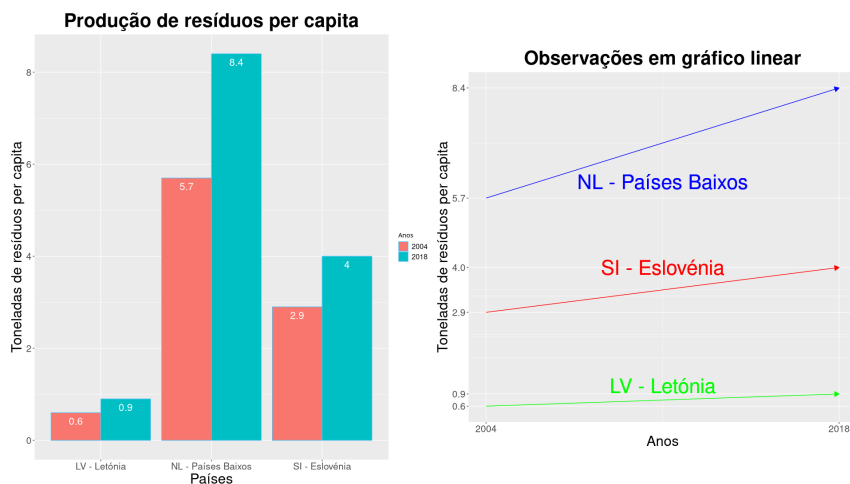


Pergunta 1

```

1 pacman::p_load(pacman, rio, tidyverse, ggplot2, cowplot) # Load req. packages
2 ## Dados
3 residuos <- import("~/Projects/git/personal/PEst/Code_env/ResiduosPerCapita.xlsx")
4 req <- c("NL - Países Baixos", "SI - Eslovénia", "LV - Letónia")
5 temp <- filter(residuos, grepl(paste(req, collapse="|"), residuos$...1))
6 result <- temp[, c(1,2,3)] # Selecionar as colunas
7 aux <- data.frame(
8   Anos=c(2004,2018),
9   NL=as.numeric(c(result[1,2:3])),
10  SI=as.numeric(c(result[2,2:3])),
11  LV=as.numeric(c(result[3,2:3])),
12  stringsAsFactors = FALSE
13 )
14 names(result)[names(result) == "...1"] <- "Países"
15 names(result)[names(result) == "...3"] <- "Anos"
16 names(result)[names(result) == "Produção de resíduos per capita"] <- "2004"
17 # Formatar os dados para o gráfico de barras
18 data <- gather(result, "Anos", "Toneladas de resíduos per capita", 2:3)
19 data[3] = as.numeric(c(data[[3]]))
20 ## Gráfico de barras
21 p1 <- ggplot(data=data, aes(x=Países, y=`Toneladas de resíduos per capita`, fill=Anos)) +
22   geom_bar(stat="identity", position=position_dodge()) +
23   ggtitle("Produção de resíduos per capita") +
24   theme(plot.title=element_text(hjust=0.5, size=32, face="bold"),
25         axis.title=element_text(size=24),
26         axis.text=element_text(size=16),
27         legend.title=element_text(size=16),
28         legend.text=element_text(size=12)) +
29   geom_text(aes(label=Toneladas de resíduos per capita`, vjust=1.6, color="white",
30               position=position_dodge(0.9), size=6)
31 ## Gráfico auxiliar
32 p2 <- ggplot(data=aux, aes(x=Anos), fill=NULL) + ggtitle("Observações em gráfico linear") +
33   theme(plot.title=element_text(hjust=0.5, size=32, face="bold"),
34         axis.title=element_text(size=24),
35         axis.text=element_text(size=16),
36         aspect.ratio=0.9) +
37   scale_x_continuous(breaks=as.numeric(unlist(aux[,1]))) +
38   scale_y_continuous(breaks=as.numeric(unlist(aux[,2:4]))) +
39   geom_line(aes(y=NL), color="red", arrow = arrow(length=unit(0.30,"cm"), type = "closed")) +
40   geom_line(aes(y=SI), color="green", arrow = arrow(length=unit(0.30,"cm"), type = "closed")) +
41   geom_line(aes(y=LV), color="blue", arrow = arrow(length=unit(0.30,"cm"), type = "closed")) +
42   geom_text(aes(x=2011, y=4.0), label="SI - Eslovénia", color="red", size=12) +
43   geom_text(aes(x=2011, y=1.1), label="LV - Letónia", color="green", size=12) +
44   geom_text(aes(x=2011, y=6.1), label="NL - Países Baixos", color="blue", size=12) +
45   labs(x="Anos", y="Toneladas de resíduos per capita")
46 ## Plot do gráfico de barras e gráfico auxiliar
47 plot <- plot_grid(p1, p2, nrow=1, ncol= 2, labels=NULL); plot
48 ## Média de todos os países
49 mean04 = mean(as.numeric(residuos$`Produção de resíduos per capita`[7:36]))
50 mean18 = mean(as.numeric(residuos$...3[7:36]))
51 sprintf("Média prod. resíduos em 2004: %.2f", mean04)
52 sprintf("Média prod. resíduos em 2018: %.2f", mean18)

