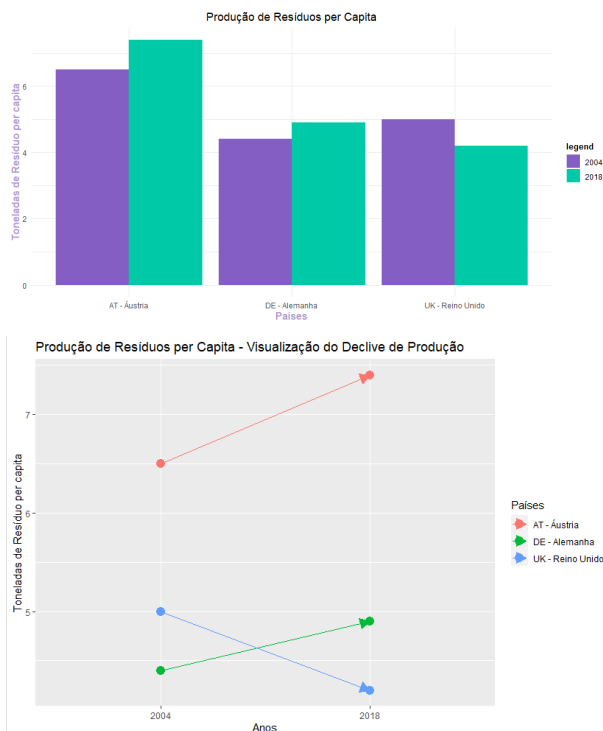


```

1 library(tidyr)
2 library(readxl)
3 library(ggplot2)
4
5 #Define a diretoria com a localização do ficheiro a ser lido
6 setwd("C:\\Users\\TERESA NOGUEIRA\\Desktop\\R_excel_files\\1_exercise")
7
8 #Obtem os dados do ficheiro excel
9 data = read_excel("ResiduosPerCapita.xlsx", range = cell_rows(c(12:43)))
10
11 #Tratamento de dados -> compacta e retira dados desnecessários, insere a matriz resultante numa
    data frame
12 data <- data[-c(1,4:30),]
13 data.df <- as.data.frame(data)
14 colnames(data.df) <- c("Países", "2004", "2018")
15
16 final.df <- gather(data.df, "Anos", "Toneladas de Resíduo per capita", 2:3)
17
18 #Gera o grafico de barras
19 ggplot(final.df, aes(x = Países, y = `Toneladas de Resíduo per capita`, fill = Anos)) +
20   geom_bar(stat = "identity", position = position_dodge()) +
21   labs(title = "Produção de Resíduos per Capita") +
22
23 #Tratamento estético do gráfico
24   scale_fill_manual("legend", values = c("2004" = "#845EC2", "2018" = "#00C9A7")) +
25   theme_minimal() +
26   theme(plot.title = element_text(hjust = 0.5),
27         axis.title.x = element_text(face="bold", colour="#B39CD0", size = 12),
28         axis.title.y = element_text(face="bold", colour="#B39CD0", size = 12),
29         legend.title = element_text(face="bold", size = 10))
30
31 #Gráfico extra para visualização do declive de Produção
32 ggplot(final.df, aes(x = Anos, y = `Toneladas de Resíduo per capita`, color = Países)) +
33   geom_point(size = 4) +
34   geom_line(aes(group=interaction(Paises)), arrow = arrow(length=unit(0.40,"cm"), type = "
    closed")) +
35   labs(title = "Produção de Resíduos per Capita - Visualização do Declive de Produção")

```

## Pergunta 1



## GRÁFICO DE BARRAS

O anterior gráfico de barras pretende evidenciar a comparação entre valores de diferentes subgrupos de dados, nomeadamente a produção de resíduos per capita de 3 países entre 2014 e 2018. Da avaliação direta do gráfico verificamos que a produção aumentou para a Áustria e para Alemanha e diminuiu para o Reino Unido.

