, \dots, X_n}(x_1, \dots, x_n) \\ &= \prod_{i=1}^n f_{X_i}(x_i) \quad X_i \text{ ind} \\ &= \prod_{i=1}^n f_X(x_i) \quad X_i \sim X \end{aligned}$$

Estatística

Função da amostra, $T = T(\underline{X})$, que não depende de parâmetros desconhecidos.

Estatística

Função da amostra, $T = T(\underline{X})$, que não depende de parâmetros desconhecidos.

Estatística		Valor observado da estatística	
Mínimo da a.a.	$X_{(1)} = \min_{i=1,\dots,n} X_i$	mínimo da amostra	$x_{(1)} = \min_{i=1,\dots,n} x_i$
Máximo da a.a.	$X_{(n)} = \max_{i=1,\dots,n} X_i$	máximo da amostra	$x_{(n)} = \max_{i=1,\dots,n} x_i$
Amplitude da a.a.	$R = X_{(n)} - X_{(1)}$	amplitude da amostra	$r = x_{(n)} - x_{(1)}$
Média da a.a.	$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$	média da amostra	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
Var. corrigida da a.a.	$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$	var. corrigida da am.	$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
Var. não corrig. da a.a.	$(S')^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$	var. não corrig. da am.	$(s')^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

Estimação paramétrica

- Seja (X_1, \dots, X_n) uma amostra de uma população cuja função densidade pertence à família $\mathcal{F}_\theta = \{f(x|\theta) : \theta \in \Theta\}$, de que se desconhece apenas o verdadeiro valor do parâmetro.
- A estimação paramétrica procura responder à seguinte questão: como utilizar a informação dada pela amostra para estimar o valor desconhecido de θ ?
- Para resolver esta questão torna-se necessário ter em consideração dois aspetos fundamentais: a **precisão** e a **confiança**.

Estimador e estimativa pontual

- **Estimador:** a estatística $T = T(\underline{X})$ diz-se um estimador do parâmetro desconhecido θ , caso tome valores exclusivamente no espaço paramétrico Θ .
- **Estimativa:** é o valor concreto assumido pelo estimador para uma amostra particular $\underline{x} = (x_1, \dots, x_n)$, $t = T(\underline{x})$.

Exemplo

Seja $X \sim Be(\theta)$, com $\theta \in \Theta = [0, 1]$. É $T(X) = \frac{1}{n} \sum_{i=1}^n X_i$ um estimador para θ ?

- $T(X)$ só depende de X ;
- $T(X)$ toma valores em $\{0, \frac{1}{n}, \frac{2}{n}, \dots, 1\} \in [0, 1]$.

Importante

Para introduzir o **método da máxima verosimilhança** é indispensável definir primeiro o conceito de função de verosimilhança.

Função de densidade conjunta

Se X_1, \dots, X_n é uma a.a. casual de uma população com função densidade (ou função de probabilidade) dada por $f(x|\theta)$, a expressão

$$f(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta), \quad (x_1, x_2, \dots, x_n) \in \mathbb{R}^n,$$

define a função densidade conjunta das variáveis que constituem a amostra, isto é, designa para um dado $\theta \in \Theta$ a densidade associada com a mostra particular (x_1, \dots, x_n) .

Função de verosimilhança

Fixado (x_1, \dots, x_n) isto é observada uma amostra concreta, a mesma expressão interpretada como função do parâmetro θ define a função de verosimilhança, L . Pode escrever-se

$$L(\theta|x_1, x_2, \dots, x_n) \equiv L(\theta) = \prod_{i=1}^n f(x_i|\theta), \quad (\theta \in \Theta).$$

Método da máxima verosimilhança

Dada a amostra (x_1, \dots, x_n) , o método da máxima verosimilhança consiste em procurar se existe uma estimativa

$$\hat{\theta} \equiv \hat{\theta}(x_1, \dots, x_n) = \hat{\theta}(\underline{x}),$$

tal que

$$L(\hat{\theta}|x_1, x_2, \dots, x_n) \geq L(\theta|x_1, x_2, \dots, x_n), \quad \forall \theta \in \Theta.$$

Estimativa de máxima verosimilhança

Obtida a amostra (x_1, \dots, x_n) , a estimativa de máxima verosimilhança do parâmetro desconhecido corresponde ao **ponto de máximo** da função de verosimilhança ou, equivalentemente, ao ponto de máximo do logaritmo da função de verosimilhança, isto é

$$L(\hat{\theta}|\underline{x}) = \max_{\theta \in \Theta} L(\theta|\underline{x}), \text{ ou } \ln L(\hat{\theta}|\underline{x}) = \max_{\theta \in \Theta} \ln L(\theta|\underline{x}).$$

Estimador de máxima verosimilhança

O estimador de MV de θ obtém-se por substituição de \underline{x} por \underline{X} na expressão geral da estimativa de MV $\hat{\theta}(\underline{x})$.