```

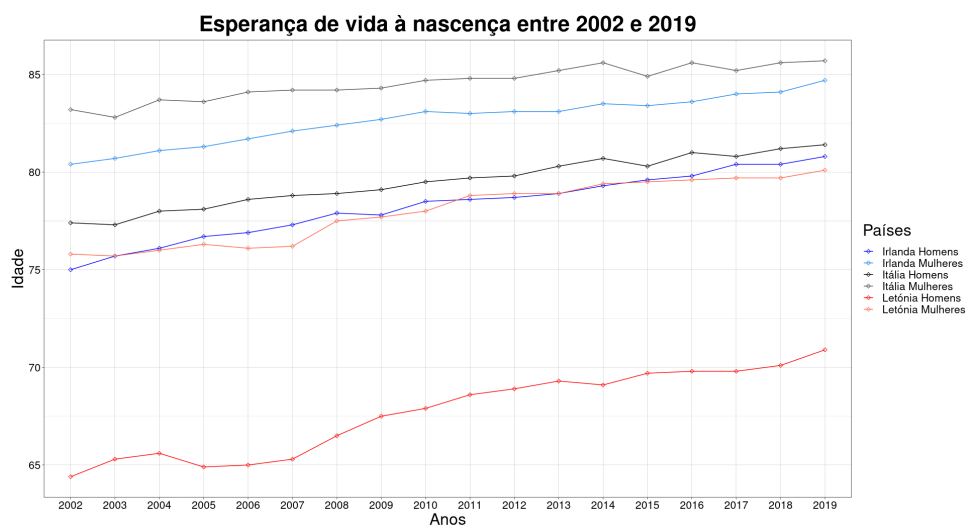


Comentário: Por observação direta do diagrama de barras e do gráfico linear dos três países europeus no ano 2004 e 2018, é necessariamente possível concluir que ocorreu um aumento de produção de resíduos nos três países. Salienta-se que os Países Baixos foram o maior produtor de resíduos (com o maior acréscimo também), enquanto a Letónia o menor.

Naturalmente, foi calculada a média da produção de resíduos de todos os países nos anos 2004 e 2018 por efeitos comentaristas. Os valores obtidos foram 6.14 e 6.21 toneladas de resíduos per capita, respetivamente para 2004 e 2018. Note-se que a Letónia se encontra tremendamente abaixo destes valores, e que em 2004 o grupo dos três países encontrava-se abaixo da média. Em 2018, os Países Baixos tiveram uma produção acima da média.