## DECLIVE DE PRODUÇÃO

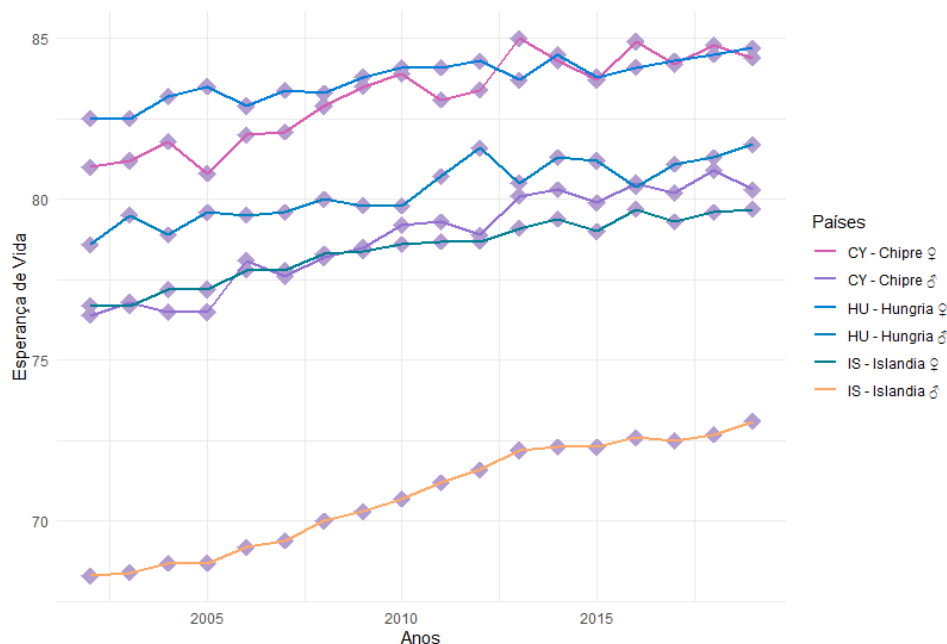
Por sua vez, o gráfico de declives ajuda apenas a tornar óbvio o que já é visualizado no anterior, sabemos que a produção em ambos os anos é mais elevada na Áustria (superior a 6 toneladas) e é onde se verifica o crescimento mais acentuado (maior declive visualizado). Os restantes países permanecem no mesmo intervalo de produção para ambos os anos, [ 0, 5] toneladas, apresentando flutuações de produção semelhantes, mas contrárias, dado que a produção de um aumenta e a do outro diminui.

```

1 library(tidyr)
2 library(readxl)
3 library(ggplot2)
4
5 #Define a diretoria com a localização do ficheiro a ser lido
6 setwd("C:\\Users\\TERESA NOGUEIRA\\Desktop\\R_excel_files\\2_exercise")
7
8 #Obtém os dados do ficheiro excel
9 data = read_excel("EsperancaVida.xlsx", range = cell_rows(c(9:70)))
10
11 #Tratamento de dados -> compacta e retira dados desnecessários, insere a matriz resultante numa
    data frame
12 data <- data[-c(1:42),-c(2:42, 44:52, 54:65, 67:76, 78:86, 88:99, 101:103)]
13
14 data.df <- as.data.frame(data)
15
16 colnames(data.df) <- c("Anos","CY - Chipre ♂","IS - Islandia ♂","HU - Hungria ♂",
17                       "CY - Chipre ♀","IS - Islandia ♀","HU - Hungria ♀")
18
19 final.df <- gather(data.df,"Países","Esperança de Vida",2:7)
20
21 #Constrói o gráfico temporal
22 ggplot(final.df, aes(x=Anos, y=`Esperança de Vida`, color = Países)) +
23   geom_point(shape = 18, colour = "#B39CD0", size = 4, stroke = 2) +
24   geom_line(size = 1) +
25   xlab("Anos")+
26   theme_minimal() +
27   scale_color_manual(values=c("#D65DB1", "#9270D3", "#007ED9",
28                             "#0082C1", "#007F93", "#FFA967"))

```

## Pergunta 2



## GRÁFICO TEMPORAL

O anterior gráfico temporal pretende analisar a evolução da esperança de vida dos homens e das mulheres da população de 3 países diferentes. A partir de uma análise meramente visual concluímos que a tendência de crescimento é positiva nos 3 países, que o género feminino tem uma expectativa de vida superior ao masculino e que a Islândia apresenta o conjunto de valores mais baixos de vida útil dos 3, especialmente para os homens, oposto a este país, evidencio o Chipre com uma expectativa de vida de 84.4 anos para as mulheres.

