

# Planejamento e Análise de Experimentos (EEE933)

## Estudo de Caso 1

Pedro Vinícius, Samara Silva e Savio Vieira

10 de Agosto de 2020

### Introdução

Uma versão atual de um software conhecido apresenta uma distribuição dos custos de execução com média populacional  $\mu_c = 50$  e variância populacional  $\sigma_c^2 = 100$ . Posteriormente, uma nova versão do software é desenvolvida, no qual deseja-se investigar prováveis melhorias de desempenho. Com esse intuito, duas análises estatísticas são propostas: (i) uma sobre o custo médio e (ii) uma sobre a variância do custo.

Para ambos os casos, as hipóteses nulas foram definidas de maneira conservadora, partindo-se do pressuposto de que os parâmetros populacionais conhecidos foram mantidos na nova versão. A partir disso, diversas etapas foram conduzidas até a conclusão dos experimentos, entre elas a coleta de dados da distribuição dos custos do software novo, a análise exploratória dessa amostra, a inferência por meio dos testes estatísticos, a validação das premissas consideradas e as conclusões. As próximas seções contém o detalhamento técnico de cada uma dessas etapas.

### Parte 1: Teste Sobre o Custo Médio

#### Planejamento dos Experimentos

No que se refere a primeira parte do estudo de caso, o teste terá que dispor de um nível de significância  $\alpha = 0.01$ , um tamanho de efeito de mínima relevância  $\delta^* = 4$  e uma potência desejada  $\pi = 1 - \beta = 0.8$ . As hipóteses estatísticas foram definidas com o intuito de responder às questões propostas abaixo:

- Há alguma diferença entre o custo médio da versão nova do software e o custo médio da versão corrente?
- Caso haja, qual a melhor versão em termos de custo médio?

Em concordância com a proposta de comparação de custo médio entre as versões, as hipóteses de teste podem ser formuladas sobre o parâmetro média:

$$\begin{cases} H_0 : \mu_n = 50 \\ H_1 : \mu_n < 50 \end{cases}$$

onde a hipótese nula implica na igualdade entre os custos médios das versões e a hipótese alternativa unilateral na superioridade da nova versão em média.

A fase subsequente desse experimento consiste em gerar uma amostra representativa do desempenho da nova versão do software. Para tal fim, é necessário especificar o tamanho dessa amostra, considerando as propriedades preestabelecidas do teste. A priori, o Poder do Teste é bastante conveniente, porém implica em um grande dilema. O cálculo do tamanho amostral requer uma estimativa da variância, que só é obtida através das observações contidas na amostra. As possibilidades mais práticas de se conduzir o experimento nesse caso são [2]:

1. Utilização de conhecimento do problema para se obter uma estimativa (inicial) da variância;
2. Condução do estudo com um tamanho amostral predefinido, como  $N = 30$ , o que poderia violar a potência desejada;

3. Realização de um estudo piloto para estimar a variância dos dados a partir do tamanho de efeito de mínima relevância  $\delta^*$ .

Considerando as vantagens e desvantagens de cada uma, optou-se por utilizar a primeira abordagem. Por mais que essa estimativa seja sobre-estimada e os prováveis ganhos não sejam observados ao término do estudo, uma vez que se espera ganhos de variância da nova versão do software em relação à versão atual, pode-se considerar igualdade de variâncias como uma estimativa inicial, ou seja,  $\sigma_n^2 \approx \sigma_c^2 = 100$ . Entretanto, essa premissa será avaliada posteriormente na análise exploratória dos dados.

Diante da estimativa inicial da variância amostral, o Poder do Teste pode ser finalmente realizado. Esse teste é originalmente usado para mensurar o controle do teste de hipóteses sobre o erro do tipo II ( $\beta$ ), isto é,  $P(\text{rejeitar } H_0 | H_0 \text{ é falso})$ . No entanto, tal teste também pode ser utilizado para estimar outros parâmetros amostrais, como tamanho de efeito  $\delta^*$ , nível de significância  $\alpha$ , tamanho da amostra  $N$ , potência  $\pi$  e desvio padrão amostral  $\sigma_n$  [9]. No presente estudo, ele é utilizado para estimar o tamanho amostral  $N$ .

```
# Poder do Teste para estimar o tamanho amostral
(params <- power.t.test(delta = 4,
  sd = 10,
  sig.level = 0.01,
  power = 0.8,
  type = "one.sample",
  alternative = "one.sided"))
```

```
##
##      One-sample t test power calculation
##
##              n = 65.45847
##              delta = 4
##              sd = 10
##      sig.level = 0.01
##              power = 0.8
##      alternative = one.sided
```

```
N <- ceiling(params$n)
```

Assim, tem-se uma estimativa de  $N = 66$  observações com o arredondamento superior.