Pergunta 2

```
1 pacman::p_load(pacman, rio, tidyverse, ggplot2, hrbrthemes) # Load req. packages
2 ## Dados
3 esperanca_vida <- import("~/Projects/git/personal/PEst/Code_env/EsperancaVida.xlsx")
4 result <- esperanca_vida[c(48:65), c(1, 54, 55, 56, 88, 89, 90)]
5 result$...54 <- as.numeric(esperanca_vida$...54[48:65])
6 result$...55 <- as.numeric(esperanca_vida$...55[48:65])
7 result$...56 <- as.numeric(esperanca_vida$...56[48:65])
8 result$...88 <- as.numeric(esperanca_vida$...88[48:65])
9 result$...89 <- as.numeric(esperanca_vida$...89[48:65])
10 result$...90 <- as.numeric(esperanca_vida$...90[48:65])
11 # Rename das colunas
12 names(result)[names(result) == "...1"] <- "Anos"
13 names(result)[names(result) == "...56"] <- "Letónia Homens"
14 names(result)[names(result) == "...55"] <- "Itália Homens"
15 names(result)[names(result) == "...54"] <- "Irlanda Homens"
16 names(result)[names(result) == "...90"] <- "Letónia Mulheres"
17 names(result)[names(result) == "...89"] <- "Itália Mulheres"
18 names(result)[names(result) == "...88"] <- "Irlanda Mulheres"
19 gtemporal <- gather(result, "Países", "Esperança de vida à nascença", 2:7)
20 ## Gráfico temporal
21 ggplot(gtemporal, aes(x="Anos", y="Esperança de vida à nascença", color="Países",
22 group="Países")) +
23   geom_point(size=2, shape=23) +
24   geom_line() +
25   labs(x="Anos", y="Idade") +
26   scale_color_manual(values=c('blue', 'dodgerblue2', 'black', 'gray31', 'red',
27 'tomato1')) +
28   labs(color="Países", fill="Países") +
29   ggtitle("Esperança de vida à nascença entre 2002 e 2019") +
30   theme_linedraw() +
31   theme(plot.title = element_text(hjust = 0.5, size=32, face="bold"),
32         axis.title=element_text(size=24),
33         axis.text=element_text(size=16),
34         legend.title=element_text(size=24),
35         legend.text=element_text(size=16))
```



Comentário: Independentemente do género e do país, é natural concluir que a esperança média de vida tem aumentado ao longo dos anos, graças aos dados desencadeados. É verificável que a aparente trend linear se mantém e que não se apresentam comportamentos cíclicos nem sazonais - o fenómeno é expectável, visto que o país de nascença afeta estridulamente a esperança de vida¹.

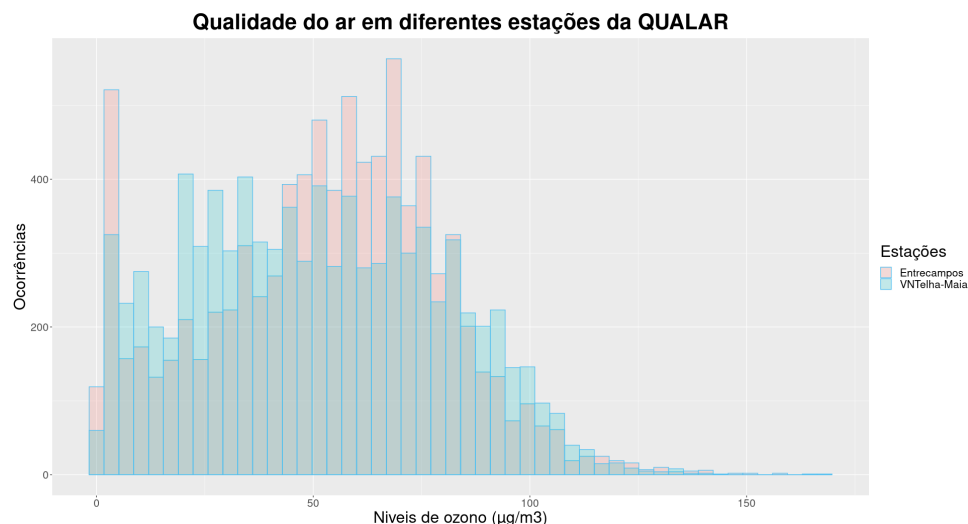
Por observação direta do gráfico temporal apresentado em seguida, destaca-se também que a esperança média de vida à nascença é superior para as mulheres em relação à dos homens em todos os casos.

É também salientado que a Itália é o país com maior esperança média de vida do grupo analisado, seguido da Irlanda e, por final, pela Eslováquia (especialmente no caso dos homens, em que a diferença é aproximadamente 10 anos, em comparação com os outros dois países).