Posso ainda verificar que não existem outliers nem shifts súbitos nos dados (mantém a trend linear) e que não apresentam comportamentos cíclicos nem sazonais, o esperado já que a esperança de vida é maioritariamente dependente (do desenvolvimento) do país (um outro fator mais evidente, especialmente no gráfico em questão é o género).

"Where we live seems to influence how long we might live."<sup>1</sup>

<sup>1</sup><https://www.disabled-world.com/fitness/longevity/>

```

1 library(tidyr)
2 library(readxl)
3 library(ggplot2)
4
5 #Define a diretoria com a localização do ficheiro a ser lido
6 setwd("C:\\Users\\TERESA NOGUEIRA\\Desktop\\R_excel_files\\3_exercise")
7
8 #Obtém os dados do ficheiro excel
9 data = read_excel("QualidadeAR03.xlsx", range = cell_cols("C:E"), col_types = c("text", "text", "text"))
10
11 #Tratamento de dados -> compacta, retira dados desnecessários e converte a matriz para numérico
12 #, insere a matriz resultante numa data frame
13 data <- data[, -2]
14
15 data_numeric = apply(as.matrix.noquote(data), 2, as.numeric)
16
17 final.df = as.data.frame(data_numeric)
18 final.df <- gather(final.df, Estações, Value, 1:2)
19
20 #Constrói o histograma
21 ggplot(final.df, aes(x = Value, fill = Estações)) +
22   geom_histogram(color = 1, alpha = 0.65, position = "identity", bins = 40, ) +
23   scale_fill_manual(values = c("#845EC2", "#2C73D2")) +
24   xlab("Níveis de Ozono (µg/m3)") +
25   ggtitle("Qualidade do ar em diferentes estações da QUALAR") +
26   theme_minimal()
27
28 #Constrói o plot de densidade
29 ggplot(final.df, aes(x=Value, group=Estações, fill=Estações)) +
30   geom_density(adjust=1.5) +
31   scale_fill_manual(values = c("#845EC2", "#F3C5FF")) +
32   xlab("Níveis de Ozono (µg/m3)") +
33   theme_minimal()

```

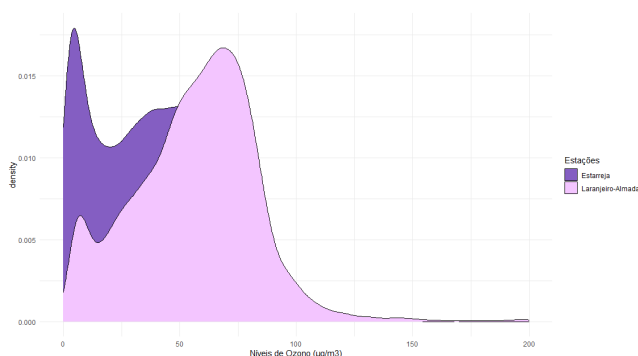
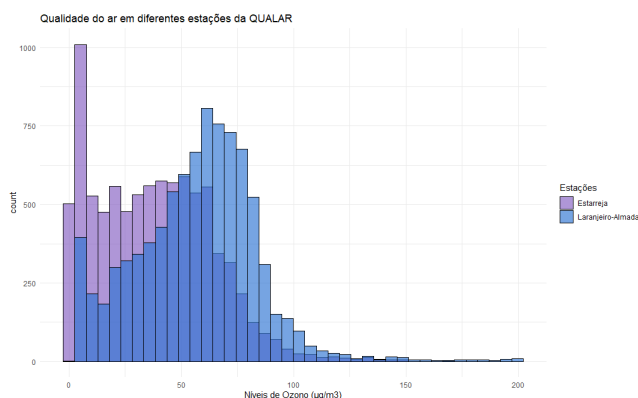
### Pergunta 3