Exercícios

- Utilize o método da máxima verosimilhança para encontrar estimadores para os parâmetros, no caso das seguintes distribuições: (i) $Po(\lambda)$, (ii) $N(\mu, \sigma^2)$ com σ^2 conhecida, (iii) $N(\mu, \sigma^2)$ com μ conhecido.
- Seja X_1, \dots, X_n uma amostra casual de uma população com função densidade dada por $f(x|\theta) = \theta x^{\theta-1}$, ($0 < x < 1$), com $\theta > 0$. Calcule um estimador para θ utilizando o método da máxima verosimilhança.

Resolução

Utilize o método da máxima verosimilhança para encontrar estimadores para os parâmetros, no caso das seguintes distribuições:

(i) $Po(\lambda)$, (ii) $N(\mu, \sigma^2)$ com σ^2 conhecida, (iii) $N(\mu, \sigma^2)$ com μ conhecido.

Caso I: Função de verosimilhança,

$$L(\theta|x_1, x_2, \dots, x_n) \equiv L(\theta) = \prod_{i=1}^n f(x_i|\theta), \quad (\theta \in \Theta).$$

Neste caso $\theta = \{\lambda\}$.

Resolução

$$L(\lambda) = \prod_{i=1}^n e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} = e^{-n\lambda} \frac{\lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!}.$$

Função de log-verosimilhança:

$$\ln L(\lambda) = -n\lambda + \left(\sum_{i=1}^n x_i \right) \ln \lambda - \ln \left(\prod_{i=1}^n x_i! \right).$$

Resolução

Maximização:

$$\frac{\partial \ln L(\lambda)}{\partial \lambda} = -n + \frac{\sum_{i=1}^n x_i}{\lambda} = 0 \Rightarrow \frac{\sum_{i=1}^n x_i}{\lambda} = n,$$

pelo que a estimativa de máxima verosimilhança de λ é $\hat{\lambda} = \bar{x}$.

O estimador de máxima verosimilhança de λ é $\hat{\lambda}_{MV} = \bar{X}$.

Resolução

Caso II: neste caso $\theta = \{\mu\}$.

$$L(\mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n e^{-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}}.$$

Função de log-verosimilhança:

$$\ln L(\mu) = \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n - \frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}.$$

Resolução

Maximização:

$$\frac{\partial \ln L(\mu)}{\partial \mu} = \frac{-2 \sum_{i=1}^n (x_i - \mu)(-1)}{2\sigma^2} = \frac{\sum_{i=1}^n (x_i - \mu)}{\sigma^2} = 0$$

o que implica que

$$\sum_{i=1}^n (x_i - \mu) = 0 \Leftrightarrow \sum_{i=1}^n x_i - n\mu = 0 \Leftrightarrow \sum_{i=1}^n x_i = n\mu,$$

pelo que a estimativa de máxima verosimilhança de μ é $\hat{\mu} = \bar{x}$.

O estimador de máxima verosimilhança de μ é $\hat{\mu}_{MV} = \bar{X}$.

Exercícios

- Considere uma amostra aleatória de uma distribuição Uniforme no intervalo $(0, \theta)$. Calcule um estimador para θ utilizando o método da máxima verosimilhança.
- Seja X_1, \dots, X_n uma amostra casual de uma população com função densidade dada por

$$f(x|\theta) = \begin{cases} 1 & \theta - 1/2 \leq x \leq \theta + 1/2 \\ 0 & \text{outros casos} \end{cases},$$

para $\theta \in (-\infty, \infty)$. Calcule um estimador para θ utilizando o método da máxima verosimilhança.

Resolução

Considere uma amostra aleatória de uma distribuição Uniforme no intervalo $(0, \theta)$. Calcule um estimador para θ utilizando o método da máxima verosimilhança.

$$f(x; \theta) = \begin{cases} \frac{1}{\theta} & 0 < x < \theta \\ 0 & \text{outros casos} \end{cases}.$$

Função de distribuição conjunta:

$$f(x_1, \dots, x_n) = \frac{1}{\theta^n} \prod_{i=1}^n I_{(0, \theta)}(x_i) = \frac{1}{\theta^n} I_{(0, \theta)}(x_{(n)}) I_{(0, x_{(n)})}(x_{(1)})$$

com $x_{(1)} = \min(x_1, \dots, x_n)$ e $x_{(n)} = \max(x_1, \dots, x_n)$.

Resolução

Estimador de máxima verosimilhança para θ

$$\hat{\theta}_{MV} = X_{(n)},$$

com $X_{(n)} = \max(X_1, \dots, X_n)$.