¹Eric Neumayer and Thomas Plümper. Inequalities of income and inequalities of longevity: A cross-country study. American Journal of Public Health, 106(1):160–165, 2016. PMID: 26562120.

Pergunta 3

```
1 pacman::p_load(pacman, rio, tidyverse, ggplot2, hrbrthemes) # Load req. packages
2 ## Dados -----
3 qualidade_ar <- import("~/Projects/git/personal/PEst/Code_env/QualidadeAR03.xlsx")
4 result <- qualidade_ar[, c(2,10)]
5 result$Entrecampos <- as.numeric(qualidade_ar$Entrecampos)
6 result$VNTelha-Maia <- as.numeric(qualidade_ar$VNTelha-Maia)
7 dataframe <- gather(result, "VNTelha-Maia", "Entrecampos", 1:2)
8 ## Histograma -----
9 ggplot(dataframe, aes(x = `Entrecampos`, fill = `VNTelha-Maia`)) +
10   geom_histogram(position = "identity", alpha = 0.2, bins = 50) +
11   labs(fill="Estações", x="Níveis de ozono (µg/m³)", y="Ocorrências") +
12   ggtitle("Qualidade do ar em diferentes estações da QUALAR") +
13   theme(plot.title = element_text(hjust = 0.5, size=32, face="bold"),
14         axis.title=element_text(size=24),
15         axis.text=element_text(size=16),
16         legend.title=element_text(size=24),
17         legend.text=element_text(size=16))
18
19 ## Valor médio da qualidade do ar nas estações -----
20 E_avg = mean(result$Entrecampos); print(E_avg)
21 V_avg = mean(result$VNTelha-Maia); print(V_avg)
```



Comentário: A partir do histograma formulado através das observações horárias de níveis de ozono em microgramas por metro cúbico em duas estações específicas - Entrecampos e VNTelha-Maia - no ano de 2020, é possível obter a conclusão de que a estação VNTelha-Maia teve mais horas com níveis baixos de ozono (0 a 50 $\mu\text{g}/\text{m}^3$) em comparação à estação Entrecampos. No entanto, para níveis intermédios de ozono (50 a 100 $\mu\text{g}/\text{m}^3$), suscita-se o contrário, i.e., a estação Entrecampos tem um nível predominante de ocorrências. A níveis elevados de ozono (100 a 200 $\mu\text{g}/\text{m}^3$), ambas as estações apresentam um número de ocorrências bastante diminuto.

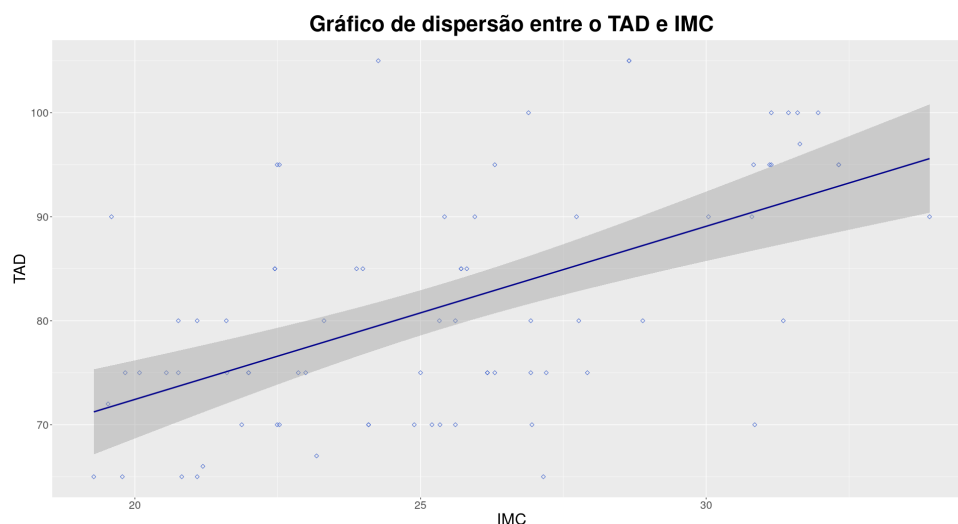
Com o observado, seria de esperar que a estação Entrecampos fosse a mais poluída. Tal é corroborado, ao analisar os níveis de poluição médias das duas estações ao longo do ano, visto que Entrecampos apresenta um valor de 52.7 $\mu\text{g}/\text{m}^3$ e VNTelha-Maia um valor de 50.9 $\mu\text{g}/\text{m}^3$.