#### HISTOGRAMA

Por análise direta do histograma relativo à qualidade de ar de duas estações de QUALAR podemos numa primeira observação admitir que ambas possuem distribuições assimétricas, o expectável, tendo em conta que a variável em observação está dependente de um conjunto de fatores relativamente versáteis (nomeadamente, o meio geográfico, o paradigma social, disposição arquitetónica do local, etc). Numa segunda observação verificamos que a mancha de concentração dos níveis de ozono para a Estarreja encontra-se especialmente no intervalo  $[0, 50]$   $\mu\text{g}/\text{m}^3$  e a de Laranjeira-Almada em  $[50, 100]$   $\mu\text{g}/\text{m}^3$ . É então de fácil afirmação dizer que a Estarreja será a menos poluída do duo. De relevância, aponto apenas o pico de concentração perto de 5  $\mu\text{g}/\text{m}^3$

#### GRÁFICO DE DENSIDADE

O respetivo gráfico de densidade corrobora a afirmação anterior e garante uma melhor visualização da mancha de concentração (através da inspeção da curva contínua das amostras), de facto, Laranjeira-Almada aparenta possuir resultados mais drásticos.

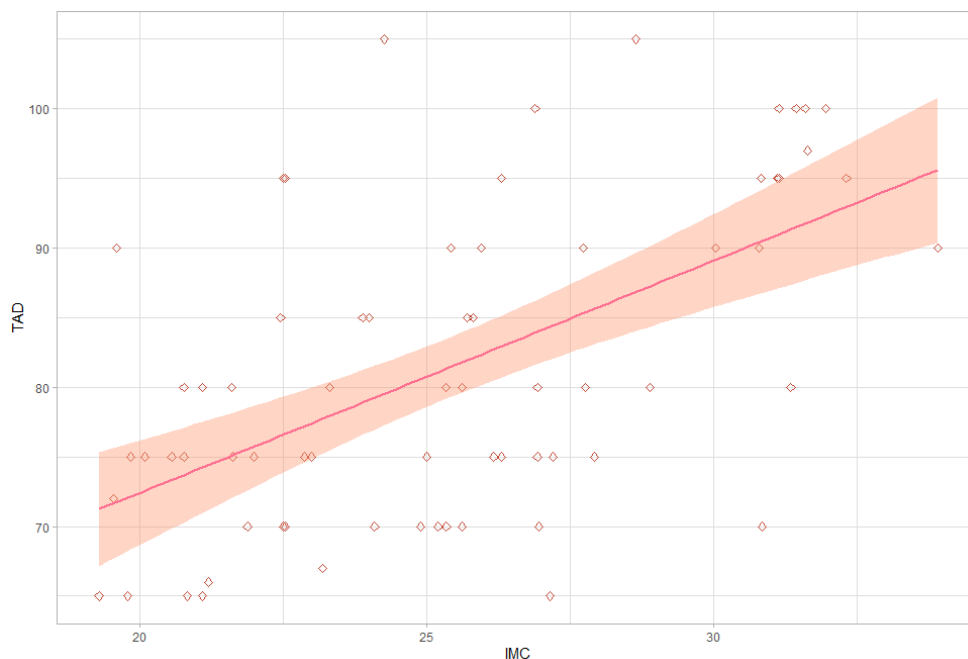


```

1 library(tidyr)
2 library(readxl)
3 library(ggplot2)
4
5 #Define a diretoria com a localização do ficheiro a ser lido
6 setwd("C:\\Users\\TERESA NOGUEIRA\\Desktop\\R_excel_files\\4_exercise")
7
8 #Obtém os dados do ficheiro excel
9 data = read_excel("Utentes.xlsx", range = cell_cols("C:D"))
10
11 final.df = as.data.frame(data)
12
13 #Cálculo da covariância e do coeficiente de correlação linear para comentários posteriores
14 IMC = final.df$IMC
15 TAD = final.df$TAD
16 cov(IMC, TAD)
17 cor(IMC, TAD)
18
19 #Constrói o Gráfico de Dispersão
20 ggplot(final.df, aes(x = IMC, y = TAD)) +
21   geom_point(size = 2, shape = 23, colour = "#C34A36") +
22   geom_smooth(formula = y ~ x, method = "lm", colour = "#FF6F91", fill="#FF9671") +
23   theme_minimal()