Propriedades dos estimadores de MV

- **Invariância:** Seja h uma função bijetiva de θ . Então

$$\widehat{h(\theta)} = h(\hat{\theta}).$$

- **Consistência:** Esta propriedade pode ser informalmente traduzida no seguinte comportamento probabilístico: à medida que aumentamos a dimensão da a.a. X_1, \dots, X_n , o estimador de máxima verosimilhança de θ dispersa-se cada vez menos em torno do verdadeiro valor de θ .

Distribuição amostral

A distribuição de uma estatística, estimador ou sua função é denominada de **distribuição amostral** (ou **distribuição por amostragem**).

Distribuição amostral

A distribuição de uma estatística, estimador ou sua função é denominada de **distribuição amostral** (ou **distribuição por amostragem**).

Exemplo

Distribuição do máximo e do mínimo de uma a.a. X_1, \dots, X_n com função de densidade F .

$$X_{(1)} = \min(X_1, \dots, X_n), \quad X_{(n)} = \max(X_1, \dots, X_n).$$

Importante

- Uma estimativa pontual de um parâmetro não contém informação sobre a **precisão** do valor obtido.
- Uma forma mais completa de abordar a questão consiste em construir estimativas na forma de **intervalos** e conhecer a **probabilidade do intervalo conter o verdadeiro valor do parâmetro**.

Importante

- Assim, em vez de propor apenas uma estimativa isolada $\hat{\theta}$ de θ , faz-se acompanhar esta de um certo intervalo (θ_1, θ_2) .
- Em muitos casos, o intervalo é da forma $(\hat{\theta} - \epsilon, \hat{\theta} + \epsilon)$, em que o valor de ϵ pode ser considerado uma medida de precisão à estimativa $\hat{\theta}$.

Intervalo Aleatório e Intervalo de Confiança

Seja (X_1, \dots, X_n) uma amostra casual de uma população com função densidade $f(x|\theta)$, e $\theta \in \Theta$. Considerem-se duas estatísticas $T_1 \equiv T_1(X_1, \dots, X_n)$ e $T_2 \equiv T_2(X_1, \dots, X_n)$ a verificar $T_1 < T_2$. Um intervalo aleatório para θ é um intervalo (T_1, T_2) tais que

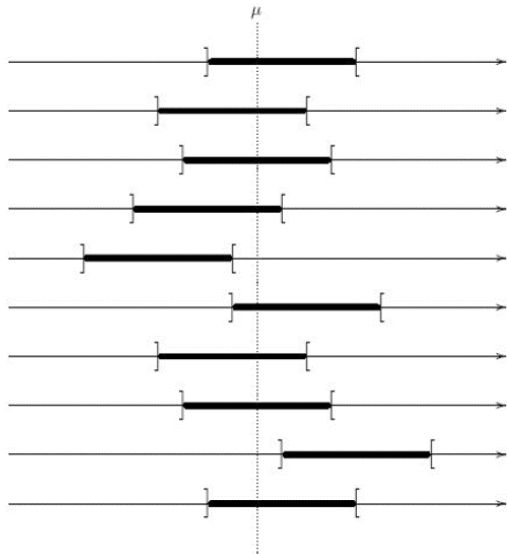
$$P(T_1 < \theta < T_2) = 1 - \alpha, \forall \theta \in \Theta, \alpha \in (0, 1).$$

Intervalo Aleatório e Intervalo de Confiança

Seja (X_1, \dots, X_n) uma amostra casual de uma população com função densidade $f(x|\theta)$, e $\theta \in \Theta$. Considerem-se duas estatísticas $T_1 \equiv T_1(X_1, \dots, X_n)$ e $T_2 \equiv T_2(X_1, \dots, X_n)$ a verificar $T_1 < T_2$. Um intervalo aleatório para θ é um intervalo (T_1, T_2) tais que

$$P(T_1 < \theta < T_2) = 1 - \alpha, \forall \theta \in \Theta, \alpha \in (0, 1).$$

Quando se dispõe de uma amostra particular $x = (x_1, \dots, x_n)$, sejam $t_1 = T_1(x)$ e $t_2 = T_2(x)$ os valores assumidos pelas estatísticas T_1 e T_2 . Então, (t_1, t_2) chama-se intervalo de confiança a $(1 - \alpha) * 100\%$ para θ . O valor $1 - \alpha$ traduz o grau de confiança do intervalo (t_1, t_2) .



Variável fulcral

Seja (X_1, \dots, X_n) uma amostra casual de uma população com função densidade $f(x|\theta)$. A função das observações e de θ , $Z(X_1, \dots, X_n, \theta)$, diz-se uma variável fulcral se a respetiva função densidade f_Z é independente de θ .

