

No ficheiro de dados [econ.xlsx](#) encontra informação relativa a dados económicos dos Estados Unidos providenciados pela empresa [FRED](#). Este conjunto de dados possui as seguintes variáveis: tempo (Data do registo); gcp (gastos de consumo pessoal, em biliões de dólares); pop (população total); tpp (taxa de poupança pessoal); ddesemp (duração mediana do desemprego, em semanas); ndesemp (número de desempregados, em milhares).

Considere as variáveis  $x_1$  – **pop** e  $x_2$  – **tpp** para os anos superiores ou iguais a **1989**. Com recurso ao pacote **ggplot** produza um único gráfico que lhe permita fazer uma análise da evolução dessas duas variáveis para esses anos.

Uma vez que as variáveis podem não ter a mesma escala, antes de construir o gráfico proceda do seguinte modo:

- Selecione os dados a usar.
- Faça a seguinte transformação aos dados associados a cada variável

$$X_k : z_{ik} = \frac{x_{ik} - \bar{x}_k}{sx_k}, \quad i = 1, 2, \dots, n,$$

onde  $n$  é a dimensão da amostra,  $x_k$  e  $s_{x_k}$  correspondem, respectivamente, à média e desvio-padrão da amostra associada à variável  $x_k$ .

Submeta um ficheiro em formato PDF, com uma única página A4, que inclua:

- O código em **R**.  
**Nota:** no código devem também constar os comandos para leitura e seleção dos dados do ficheiro.
- O gráfico que achar mais adequado para analisar a evolução dessas variáveis nesse período de tempo.

O ficheiro [TIME\\_USE\\_24092022.csv](#) contém uma compilação de dados enviados por diversos países para a OCDE ([Organização para a Cooperação e Desenvolvimento Económico](#)) sobre o tempo médio diário (em minutos) despendido pelas pessoas entre os 15 e os 64 anos em diferentes tipos de ocupações.

- Leia o ficheiro de dados no **R** e elimine todos os registos referentes à África do Sul (dados incompletos).
- Submeta um ficheiro em formato PDF com uma única página A4, que inclua, num único gráfico, dois diagramas de extremos e quantis que permitam comparar os tempos médios diários registados para **Mulheres** em duas ocupações distintas: **Outros** e **Trabalho remunerado ou estudo**.  
**Nota:** o código apresentado deve incluir os comandos para leitura e seleção dos dados do ficheiro.

O ficheiro [GENDER\\_EMP\\_19032023152556091.txt](#) contém uma compilação de dados sobre emprego enviados por diversos países para a OCDE ([Organização para a Cooperação e Desenvolvimento Económico](#)).

Com recurso ao pacote **ggplot** produza um único gráfico de barras que permita comparar os valores da variável **EMP5** (*Share of employed in part-time employment, by sex and age group*) entre homens e mulheres nos grupos etários 15–24, 25–54 e 55–64, registados em **2015** no seguinte país: **Germany**.

Por simplicidade, mantenha todo o texto no gráfico em Inglês.

Submeta um ficheiro em formato PDF com uma única página A4, que inclua:

- O código em **R**, que deve incluir os comandos para leitura e seleção dos dados do ficheiro.
- O gráfico produzido.

- Fixando a semente em 4762, gere uma amostra de dimensão  $n = 1000$  proveniente de uma distribuição Exponencial de parâmetro  $\lambda = 36$ . Os valores gerados correspondem aos tempos entre acontecimentos sucessivos.
- Considere agora a soma sucessiva destas observações, i.e., se  $x_j$  designar o  $j$ -ésimo valor gerado, então  $y_j = \sum_{i=1}^j x_i$  é o instante de ocorrência do  $j$ -ésimo acontecimento. Seja  $r = \lfloor r_{0.200} \rfloor$  o menor número inteiro maior ou igual ao instante de ocorrência do último acontecimento.
- Divida o intervalo  $[0, y_r]$  em intervalos de amplitude unitária e contabilize o número de acontecimentos que ocorreram em cada um desses subintervalos.
- Calcule a média do número de acontecimentos por subintervalo e de seguida calcule o desvio absoluto entre este valor e o valor esperado (teórico) do número de acontecimentos num subintervalo. Indique este desvio arredondado a 4 casas decimais.

Ensaio de Bernoulli Independentes, cada um dos quais com probabilidade de sucesso  $p = 0.3$ , são sucessivamente realizados. Seja  $x$  o número de insucessos até ao primeiro ensaio que resulta em sucesso. A distribuição da variável aleatória  $x$  é conhecida por distribuição geométrica de parâmetro  $p = 0.3$ , cuja função (massa) de probabilidade é dada por:

$$f_X(x) = \begin{cases} (1 - p)^x p, & x = 0, 1, 2, \dots \\ 0, & \text{caso contrário.} \end{cases}$$

Podemos gerar valores de uma distribuição geométrica a partir de uma distribuição uniforme contínua usando o **método de transformação Inversa**. Nesse sentido, requer-se a execução dos seguintes passos:

- Simula-se um valor,  $u$ , proveniente de uma distribuição uniforme no intervalo  $[0, 1]$ .
- Se  $p_X(x - 1) < u \leq p_X(x)$ , aceita-se  $x$  como um valor simulado de  $x$ , onde  $p_X(x)$  é a função de distribuição de  $x$ .

Fixando a semente em 1006, implemente este método de simulação estocástica repetindo os passos anteriores até obter uma amostra de dimensão  $n = 1000$ .