Pergunta 4

```

1 pacman::p_load(pacman, rio, tidyverse, ggplot2, hrbrthemes) # Load req. packages
2 ## Dados -----
3 utentes <- import("~/Projects/git/personal/PEst/Code_env/Utentes.xlsx")
4 result <- utentes[, c(3,4)]
5 # Covariância e coeficiente de correlação linear
6 IMC=result$IMC
7 TAD=result$TAD
8 cov=cov(IMC, TAD); sprintf("Covariância: %.2f", cov)
9 corr=cor(IMC, TAD); sprintf("Coeficiente de correlação linear: %.2f", corr)
10 ## Gráfico de dispersão -----
11 ggplot(result, aes(x=IMC, y=TAD)) +
12   geom_point(size=2, shape=23, color = 'royalblue3') +
13   ggtitle("Gráfico de dispersão entre o TAD e IMC") +
14   theme(plot.title=element_text(hjust=0.5, size=32, face="bold"),
15         axis.title=element_text(size=24),
16         axis.text=element_text(size=16))+
17   geom_smooth(method=lm, color = 'darkblue')

```



Comentário: Note-se que apesar da aproximação linear, há uma fraca correlação entre o IMC e o TAD, tal pode ser observado graficamente através dos diversos pontos dos dados recolhidos que não seguem este padrão. É então possível concluir que apesar de se associar² que um aumento do IMC leva a um aumento do TAD, os diversos pontos dispersos, provenientes dos dados recolhidos, suscitam uma relação que não é perfeitamente linear - sendo este o motivo pelo qual foi considerado que a correlação entre as duas variáveis é fraca/não ideal, por análise gráfica. A associação positiva entre os dados é amplamente estudada¹, e portanto, esperada.

Covariância $[s_{xy}]$	Coeficiente de correlação linear $[r_{xy}]$
24.94	0.57

Assim, analisando a covariância e do coeficiente de correlação linear, é facilmente corroborado o visualizado, e expectável entre o binómio de conjuntos de dados. Visto que $s_{xy} > 0$, é natural a associação positiva observada (e documentada¹); no entanto, verifica-se $r_{xy} \approx 0.57$, que ainda se encontra longe do valor ideal unitário.

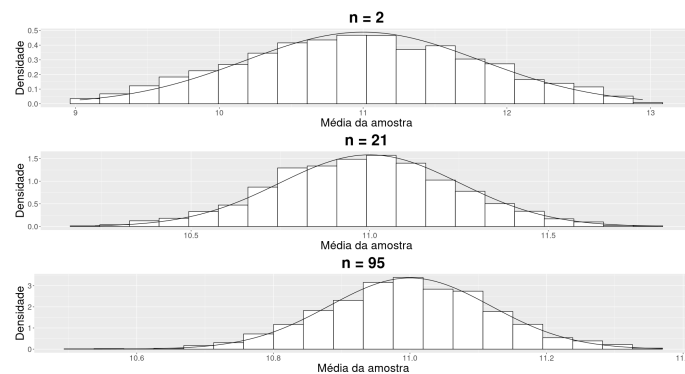
²Dua S, Bhuker M, Sharma P, Dhall M, Kapoor S. Body mass index relates to blood pressure among adults. N Am J Med Sci. 2014 Feb;6(2):89-95. doi: 10.4103/1947-2714.127751. PMID: 24696830; PMCID: PMC3968571.