Cálculo do intervalo de confiança

- Fixado o grau de confiança, procuram-se dois números a_α e b_α tais que $P(a_\alpha < Z < b_\alpha) = 1 - \alpha$. Uma vez que existe geralmente uma infinidade de pares de valores (a_α, b_α) nessas condições, a opção mais corrente consiste em escolher (a_α, b_α) tais que

$$P(Z < a_\alpha) = P(Z > b_\alpha) = \frac{\alpha}{2}.$$

- Inverter a dupla desigualdade $a_\alpha < Z < b_\alpha$ em ordem a θ

$$P(a_\alpha < Z < b_\alpha) = 1 - \alpha \rightarrow P(T_1 < \theta < T_2) = 1 - \alpha.$$

- Para (x_1, \dots, x_n) , (t_1, t_2) é o IC a $(1 - \alpha) * 100\%$ para θ .

IC para o valor esperado de uma população Normal com variância CONHECIDA

- **Variável fulcral:** $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$;
- $a_\alpha = \Phi^{-1}(\alpha/2)$, $b_\alpha = \Phi^{-1}(1 - \alpha/2)$;
- Inverter a dupla desigualdade $a_\alpha < Z < b_\alpha$ em ordem a μ

$$P(a_\alpha \leq Z \leq b_\alpha) = 1 - \alpha \iff$$

$$P(\bar{X} - \Phi^{-1}(1 - \alpha/2) * \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + \Phi^{-1}(1 - \alpha/2) * \frac{\sigma}{\sqrt{n}}) \\ = 1 - \alpha.$$

IC para o valor esperado de uma população Normal com variância CONHECIDA

$$IC(\mu)_{(1-\alpha)*100\%} = (\bar{x} - \Phi^{-1}(1 - \alpha/2) * \frac{\sigma}{\sqrt{n}}, \bar{x} + \Phi^{-1}(1 - \alpha/2) * \frac{\sigma}{\sqrt{n}}).$$

Exemplo

Uma brigada da GNR pediu ao gerente do bar onde trabalha o Rocha para colaborar numa campanha de sensibilização para o dever de não beber quando se vai conduzir. Assim, o Rocha foi encarregue de pedir aos clientes do bar para soprarem no balão à saída do bar. Ao fim de alguns dias o Rocha já tinha recolhido 101 observações sendo a média amostral 0.6mg/l . De acordo com as instruções da GNR, os valores da quantidade de álcool no sangue das pessoas são bem modelados por uma distribuição Normal de variância 0.04. Calcular um IC a 95% para a verdadeira taxa de alcoolemia dos clientes à saída do bar.

Resolução

Informação inicial: $n = 101$, $\bar{x} = 0.6$, $\sigma^2 = 0.04$.

Intervalo de confiança:

$$IC(\mu)_{(1-0.05)*100\%} = (\bar{x} - \Phi^{-1}(1-0.05/2) * \frac{\sigma}{\sqrt{n}}, \bar{x} + \Phi^{-1}(1-0.05/2) * \frac{\sigma}{\sqrt{n}}).$$

Quantil: $\Phi^{-1}(1 - 0.05/2) = \Phi^{-1}(0.975) = 1.96$ (tabela 4)

ou `qnorm(0.975, 0, 1)`

$$\begin{aligned} IC(\mu)_{(1-0.05)*100\%} &= (0.6 - 1.96 * \frac{\sqrt{0.04}}{\sqrt{101}}, 0.6 + 1.96 * \frac{\sqrt{0.04}}{\sqrt{101}}) \\ &= (0.5609, 0.6390). \end{aligned}$$

Exemplo

Considere uma amostra aleatória X_1, \dots, X_n de uma $N(\mu, 9)$. Um IC para μ é $(\bar{x} - \frac{5.88}{\sqrt{n}}, \bar{x} + \frac{5.88}{\sqrt{n}})$. Determine o grau de confiança do intervalo.

IC para o valor esperado de uma população Normal com variância DESCONHECIDA

- **Variável fulcral:** $Z = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{(n-1)}$;
- $a_\alpha = -F_{t_{(n-1)}}^{-1}(1 - \alpha/2)$, $b_\alpha = -a_\alpha$;
- Inverter a dupla desigualdade $a_\alpha < Z < b_\alpha$ em ordem a μ

$$IC(\mu)_{(1-\alpha)*100\%} =$$
$$(\bar{x} - F_{t_{(n-1)}}^{-1}(1 - \alpha/2) * \frac{s}{\sqrt{n}}, \bar{x} + F_{t_{(n-1)}}^{-1}(1 - \alpha/2) * \frac{s}{\sqrt{n}}).$$

Exercício

No exemplo anterior (bar do Rocha) seria mais realista **não** assumir um valor conhecido para a variância da distribuição pois, em geral, não é possível determinar com exatidão esse valor. Com base na amostra de valores recolhidos pelo Rocha a variância amostral corrigida foi de **0.06**. Calcular o IC de confiança a **95%** para a taxa média de alcoolemia.