```

#### Pergunta 4



Covariância $[s_{xy}]$	Coeficiente de Correlação Linear $[r_{xy}]$
24.94318	0.5662385

Por análise direta do gráfico, rapidamente concluímos existir uma relação linear crescente entre o aumento do índice de massa corporal - IMC e a Tensão Arterial Diastólica - TAD. Podemos até mesmo invocar a seguinte citação, **"Body mass index (BMI) is positively associated with both systolic blood pressure (SBP) and diastolic blood pressure (DBP)."** [2018,Linderman]<sup>2</sup> para corroborar os resultados obtidos.

Por outro lado, a inspeção da covariância e do coeficiente de Correlação linear facilmente nos indicariam a relação expectável entre os dois conjuntos de dados:  $s_{xy} > 0 \rightarrow s_{xy} = 24.94318$  logo era esperada uma associação linear positiva entre conjuntos, ainda,  $r \approx 0.5662385 \approx 1$  portanto era expectável visionar uma distribuição de pontos próxima de uma resta com declive positivo, ambas estas observações são corroboradas pelo gráfico em questão.

<sup>2</sup>George C. Linderman, B. S. (2018, August 17). Association of Body Mass index with blood pressure among 1.7 million Chinese adults. JAMA Network Open. Retrieved Junho 6, 2022, de <https://jamanetwork.com/journals/jamanetworkopen/fullarticle/2696872>

## Pergunta 6

```

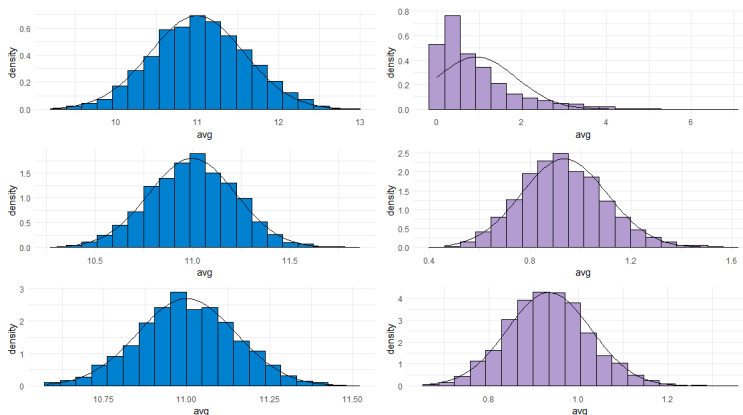
1 library(tidyr)
2 library(cowplot)
3 library(readxl)
4 library(ggplot2)
5 #—dados do problema—
6 set.seed(950)
7 a = 9
8 b = 13
9 n_samples = 1970
10 #Produz o vetor da média de amostras para cada
    n
11 vector_maker <- function(n){
12   avg = rep(0,n_samples)
13   samples = rep(0,n)
14
15   for(i in 1:n_samples){
16     samples = runif(n, a, b)
17     avg[i] = mean(samples)
18   }
19   avg.df = as.data.frame(avg)
20   return(avg.df)
21 }
22 #Produz os gráficos (o gráfico da distribuição
    exponencial não consta
23 # no código por simplicidade, é simplesmente
    uma avaliação extra)
24 plot_maker <- function(avg.df){
25   mean_of_dist = (a + b)/2
26   variance_of_dist = (b - a)^2/12