## Coleta dos Dados

O custo de execução do software novo é modelado por um algoritmo de Evolução Diferencial que faz uso de recombinação binomial, mutação aleatória e seleção elitista padrão para a minimização do *Shifted Sphere Problem*. Dessa forma, como cada observação é obtida por meio do módulo ExpDE [1], foram necessárias  $N$  execuções para construção da amostra. Os dados coletados foram salvos em um arquivo `.csv` com uma semente definida empiricamente `seed = 1007`. Uma vez que a amostra está fixa, as análises estatísticas da média e da variância poderão ser efetuadas sobre os mesmos dados.

```
data_generation <- function(n){
  mre <- list(name = "recombination_bin", cr = 0.9)
  mmu <- list(name = "mutation_rand", f = 2)
  mpo <- 100
  mse <- list(name = "selection_standard")
  mst <- list(names = "stop_maxeval", maxevals = 10000)
  mpr <- list(name = "sphere", xmin = -seq(1, 20), xmax = 20 + 5 * seq(5, 24))
```

```

sample <- c()
# Geração de N observações
for (i in 1:n){
  observation <- ExpDE(mpo, mmu, mre, mse, mst, mpr,
    showpars = list(show.its = "none"))$Fbest
  sample <- c(sample, observation)
}
return(sample)
}

```

## Análise Exploratória de Dados

Algumas primeiras propriedades da amostra, como média, mediana, valores extremos e variância, podem ser obtidas de imediato.

```

data.frame('Variância' = var(sample), 'Média' = mean(sample),
  'Mediana' = median(sample), 'Mínimo' = min(sample),
  'Máximo' = max(sample), 'Desvio' = sd(sample))

```

```

## Variância Média Mediana Mínimo Máximo Desvio
## 1 36.6344 49.63844 48.92703 38.76112 67.21031 6.052636

```

É importante ressaltar que a variância é significativamente menor que a considerada no cálculo do tamanho amostral  $N$ , assim, caso as premissas consideradas sejam válidas, a potência do teste será superior a desejada ( $\pi > 0.8$ ).

Entretanto, a fim de compreender melhor os dados em estudo, algumas representações gráficas da amostra coletada serão investigadas. Inicialmente, faz-se o uso do histograma, que representa a distribuição de frequências das observações. O histograma é uma boa ferramenta de análise visual do comportamento de variáveis contínuas em geral, auxiliando no processo de visualização de características, como a média da amostra e a variação da distribuição em torno da mesma. Para a amostra em questão, o histograma mostra uma distribuição levemente inclinada à direita, que corrobora que a média é maior do que a moda, bem como sugere a presença de outliers com custos maiores que 60 (Figura 1).

```

histogram <- ggplot(data = as.data.frame(sample), mapping = aes(x = sample))
histogram + geom_histogram(lwd = 0.3, bins = 20, color = 'black', fill = 'gray') +
  scale_x_continuous(name = 'Custo de Execução') +
  scale_y_continuous(name = 'Frequência')

```

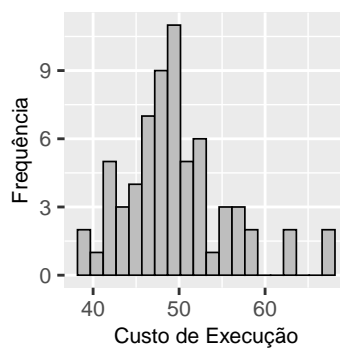


Figura 1: Histograma.

Posteriormente, o boxplot, também conhecido como gráfico de caixa, é considerado. Além de descrever simultaneamente várias propriedades importantes de um conjunto de dados, como centro, dispersão e desvio de simetria, ele também possibilita a identificação de valores atípicos [7]. Para a amostra em questão, foi possível reforçar a presença das observações em torno da mediana levemente inferior a 50 (segundo quartil). Outro aspecto importante é a presença de uma assimetria positiva na distribuição, uma vez que a mediana está mais próxima do primeiro do que do terceiro quartil. A diferença nos comprimentos das linhas que saem da caixa pressupõe novamente a cauda acentuada à direita (Figura 2).

```
boxplot <- ggplot(data = as.data.frame(sample), mapping = aes(y = sample))
boxplot + geom_boxplot(lwd = 0.3) + scale_x_continuous(name = '') +
  scale_y_continuous(name = 'Custo de Execução') +
  theme(axis.text.x = element_blank())
```

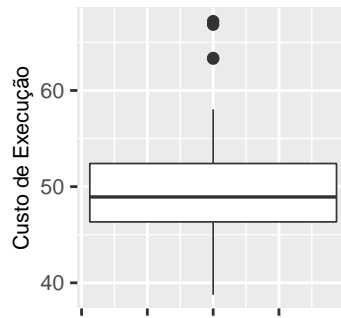


Figura 2: Boxplot.