Resolução

Informação inicial: $n = 101$, $\bar{x} = 0.6$, $s^2 = 0.06$.

Intervalo de confiança:

$$IC(\mu)_{(1-0.05)*100\%} = (\bar{x} - F_{t_{(n-1)}}^{-1}(1 - 0.05/2) * \frac{s}{\sqrt{n}}, \bar{x} + F_{t_{(n-1)}}^{-1}(1 - \alpha/2) * \frac{s}{\sqrt{n}}).$$

Quantil: $F_{t_{(100)}}^{-1}(0.975) = 1.984$ (tabela 5) ou $qt(0.975, 100)$

$$\begin{aligned} IC(\mu)_{(1-0.05)*100\%} &= (0.6 - 1.984 * \frac{\sqrt{0.06}}{\sqrt{101}}, 0.6 + 1.984 * \frac{\sqrt{0.06}}{\sqrt{101}}) \\ &= (0.5516, 0.6483). \end{aligned}$$

Exercício

De uma certa população observou-se uma amostra aleatória de dimensão $n = 151$, tendo-se obtido $\bar{x} = 75$ e $s = 10$. Construa intervalos de confiança aproximados de 95% e 99% para μ .

IC para a variância de uma população Normal com média DESCONHECIDA

- **Variável fulcral:** $Z = \frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$;
- $a_\alpha = F_{\chi_{n-1}^2}^{-1}(\alpha/2)$, $b_\alpha = F_{\chi_{n-1}^2}^{-1}(1 - \alpha/2)$;
- Inverter a dupla desigualdade $a_\alpha < Z < b_\alpha$ em ordem a σ^2

$$P(a_\alpha \leq Z \leq b_\alpha) = 1 - \alpha \iff$$

$$P\left(\frac{(n-1)S^2}{b_\alpha} \leq \sigma^2 \leq \frac{(n-1)S^2}{a_\alpha}\right) = 1 - \alpha.$$

$$IC(\sigma^2)_{(1-\alpha)*100\%} = \left(\frac{(n-1)s^2}{F_{\chi_{n-1}^2}^{-1}(1-\alpha/2)}, \frac{(n-1)s^2}{F_{\chi_{n-1}^2}^{-1}(\alpha/2)} \right).$$

Exercício

Suponha-se que o tempo de vida em horas de certas componentes eletrónicas produzidas segundo determinado processo de fabrico tem distribuição Normal (μ, σ^2) . Obtida uma amostra de 20 componentes, verificou-se a duração das mesmas, sendo $\bar{x} = 1832$ e $s = 479$. Determinar um IC a 95% para σ^2 .

Resolução

Informação inicial: $n = 20$, $\bar{x} = 1832$, $s = 479$

Quantis: $F_{\chi^2_{19}}^{-1}(0.025) = 8.907$, $F_{\chi^2_{19}}^{-1}(0.975) = 32.85$ (tabela 6) ou

`qchisq(.025, 19)` e `qchisq(.975, 19)`

Intervalo de confiança:

$$\begin{aligned} IC(\sigma^2)_{(1-0.05)*100\%} &= \left(\frac{19 * 479^2}{32.85}, \frac{19 * 479^2}{8.907} \right) \\ &= (132705.6, 489432.9). \end{aligned}$$

IC aproximado para uma probabilidade

- $(X_1, \dots, X_n) \quad X \sim Be(p)$;
- n grande ($n > 30$);
- **Variável fulcral:** $Z = \frac{\bar{X} - p}{\sqrt{\frac{\bar{X}(1-\bar{X})}{n}}} \sim_{TLC} N(0, 1)$;
- $a_\alpha = \Phi^{-1}(\alpha/2)$, $b_\alpha = \Phi^{-1}(1 - \alpha/2)$;
- Inverter a dupla desigualdade $a_\alpha < Z < b_\alpha$ em ordem a p

$$IC(p)_{(1-\alpha)*100\%} = \left(\bar{x} - \Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \sqrt{\frac{\bar{x}(1-\bar{x})}{n}}, \bar{x} + \Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \sqrt{\frac{\bar{x}(1-\bar{x})}{n}} \right)$$

Exercício

Um candidato a deputado por determinado círculo uninominal teve conhecimento de que, numa sondagem feita junto de $n = 350$ eleitores, houve 185 que declararam dar-lhe o seu voto. Determine um IC a 95% para a probabilidade p de uma pessoa escolhida ao acaso nesse círculo uninominal votar no candidato.