Indique a proporção de valores simulados que são superiores à soma da média com o desvio padrão amostrais, de entre os que são superiores à respetiva média amostral. Apresente o resultado com 4 casas decimais.

Considere a variável aleatória  $x$  que representa o primeiro algarismo de um número inteiro escrito em base decimal. Admita que  $x$  possui distribuição de Benford, com função de probabilidade dada por:

$$P(X = x) = \log_{10} \left( 1 + \frac{1}{x} \right), \quad x \in \{1, 2, \dots, 9\}.$$

- Calcule a probabilidade de  $x$  ser igual a 1 ou 8.
- Obtenha a fração de potências de dois no intervalo  $[x^{\frac{1}{10}}, x^{\frac{1}{9}}]$  cujo primeiro algarismo é igual a 1 ou 8.
- Calcule o desvio absoluto entre os valores calculados em 1. e 2.
- Indique este desvio arredondado a 4 casas decimais.

Fixando a semente em 1003, simule  $n = 2000$  amostras de dimensão  $n = 10$  de uma população normal de média nula e variância unitária. Para cada uma das amostras, calcule a soma dos quadrados dos valores observados.

Indique a diferença em valor absoluto (arredondado a 4 casas decimais), entre o quantil de probabilidade  $0.95$  da amostra das somas dos quadrados dos valores observados e o quantil correspondente à distribuição teórica da soma de quadrados de variáveis normais reduzidas independentes.

**Nota:** Use a função **quantile** com a opção **type=2**.

Considere uma variável aleatória com distribuição de Cauchy, com parâmetros de localização e escala iguais a 0 e 1,0, respectivamente.

Usando o R e fixando a semente em 12345, gere uma amostra de dimensão  $n = 100$  desta população.

Represente num único gráfico:

1. Os valores gerados ordenados por ordem crescente versus os quantis de probabilidade  $\Phi^{-1}(i/101 + 0.5)$ ,  $i = 1, \dots, 100$  desta população.
2. Os valores gerados ordenados por ordem crescente versus os quantis de probabilidade  $\Phi^{-1}(i/101 + 0.5)$ ,  $i = 1, \dots, 100$  de uma população normal com valor esperado  $\mu = 3.2$  e variância  $\sigma^2 = 4$ .
3. A recta bissectriz dos quadrantes ímpares.

Submeta um ficheiro em formato PDF, com uma única página A4, que inclua:

1. O código em R.
2. O gráfico produzido.

Para a construção de intervalos de confiança para o parâmetro  $p$  de uma distribuição de Bernoulli podemos recorrer à variável fulcral

$$Z_1 = \frac{\bar{X} - p}{\sqrt{\frac{p(1-p)}{n}}} \stackrel{d}{=} N(0, 1)$$

obtida pela aplicação do teorema do limite central a uma amostra aleatória de tamanho  $n$  suficientemente grande da referida população. Duas variantes são possíveis:

**Método 1**

Usando  $n_1$ , não é difícil mostrar que os limites do intervalo de confiança são as soluções da seguinte equação do segundo grau em  $x$

$$\bar{x}^2 - 2p\bar{x} + p^2 - z^2 \frac{p(1-p)}{n} = 0,$$

em que  $\bar{x}$  representa a média amostral e  $z = z^*(\frac{1+\alpha}{2})$  para um nível de confiança aproximado  $1-\alpha$ .

**Método 2**

Uma segunda aproximação conduz à variável fulcral

$$Z_2 = \frac{\bar{X} - p}{\sqrt{\frac{\bar{X}(1-\bar{X})}{n}}} \stackrel{d}{=} N(0, 1)$$

que permite a construção de intervalos de confiança de uma forma mais simples e habitual.

Com o objetivo de comparar os dois métodos e, em particular, avaliar a adequação da segunda aproximação, implemente os seguintes passos no R:

1. Fixe a semente em 12345 e para cada valor de  $n \in \{10, 100, 1000, 10000, 100000\}$ :
  - a. gere  $n - 500$  amostras de tamanho  $n$  de uma distribuição de Bernoulli com parâmetro  $p = 0.5$ ;
  - b. para cada amostra gerada, calcule a diferença entre os comprimentos dos intervalos de confiança construídos pelo **Método 2** e pelo **Método 1**, com um nível de confiança aproximado  $1 - 0.05$ ;
  - c. calcule a média das  $n - 500$  diferenças anteriores.
2. Construa um gráfico que ilustre a variação das diferenças médias em função do tamanho da amostra.

Submeta um ficheiro em formato PDF, com uma única página A4, que inclua:

1. O código em R.
2. O gráfico pedido.
3. Comentários sobre os resultados obtidos.

Considere uma variável aleatória  $x$  com distribuição Normal de valor esperado  $\mu$  desconhecido e variância  $\sigma^2 = 4$ . Construa um teste de hipóteses  $H_0: \mu = 0.5$  contra  $H_1: \mu \neq 0.5$ , ao nível de significância de  $\alpha = 0.05$ .

Com recurso ao R e fixando a semente em 12345, gere  $n = 500$  amostras de dimensão  $n = 40$  dessa variável, admitindo que  $\mu = 0.3$ . Aplique o teste de hipóteses que construiu para cada amostra gerada, e use o conjunto de resultados para obter uma estimativa da probabilidade do teste conduzir à não rejeição de  $H_0$ . Indique o resultado com 3 casas decimais.