Por fim, um gráfico bastante utilizado em estatística é o quantil-quantil (*QQ-Plot*), cujo objetivo é verificar a adequação da distribuição de frequência dos dados à uma distribuição de probabilidades de interesse. Caso os dados tenham distribuições idênticas, uma linha reta com inclinação unitária é traçada. A densidade com que os pontos são representados indicam a frequência de ocorrência dos mesmos [6]. No presente trabalho, o gráfico Q-Q foi utilizado para comparar a distribuição dos custos de execução da nova versão do software com uma distribuição normal. Ao gerá-lo, é evidente que os dados não seguem um padrão linear e, portanto, não provém de uma distribuição normal (Figura 3).

```
qqplot <- ggplot(data = as.data.frame(sample), mapping = aes(sample = sample))
qqplot + geom_qq_line() + geom_qq() + scale_y_continuous(name = 'Quantis da Amostra') +
  scale_x_continuous(name = 'Quantis Teóricos Normais')
```

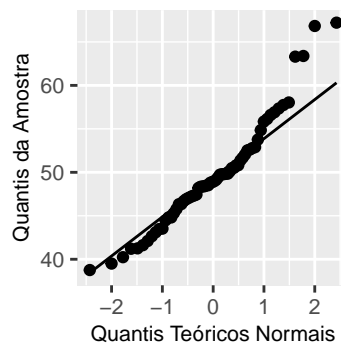


Figura 3: QQ-Plot.

## Validação de Premissas

Devido à presumível falta de normalidade dos dados, o teste de Shapiro-Wilk foi feito para confirmação. A hipótese nula  $H_0$  desse teste assume que a amostra veio de uma população com distribuição normal e para isso o parâmetro  $W$  é calculado e comparado com um valor de referência. Caso a resultante de  $W$  seja inferior ao valor tabelado, rejeita-se a hipótese de normalidade a um nível de significância  $\alpha$ . A aplicação do teste gerou como resultado  $W = 0.94772$  e o valor de  $p = 0.00755$ , que confirmam a não normalidade dos dados ( $p < 0.05$ ).

```
(shapiro_test <- shapiro.test(x = sort(sample)))
```

```
##  
## Shapiro-Wilk normality test  
##  
## data:  sort(sample)  
## W = 0.94772, p-value = 0.007555
```

Apesar disso, como o número de observações é suficientemente grande ( $N = 66 \gg 30$ ), o Teorema do Limite Central pôde ser evocado para estimação da distribuição amostral das médias. Já a premissa de independência foi assumida no processo de coleta dos dados. Caso uma amostra do software atual estivesse disponível, o teste de Durbin-Watson também caberia para avaliar uma possível autocorrelação nos resíduos.

## Análise Estatística

A partir disso, pode-se realizar o Teste t de Student sem que suas premissas sejam violadas.

```
(t_test <- t.test(x = sample,  
                 mu = 50,  
                 alternative = "less",  
                 conf.level = 0.99))
```

```
##  
## One Sample t-test  
##  
## data:  sample  
## t = -0.4853, df = 65, p-value = 0.3145  
## alternative hypothesis: true mean is less than 50  
## 99 percent confidence interval:  
##      -Inf 51.4154  
## sample estimates:  
## mean of x  
##  49.63844
```

```
cat('Intervalo de confiança:', t_test$conf.int[1:2])
```

```
## Intervalo de confiança: -Inf 51.4154
```

Como o valor  $p$  encontrado é maior que o nível de significância estabelecido ( $0.3145 > 0.05$ ), conclui-se que a hipótese nula  $H_0$  não pode ser rejeitada a um nível de confiança de 95%. No contexto do problema estudado e para as configurações de teste definidas, uma falha ao rejeitar a hipótese nula significa que não há evidências de que os desempenhos médios das duas versões do software sejam distintos e, portanto, trata-se de uma conclusão fraca. Além disso, dado que a hipótese alternativa  $H_1$  é unilateral, tem-se um intervalo de confiança aberto à esquerda  $[-\infty, 51.4154]$ .