Inferência Estatística: Testes de Hipóteses

Hipótese estatística

Uma **hipótese estatística** é uma conjectura sobre uma característica da população.

Teste de hipóteses

Um **teste de hipóteses** é um procedimento estatístico que averigua se os dados sustentam uma hipótese.

Inferência Estatística: Testes de Hipóteses

Hipótese estatística

Uma **hipótese estatística** é uma conjectura sobre uma característica da população.

Teste de hipóteses

Um **teste de hipóteses** é um procedimento estatístico que averigua se os dados sustentam uma hipótese.

As hipóteses: Num teste de hipóteses há sempre duas hipóteses:

Hipótese Nula — H_0 *vs* Hipótese alternativa — H_1

Tipos de hipótese e tipos de testes

- As várias hipóteses podem ser **simples** ou **compostas**. Uma hipótese simples apenas contempla uma possibilidade ($=$).
- Os testes podem ser **unilaterais** ou **bilaterais**. Nos testes unilaterais a hipótese alternativa apenas contempla possibilidades à direita ou à esquerda da hipótese nula.

Exemplos de testes unilaterais

- Testes unilaterais à direita:

$$H_0 : \mu = 1 \text{ vs } H_1 : \mu > 1$$

$$H_0 : \mu = 4 \text{ vs } H_1 : \mu = 7$$

- Testes unilaterais à esquerda:

$$H_0 : \mu = 1 \text{ vs } H_1 : \mu < 1$$

$$H_0 : \mu = 4 \text{ vs } H_1 : \mu = 2$$

Testes bilaterais

Nos testes **bilaterais** a hipótese alternativa contempla valores à direita e à esquerda da hipótese nula.

Exemplo de testes bilateral

- Teste bilateral:

$$H_0 : \mu = 1 \text{ vs } H_1 : \mu \neq 1$$

Estatística de teste

É uma estatística calculada a partir da amostra e que é usada para tomar a decisão acerca de rejeitar, ou não, a hipótese nula. Costuma-se representar por T .

Região de rejeição

A região de rejeição (ou região crítica, RC) é o conjunto de valores da estatística T que nos levam a rejeitar a hipótese nula.

Se o valor observado de T , (t_{obs}), pertencer à região crítica, rejeita-se H_0 a favor de H_1 . Caso contrário não se rejeita H_0 .

Inferência Estatística: Testes de Hipóteses

		Hipótese verdadeira	
		H_0	H_1
Decisão do teste	Rejeito H_0	Erro de tipo I	✓
	Não rejeito H_0	✓	Erro de tipo II

$$P(\text{Erro de tipo I}) = \alpha \quad P(\text{Erro de tipo II}) = \beta.$$

Tamanho do teste ou Nível de significância:

$$\alpha = P(\text{Erro de tipo I}) = P(\text{rejeitar } H_0 | H_0 \text{ verdadeiro}).$$

Potência do teste:

$$1 - \beta = 1 - P(\text{Erro de tipo II}) = P(\text{rejeitar } H_0 | H_1 \text{ verdadeiro}).$$

Importante

Procedimentos para a realização de um TH de tamanho $\alpha \in (0, 1)$.

Procedimento com base na região de rejeição

- 1 Identificar o parâmetro de interesse e especificar H_0 e H_1 ;
- 2 Escolher uma estatística de teste T , com distribuição conhecida (admitindo que H_0 é verdadeira);
- 3 Identificar a região de rejeição;
- 4 Calcular o t_{obs} ;
- 5 Se $t_{obs} \in RC$ então rejeita-se H_0 . Caso contrário não se rejeita;
- 6 Concluir.

Teste para a média de uma distribuição Normal com variância CONHECIDA

- 1 Definir H_0 e H_1 : $H_0 : \mu = \mu_0$ vs $H_1 : \mu <, \neq, > \mu_0$
- 2 Estatística de teste:

$$T = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \underset{\text{sob } H_0}{\sim} N(0, 1)$$

- 3 Região de rejeição:

Se $\mu \neq \mu_0$ então $RC = \{t_{obs} : |t_{obs}| > z_{1-\alpha/2}\};$

Se $\mu < \mu_0$ então $RC = \{t_{obs} : t_{obs} < z_{\alpha}\};$

Se $\mu > \mu_0$ então $RC = \{t_{obs} : t_{obs} > z_{1-\alpha}\}.$

Teste para a média de uma distribuição Normal com variância DESCONHECIDA

- 1 Definir H_0 e H_1 : $H_0 : \mu = \mu_0$ vs $H_1 : \mu <, \neq, > \mu_0$
- 2 Estatística de teste:

$$T = \frac{\bar{X} - \mu_0}{S_c / \sqrt{n}} \underset{\text{sob } H_0}{\sim} t_{n-1}$$