Pergunta 6

```

1  pacman::p_load(pacman, rio, cowplot, tidyverse, datasets, ggplot2, dplyr,
2                    data.table, hrbrthemes) # Load req. packages
3  options(digits=15)
4  amostras=1940
5  a=9; b=13
6  E=(a+b)/2 # Valor teórico para o valor esperado
7  V=((b-a)^2)/12 # Valor teórico para a variância
8  ## Simulação -----
9  sim <- function(n){
10    set.seed(731)
11    Vn=V/n
12    data=rep(0, amostras)
13    data_avg=data.frame(rep(0, amostras))
14    avg=rep(0, amostras)
15    for(i in 1:amostras){
16      data[i] = data.frame(X=runif(n, a, b))
17      avg[i] = mean(data[[i]])
18    }
19    data_avg <- data.frame(avg)
20    # Plot do gráfico
21    plot <- ggplot(data_avg, aes(x = avg)) +
22      geom_histogram(aes(y = ..density..), colour = 1, fill = "white", bins=20) +
23      labs(x="Média da amostra", y="Densidade") +
24      stat_function(fun = dnorm, args = list(mean = E, sd = sqrt(Vn))) +
25      theme(plot.title=element_text(hjust=0.5, size=32, face="bold"),
26            axis.title=element_text(size=24),
27            axis.text=element_text(size=16)) + ggtitle(paste0("n = ", n))
28    return(plot)
29  }
30  ## Histogramas sobrepostos com as curvas com a distribuição normal -----
31  plot_grid(sim(2), sim(21), sim(95), nrow=3, ncol=1)

```



Parâmetros: Seed: 731; Número de amostras: 1940; Uniforme no intervalo: [9, 13]

Comentário: Com base nas múltiplas simulações, onde foi calculada a média de uma distribuição uniforme (para três valores da dimensão de amostras), foi elaborado o histograma de densidades relativas apresentado, ao qual foi sobreposto uma curva com distribuição normal.

Visto que as variáveis aleatórias são independentes e identicamente distribuídas, estamos perante uma aplicação possível do Teorema do Limite Central, e deste modo é aceitável aproximar a distribuição uniforme a uma normal de valor esperado $E(\mathbb{X}) = \mu$ e variância $V_n(\mathbb{X}) = \sigma^2/n$. Salienta-se a aproximação vai melhorando com o aumento da dimensão das amostras, no entanto prova-se aceitável para os valores baixos requisitados (inferiores a 30).

Ao aumentar a dimensão de amostras é aparente a tendência de um estreitamento e maior semelhança com a curva com distribuição normal. Este comportamento é necessariamente esperado, visto que ocorre uma menor dispersão da média à medida que a dimensão de amostras aumenta - a variância é inversamente proporcional a este valor: $V_n(\mathbb{X}) = \sigma^2/n$.

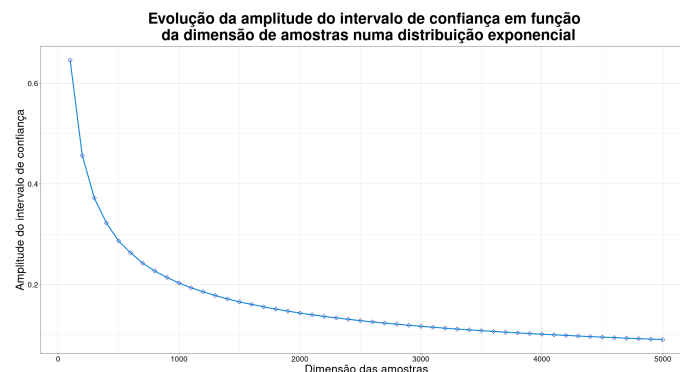
Então para amostras de maiores dimensões, a média das amostras da distribuição uniforme, aproxima-se, de facto, a uma distribuição normal centrada no valor esperado de cada amostra e com variância que diminui com o aumento da dimensão das amostras.