```

```

27   Vn = variance_of_dist/n
28
29   plot1 <- ggplot(avg.df, aes(x = avg)) +
30     geom_histogram(aes(y = ..density..),
31       colour = 1,
32       fill = "#0081CF", bins = 15) +
33     stat_function(fun = dnorm,
34       args = list(mean = mean_of_dist,
35         sd = sqrt(Vn))) +
36     theme_minimal()
37
38   return(plot1)
39 }
40
41 n = 4
42 avg_1.df <- vector_maker(n)
43 plotFreqDen1 <- plot_maker(avg_1.df)
44
45 n = 27
46 avg_2.df <- vector_maker(n)
47 plotFreqDen2 <- plot_maker(avg_2.df)
48
49 n = 61
50 avg_3.df <- vector_maker(n)
51 plotFreqDen3 <- plot_maker(avg_3.df)
52
53 plot_grid(plotFreqDen1, plotFreqDen2,
54   plotFreqDen3, nrow = 3, ncol = 1)

```



## COMPORTAMENTO DAS DISTRIBUIÇÕES UNIFORME E EXPONENCIAL

Na lateral encontram-se 6 gráficos correspondentes à distribuição da média de duas populações, uma de distribuição normal e outra de distribuição exponencial, para diferentes valores de  $n$  (4, 27 e 61; 1, 30 e 100, respetivamente). A observação direta e a natureza do conjunto ilícita um análise do teorema do limite Central:

Invocando o teorema do limite central, sabemos que a média de uma amostra de dados estará mais próxima da média da população geral em questão, à medida que o tamanho da amostra aumentar, independentemente da distribuição real dos dados. Afirma assim que com a extensão do tamanho da amostra, a distribuição da sua média aproxima-se de uma distribuição normal.

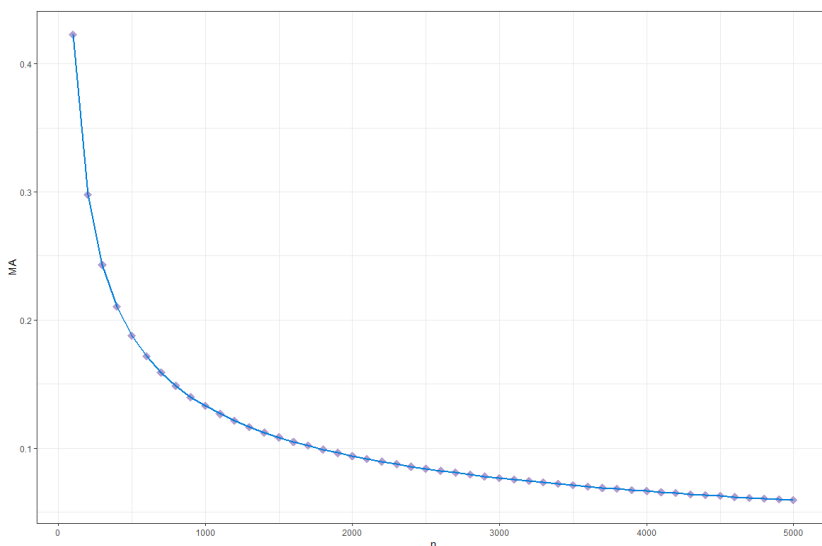
É exatamente isto que verificamos para os 6 conjuntos de amostras: para os primeiros 3, cuja população em análise é dotada de uma distribuição uniforme, é esperado verificar já a partir de  $n = 4$  uma boa sobreposição entre a curva normal teórica e as probabilidades obtidas, verificando-se apenas um afunilamento do gráfico com o aumento do tamanho da amostra (tal era esperado já que se a distribuição subjacente for simétrica, não precisamos de um tamanho de amostra muito grande para atingir a distribuição normal), basta analisar o eixo das ordenadas para verificar o aumento do concentração perto do valor esperado da média,  $E[X] = \frac{a+b}{2} = \frac{13+9}{2} = 11.0$ . Para o segundo conjunto de 3 gráfico, o mesmo fenómeno verifica-se, pesa embora apenas a partir de  $n = 30$  (a distribuição não possui simetria), onde passará a aproximar-se cada vez mais da média amostral teórica, (para  $\lambda = 1.07$ ,  $E[X] = 1/\lambda \approx 0.93$ ).