- 3 Região de rejeição:

Se $\mu \neq \mu_0$ então $RC = \{t_{obs} : |t_{obs}| > t_{1-\alpha/2, n-1}\};$

Se $\mu < \mu_0$ então $RC = \{t_{obs} : t_{obs} < t_{\alpha, n-1}\};$

Se $\mu > \mu_0$ então $RC = \{t_{obs} : t_{obs} > t_{1-\alpha, n-1}\}.$

Exemplo

Um serviço adquiriu, para equipar as suas instalações, lâmpadas da marca A, aceitando a afirmação do fabricante de que a vida média das lâmpadas é 1650 horas. Numa experiência feita com 50 lâmpadas obteve-se $\bar{x} = 1575$ horas e $s_c = 120$ horas. Pode admitir-se correta a afirmação do fabricante, admitindo que os tempos de vida têm distribuição Normal?

- 1 Hipóteses: $H_0 : \mu = 1650$ vs $H_1 : \mu \neq 1650$
- 2 Estatística de teste:

$$T = \frac{\bar{X} - \mu_0}{S_c / \sqrt{n}} \underset{\text{sob } H_0}{\sim} t_{n-1}$$

Exemplo (cont)

1 Região de rejeição:

$$\begin{aligned} RC &= \{t_{obs} : |t_{obs}| > t_{1-\alpha/2, n-1}\} \\ &= \{t_{obs} : |t_{obs}| > t_{0.975, 49}\} \\ &= \{t_{obs} : |t_{obs}| > 2.009\}. \end{aligned}$$

2 Valor observado da estatística de teste:

$$t_{obs} = \frac{1575 - 1650}{120/\sqrt{50}} = -4.41.$$

Teste para a média numa população genérica, $n \geq 30$

- 1 Definir H_0 e H_1 : $H_0 : \mu = \mu_0$ vs $H_1 : \mu <, \neq, > \mu_0$
- 2 Estatística de teste:

$$T = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \underset{\text{sob } H_0}{\overset{\circ}{\sim}} N(0, 1), \sigma^2 \text{ conhecida}$$

$$T = \frac{\bar{X} - \mu_0}{S_c/\sqrt{n}} \underset{\text{sob } H_0}{\overset{\circ}{\sim}} N(0, 1), \sigma^2 \text{ desconhecida}$$

Teste para a variância de uma distribuição Normal com média DESCONHECIDA

① Definir H_0 e H_1 : $H_0 : \sigma^2 = \sigma_0^2$ vs $H_1 : \sigma^2 <, \neq, > \sigma_0^2$

② Estatística de teste:

$$T = \frac{(n-1)S_c^2}{\sigma_0^2} \underset{\text{sob } H_0}{\sim} \chi_{n-1}^2$$

③ Região de rejeição: se $\sigma^2 \neq \sigma_0^2$ então

$$RC = \{t_{obs} : t_{obs} < \chi_{\alpha/2, n-1}^2 \text{ ou } t_{obs} > \chi_{1-\alpha/2, n-1}^2\};$$

$$\text{Se } \sigma^2 < \sigma_0^2 \text{ então } RC = \{t_{obs} : t_{obs} < \chi_{\alpha, n-1}^2\};$$

$$\text{Se } \sigma^2 > \sigma_0^2 \text{ então } RC = \{t_{obs} : t_{obs} > \chi_{1-\alpha, n-1}^2\}.$$

Exercício

O diâmetro interior das porcas feitas numa determinada fábrica é uma variável aleatória $X \sim N(\mu, \sigma^2)$. Pretende-se testar

$$H_0 : \sigma^2 = 0.7 \text{ vs } H_1 : \sigma^2 > 0.7$$

Determine a região de rejeição a partir de uma amostra de dimensão 16. Considere $\alpha = 0.05$. Na amostra de 16 porcas obteve-se $s_c^2 = 0.96$. Qual a conclusão do teste com base nesta amostra?

Exercício

- Região de rejeição:

$$\begin{aligned}RC &= \{t_{obs} : t_{obs} > \chi^2_{1-\alpha, n-1}\} \\&= \{t_{obs} : t_{obs} > \chi^2_{0.95, 15}\} \\&= \{t_{obs} : t_{obs} > 25.00\}\end{aligned}$$

- Valor observado da estatística de teste:

$$t_{obs} = \frac{15 * 0.96}{0.70} = 20.57.$$

Teste para uma proporção

- 1 Definir H_0 e H_1 : $H_0 : p = p_0$ vs $H_1 : p <, \neq, > p_0$
- 2 Estatística de teste:

$$T = \frac{\hat{p} - p_0}{\sqrt{p_0(1 - p_0)/n}} = \frac{\frac{X}{n} - p_0}{\sqrt{p_0(1 - p_0)/n}} \underset{\text{sob } H_0}{\overset{\circ}{\sim}} N(0, 1)$$