Pergunta 9

```

1  pacman::p_load(pacman, rio, tidyverse, datasets, ggplot2, dplyr, data.table,
2                    hrbrthemes) # Load req. packages
3  set.seed(272)
4  amostras=850
5  size=5000/100
6  lambda=1.77
7  options(digits=15)
8  ## Extremo da normal de lambda approx. -----
9  gama=0.93
10 a=qnorm((1+gama)/2)
11 q=(1-gama)/2
12
13 MA <- rep(0, size)
14 vec <- rep(0, size)
15 avec <- rep(0, amostras)
16
17 for(m in 1:size){
18   n=m*100
19   mle <- rep(0, amostras)
20   # Simular a situação amostras vezes
21   for (i in 1:amostras){
22     mle[i] <- mean(rexp(n, lambda)) # Valor esperado (1/lambda)
23     avec[i] <- (2*a)/(mle[i]*sqrt(n))
24   }
25   ## Amplitude
26   MA[m]= mean(avec)
27   vec[m] <- n # Array com no. de amostras utilizado
28 }
29 data = data.frame(vec, MA)
30 ## Gráfico
31 ggplot(data, aes(x=vec, y=MA)) + theme_linedraw() +
32   ggtitle("Evolução da amplitude do intervalo de confiança em função \n da dimensão de amostras numa distribuição
33     exponencial") +
34   geom_point(size=3, shape=23, color="dark blue") + geom_line(size = 1 , colour = "#007ED9" ) +
35   labs(x="Dimensão das amostras", y="Amplitude do intervalo de confiança") +
36   theme(plot.title=element_text(hjust=0.5, size=32, face="bold"),
37         axis.title=element_text(size=24),
38         axis.text=element_text(size=16))

```



Parâmetros: Semente: 272; Número de amostras (m): 850; $\lambda = 1.77$;
Nível de confiança $(1 - \alpha)$: 0.93

Comentário: Por observação direta da evolução da amplitude do intervalo de confiança em função da dimensão de amostras, conclui-se, que de facto ocorre um decaimento de amplitude à medida que a dimensão de amostras da distribuição exponencial aumenta. Tal é corroborado teoricamente, visto que seria expectável um decaimento na ordem de $1/\sqrt{n}$.

Dado que a distribuição exponencial apresenta amostras de dimensão superior a 30 e as variáveis aleatórias são independentes e identicamente distribuídas, o Teorema do Limite Central, é aplicável, e então $\bar{X} \sim N(\mu, \sigma^2/n)$. Como tal, a normal padrão é facilmente obtida: $Z = (\bar{X} - E(\bar{X})) / (\sqrt{V(\bar{X})}/n)$.

A amplitude do intervalo de confiança λ obtem-nos, então, a amplitude do intervalo: $A = \frac{2a}{\bar{x}\sqrt{n}}$, em que $a = \Phi^{-1}(\frac{1+(1-\alpha)}{2})$.

Pergunta 10

```

1 pacman::p_load(pacman, rio, tidyverse,
2   datasets, ggplot2, dplyr, data.table,
3   hrbrthemes) # Load req.
4   packages
5 options(digits=15)
6 set.seed(50)
7 amostras=1000
8 lambda1=4.57
9 lambda2=0.19
10 gama=0.98
11 rc=0.20 # Contaminados
12 size=2500/100
13 MA_p <- rep(0, size)
14 MA_c <- rep(0, size)
15 vec <- rep(0, size)
16 mlep <- rep(0, amostras)
17 mlec <- rep(0, amostras)
18 avec_p <- rep(0, amostras)
19 avec_c <- rep(0, amostras)
20 a=qnorm((1+gama)/2) # Extremo da normal de
21   lambda aproximada
22 ## Variar o valor de dimensão das amostras
23 for(m in 1:25){
24   n=m*100
25   # Simular a situação amostras vezes
26   for(i in 1:amostras){
27     data_p <- rexp(n, lambda1) #vetor não
28     contaminado
29     temp1 <- rexp(floor(n*rc), lambda2) #
30     vetor contaminado
31     #mlep é o valor esperado da data não
32     contaminada
33     mlep[i] <- mean(data_p)
34     avec_p[i] <- (2*a)/(mlep[i]*sqrt(n))
35     # Criar o vetor contaminado substituindo
36     os primeiros 20% dos dados puros
37     pelos contaminados
38     temp2 = data_p
39     for(k in 1:floor(n*rc)){
40       temp2[k] <- temp1[k]
41     }
42     data_c = temp2

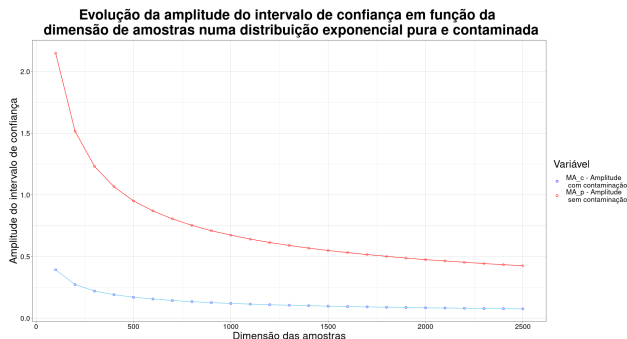
```