**Nota:** Esta concentração/afunilamento é expectável já que a variância é inversamente proporcional ao tamanho da amostra  $\rightarrow Var[\bar{X}] = \frac{\sigma^2}{n}$  e  $Var[\bar{X}] = \frac{1}{\lambda^2 n}$  para cada uma das respetivas populações.

Dados do problema  $\rightarrow \lambda = 1.17, m = 1350, \text{seed} = 269$

```
1 library(tidyr)
2 library(readxl)
3 library(ggplot2)
4 #---dados do problema---
5 set.seed(269)
6 m = 1350
7 lambda = 1.07
8 #-----
9 n = rep(0,50)
10 n[1] = 100
11 sample_number = 100
12 i = 2
13 #Cálcula o vetor com o tamanho das samples
14 while (sample_number < 5000){
15   n[i] = n[1] + sample_number
16   sample_number = sample_number + 100
17   i = i + 1
18 }
19 #Cálcula o vetor de intervalos para cada n e de a sua média
20 intervalo_maker <- function(n){
21   x = 0
22   for (i in 1:m){
23     x[i] = mean(rexp(n,lambda))
24   }
25   gama = 0.95
26   a = qnorm((1+gama)/2)
27   amplitude = (2*a)/(x*sqrt(n))
28   result = mean(amplitude)
29   return(result)
30 }
31 #Produz o vetor de médias da amplitude
32 z = 1
33 MA = rep(0,50)
34 for(i in 1:50){
35   MA[z] = intervalo_maker(n[i])
36   z = z + 1
37 }
38
39 data.df <- data.frame(n, MA)
40
41 ggplot(data.df, aes(x=n, y = MA)) +
42   geom_point(shape = 18, colour = "#B39CD0", size = 3, stroke = 2) +
43   geom_line(size = 1, colour = "#007ED9") +
44   theme_bw()
```

### Pergunta 9



### MÉDIA DAS AMPLITUDES DOS INTERVALOS DE CONFIANÇA

Por observação direta da progressão da amplitude de confiança em função da dimensão das amostras, admitimos que esta sofre um declínio com o aumento da extensão amostral. Tal é expectável uma vez que para obter a média das amplitudes utilizamos a seguinte fórmula:

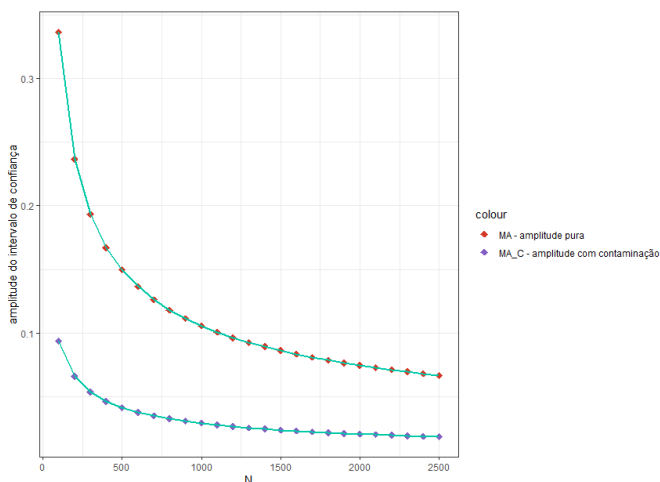
$MA = \frac{2a}{\bar{x}\sqrt{n}}$  onde  $\bar{x}$  é a média da amostra gerada,  $a = \frac{1+\gamma}{2}$  m é o tamanho do vetor final obtido de x e subsequentemente da própria amplitude.

É então esperado um decaimento na ordem dos  $\frac{1}{\sqrt{n}}$  valor este suportado pela regressão calculada.