- 3 Região de rejeição:

Se $p \neq p_0$ então $RC = \{t_{obs} : |t_{obs}| > z_{1-\alpha/2}\};$

Se $p < p_0$ então $RC = \{t_{obs} : t_{obs} < z_{\alpha}\};$

Se $p > p_0$ então $RC = \{t_{obs} : t_{obs} > z_{1-\alpha}\}.$

Exemplo

Numa sondagem à opinião pública, em dado círculo eleitoral, foram inquiridas 1000 pessoas e houve 43% que se disseram favoráveis a determinado partido político. Rejeita-se a hipótese de este partido ter 50% das preferências naquele círculo?

$$H_0 : p = 0.5 \text{ vs } H_1 : p \neq 0.5$$

- ① Região de rejeição: $RC = \{t_{obs} : |t_{obs}| > z_{0.975} = 1.96\}$
- ② Valor observado da estatística de teste:

$$t_{obs} = \frac{0.43 - 0.50}{\sqrt{0.5 \times 0.5 / 1000}} = -4.42$$

Procedimento com base no p -value

- 1 Identificar o parâmetro de interesse e especificar H_0 e H_1 ;
- 2 Calcular o valor observado da estatística de teste;
- 3 Determinar o p -value do teste;

O p -value do teste é a probabilidade de observar um valor da estatística de teste tanto ou mais afastado que o valor observado na amostra, assumindo que H_0 é verdadeira.

- 4 Rejeitar H_0 se $p\text{-value} \leq \alpha$. Não rejeitar H_0 se $p\text{-value} > \alpha$.

Cálculo do p -value

- Quando a RC é da forma $T > c$ o $p\text{-value} = P(T > t_{obs} | H_0)$;
- Quando a RC é da forma $T < c$ o $p\text{-value} = P(T < t_{obs} | H_0)$;
- Quando a RC é da forma $T < c_1$ ou $T > c_2$ (com igual probabilidade nos dois casos) o $p\text{-value}$ é igual a

$$\begin{cases} 2P(T < t_{obs} | H_0) & \text{se } t_{obs} \text{ for reduzido} \\ 2P(T > t_{obs} | H_0) & \text{se } t_{obs} \text{ for elevado} \end{cases};$$

- Dizer que t_{obs} é reduzido (elevado) significa dizer que a estimativa que se obtém para o parâmetro a testar é inferior (superior) ao valor especificado em H_0 .

Exemplo

Na empresa onde o Rocha trabalha de dia, faz-se regularmente o controlo de qualidade dos medicamentos aí produzidos. Um desses medicamentos é constituído por cápsulas que devem ter 20mg de fluoxetina cada uma. Regularmente são retiradas amostras de $n = 100$ cápsulas com as quais se realiza o seguinte teste:

$$H_0 : \mu = 20 \text{ vs } H_1 : \mu \neq 20,$$

onde μ representa a verdadeira média da quantidade de fluoxetina contida nas cápsulas. Suponhamos que uma certa amostra forneceu $\bar{x} = 20.04$ e $s_c = 0.2$. Para o nível de significância $\alpha = 0.001$ que podemos concluir?

Exemplo, cont

- **Procedimento com base na região de rejeição**

Neste caso a região de rejeição é do tipo:

$$RC = \{t_{obs} : |t_{obs}| > t_{0.9995,99} = 3.39\}.$$

Valor observado da estatística de teste:

$$t_{obs} = \frac{20.04 - 20}{0.2/\sqrt{100}} = 2.$$

Como t_{obs} não pertence à região crítica não se rejeita H_0 .

Exemplo, cont

- **Procedimento com base no intervalo de confiança**

Um intervalo de confiança para μ a 99.9% é dado por

$$\begin{aligned} IC(\mu)_{99.9\%} &= \left(\bar{x} - t_{0.9995,99} \frac{s_c}{\sqrt{n}}, \bar{x} + t_{0.9995,99} \frac{s_c}{\sqrt{n}} \right) \\ &= \left(20.04 - 3.39 \times \frac{0.2}{10}, 20.04 + 3.39 \times \frac{0.2}{10} \right) \\ &= (19.97, 20.11). \end{aligned}$$

Como $20 \in IC(\mu)_{99.9\%}$, não se rejeita H_0 .

Exemplo, cont

- **Procedimento com base no p -value**

Neste caso $\bar{x} = 20.04 > 20$, pelo que o p -value é

$$2P(T > t_{obs}|H_0) = 2P(T > 2|H_0) = 0.048.$$

Como este valor é superior ao $\alpha = 0.001$ não se rejeita H_0 .