Continuação (...)

```

35   mlec[i] <- mean(data_c) # Valor esperado
36   da data contaminada
37   avec_c[i] <- (2*a)/(mlec[i]*sqrt(n))
38 }
39 MA_p[m]= mean(avec_p) # Amplitude da
40 amostra não contaminada
41 MA_c[m]= mean(avec_c) # Amplitude da
42 amostra não contaminada
43 vec[m] <- n # Array com no. de amostras
44 utilizado
45 }
46 data_1 = data.frame(vec, MA_p)
47 data_2 = data.frame(vec, MA_c)
48 ## Gráfico -----
49 ggplot() + theme_linedraw() +
50   ggtitle("Evolução da amplitude do intervalo
51   de confiança em função da \n dimensão
52   de amostras numa distribuição
53   exponencial pura e contaminada") +
54   geom_point(data = data_1, aes(vec, MA_p,
55   color = "MA_p - Amplitude \n sem
56   contaminação"),
57   shape=1, size=2, alpha = 1) +
58   geom_line(data=data_1, aes(
59   x=vec, y=MA_p), color="red"
60   ) +
61   geom_point(data = data_2, aes(vec, MA_c,
62   color = "MA_c - Amplitude \n com
63   contaminação"),
64   shape=1, size=2, alpha = 1) +
65   geom_line(data=data_2, aes(
66   x=vec, y=MA_c)) +
67   labs(x="Dimensão das amostras", y="
68   Amplitude do intervalo de confiança",
69   color='Variável', fill='Variável') +
70   theme(plot.title=element_text(hjust=0.5,
71   size=32, face="bold"),
72   axis.title=element_text(size=24),
73   axis.text=element_text(size=16),
74   legend.title=element_text(size=24),
75   legend.text=element_text(size=16)) +
76   scale_color_manual(values=c('blue', 'red'))

```



Parâmetros: Semente: 50;
Amostras m: 1000; $\lambda = 4.57$;
 $\lambda_C = 0.19$; $\epsilon = 20\%$;
 $(1 - \alpha) = 0.98$

Comentário: À semelhança do exercício anterior, efetuou-se a média de amplitude de vários intervalos de confiança de uma distribuição exponencial com diferentes dimensões. Porém neste caso,

compara-se uma distribuição exponencial "pura" com $\lambda = \lambda_1 = 4.57$ com uma distribuição exponencial "contaminada" com $\lambda_C = \lambda_2 = 0.19$ em que 20% foi substituído de modo a modelar a distribuição dos *outliers*.

Naturalmente, como no exercício anterior, observa-se diretamente no gráfico um decréscimo do intervalo de confiança com o aumento da dimensão das amostras (decaimento da ordem de $1/\sqrt{n}$). Na situação com amostras contaminadas por um λ_C baixo, faz-se o reparo de que a amplitude do intervalo de confiança é menor do que na situação com as amostras inalteradas ("puras").

Através do Teorema do Limite Central, é trivial concluir que a amplitude do intervalo de confiança de λ numa distribuição exponencial é dado por $A = \frac{2a}{\bar{x}\sqrt{n}}$ (com $a = \Phi^{-1}(\frac{1+(1-\alpha)}{2})$), i.e., a amplitude depende inversamente da média das amostras e da dimensão das amostras. O aumento da média das amostras também leva a uma redução na largura do intervalo de confiança, o que é igualmente observado no gráfico.