Pergunta 10: Dados do problema  $\rightarrow \lambda = 0.92$ ,  $\lambda_c = 0.08$ ,  $m = 1050$ , seed = 269

```
1 library(tidyr)
2 library(readxl)
3 library(ggplot2)
4 #—dados do problema—
5 set.seed(269)
6 m = 1050
7 lambda = 0.92
8 lambda_c = 0.08
9 #
10 n = rep(0,25)
11 n[1] = 100
12 sample_number = 100
13
14 #declara o vetor das dimensões
15 i = 2
16 while (sample_number < 2500){
17   n[i] = n[1] + sample_number
18   sample_number = sample_number + 100
19   i = i + 1
20 }
21
22 #provoca a contaminação do vetor puro
23 cont <- function(x_pure, x_other, n){
24   for(i in 1:floor(n*0.25)){
25     x_pure[i] = x_other[i]
26   }
27   return(x_pure)
28 }
29
30 #Calcula as amplitudes dos intervalos de
  confiança
31 intervalo_maker <- function(n, x, lambda){
32   gama = 0.93
33   a = qnorm((1+gama)/2)
34   amplitude = (2*a)/(x*sqrt(n))
35   result = mean(amplitude)
36   return(result)
37 }
38
39 #Contrói os vetores de MA e MA_C
40 x_pure = rep(0,m)
```

```
41 x_other = rep(0,m)
42 x_contaminated = rep(0,m)
43 MA = rep(0,25)
44 MA_C = rep(0,25)
45 z = 1
46 for(i in 1:25){
47   for(c in 1:m){
48     x = rexp(n[i], lambda)
49     y = rexp(n[i], lambda_c)
50     j = cont(x,y, n[i])
51
52     x_pure[c] = mean(x)
53     x_contaminated[c] = mean(j)
54   }
55
56   MA[z] = intervalo_maker(n[i], x_pure, lambda)
57   MA_C[z] = intervalo_maker(n[i],
58                             x_contaminated, lambda_c)
59   z = z + 1
60 }
61 }
62
63 #Constrói o gráfico das regressões lineares
64 data.df <- data.frame(n, MA, MA_C)
65
66 ggplot(data.df) + geom_point(aes(x = n, y = MA,
67   color = "MA - amplitude pura"),
68   shape = 18, size = 2, stroke = 2) +
69   geom_point(aes(x = n, y = MA_C,
70     color = "MA_C - amplitude com contaminação"),
71     shape = 18, size = 2, stroke = 2) +
72     scale_color_manual (
73       values = c("#D43725", "#845EC2")) +
74     geom_line(aes(x = n, y = MA),
75       size = 1, colour = "#00C9A7") +
76     geom_line(aes(x = n, y = MA_C),
77       size = 1, colour = "#00C9A7") +
78     xlab("N") +
79     ylab("amplitude do intervalo de confiança")
80 + theme_bw()
```



## MÉDIA DAS AMPLITUDES DOS INTERVALOS DE CONFIANÇA COM E SEM CONTAMINAÇÃO

De forma semelhante ao exercício 9, é observado um declínio da média dos intervalos de confiança com a extensão do tamanho das amostras, nomeadamente na ordem de  $\frac{1}{\sqrt{n}}$ , já que (novamente) admitimos que

$$MA = \frac{2a}{x\sqrt{n}}.$$

De se fazer notar a discrepância entre a regressão pura e a regressão contaminada, fruto da aplicação de um  $\lambda_c$  ( $\lambda_c = 0.08 < \lambda = 0.92$  menor com contaminação da amostra inicial de 25%).

**Nota:** a regressão contaminada é menor do que a pura graças à aplicação de um *failure rate* menor, já como foi explicitado no parágrafo supra, tal deve-se ao próprio modelo da distribuição:  $f(x) = \lambda \cdot e^{-\lambda \cdot x}$ . Para o modelo contaminado teremos uma média de amostras dopada,  $E[\bar{X}] = \frac{1}{\lambda} = \frac{1}{0.75\lambda_{puro} + 0.25\lambda_c} > \frac{1}{\lambda_{puro}}$  que consequentemente provocará um decréscimo das amplitudes calculadas.