## Potência do Teste

Conforme relatado anteriormente, como a estimativa da variância amostral foi substancialmente superior à variância calculada sobre os dados, é bem provável que a potência do teste seja superior aos 80% utilizados para o cálculo do tamanho amostral  $N$ . Dessa forma, uma estimativa da potência para detecção de diferenças maiores ou iguais ao tamanho de efeito  $\delta^*$  pode ser obtida pelo Poder do Teste, utilizando o desvio padrão amostral  $\sigma = 6.0526$ , o tamanho amostral  $N = 66$  e conservando as demais propriedades.

```
(params <- power.t.test(delta = 4,
  sd = sqrt(var(sample)),
  sig.level = 0.01,
  n = 66,
  type = "one.sample",
  alternative = "one.sided"))
```

```
##
##      One-sample t test power calculation
##
##              n = 66
##            delta = 4
##            sd = 6.052636
##          sig.level = 0.01
##            power = 0.9982989
##      alternative = one.sided
```

```
cat('Potência obtida:', params$power)
```

```
## Potência obtida: 0.9982989
```

Em concordância com o que se esperava, a potência obtida  $\pi = 0.9982$  é superior à desejada.

Outra perspectiva para esse caso é que, consequentemente, o número de observações na amostra foi sobre-estimado. Assim, uma estimativa para  $N$  capaz de garantir exatamente a potência desejada  $\pi = 0.8$  pode ser obtida pelo Poder de Teste ao utilizar o desvio padrão amostral em conjunto com os demais parâmetros.

```
(params <- power.t.test(delta = 4,
  sd = sqrt(var(sample)),
  sig.level = 0.01,
  power = 0.8,
  type = "one.sample",
  alternative = "one.sided"))
```

```
##
##      One-sample t test power calculation
##
##              n = 25.752
##            delta = 4
##            sd = 6.052636
##          sig.level = 0.01
##            power = 0.8
##      alternative = one.sided
```

```
cat('Tamanho amostral:', ceiling(params$n))
```

```
## Tamanho amostral: 26
```

Ou seja, o tamanho amostral necessário para o experimento é de apenas 26 observações, 40 observações a menos do que foi utilizado na estimativa inicial.

## Parte 2: Teste Sobre a Variância do Custo

### Planejamento dos Experimentos

Em relação à segunda parte deste experimento, quanto a avaliação da melhora da variância, foi considerado como hipótese nula  $H_0$  a asserção de que o valor de variância conhecido foi mantido e como hipótese alternativa a que este valor foi reduzido.

$$\begin{cases} H_0 : \sigma_n^2 = 100 \\ H_1 : \sigma_n^2 < 100 \end{cases}$$

Existem vários testes estatísticos para comparação de variância, como o teste F, o teste de Bonett, o teste Chi-quadrado, o teste de Levene, dentre outros. Contudo, além das particularidades de cada um deles, os três primeiros são indicados para amostras normais. Tendo em vista que o TLC não é aplicável para o cálculo da variância, a opção recomendada é o uso de um teste não-paramétrico, como o teste de Levene, já que os dados não pertencem a uma distribuição normal.

O Teste de Levene tem como hipótese nula a proposição de que os grupos que estão sendo comparados têm variâncias iguais, enquanto que a hipótese alternativa é que pelo menos um par desses grupos não possui igualdade de variância. Entretanto, essa categorização que não comporta a situação em questão, visto que o intuito é a comparação da amostra da versão atual com o valor de referência da versão nova. Nesse âmbito, adequa-se o teste não-paramétrico que receba somente uma amostra e a compare em com variância populacional dada. Sendo assim, mesmo sendo indicado para dados não paramétricos, o teste de Levene não é apropriado para o problema em questão.

Em virtude disso, e diante da escassez de métodos não-paramétricos compatível aos dados obtidos, a equipe decidiu testar outras abordagens, como a transformação de dados logarítmica e quadrática. As transformações monotônicas não lineares log 10 e raiz quadrada a mudam a distância entre os valores na distribuição e, portanto, a forma da distribuição [10]. A primeira calcula o logaritmo de um grupo de números contidos na calda à direita da distribuição, com intuito de reduzir a inclinação positiva, enquanto a segunda toma a raiz quadrada de grandes valores, levando os valores altos para mais próximo do centro. Contudo, esses dois tipos de transformação são convenientes quando os dados possuem variâncias distintas e inclinação à direita, contudo, o histograma mostra que a inclinação dos dados é à esquerda [5]. As tentativas de uso desses métodos mostraram que o comportamento dos dados não correspondia ao esperado .

Por esse motivo, o grupo prosseguiu para outra metodologia: a de reamostragem dos dados. Essa técnica consiste em realizar estimativas por meio de repetidas amostragens dentro de uma mesma amostra. A literatura apresenta alguns tipos, a saber: Teste de Aleatorização, Validação Cruzada, Jackknife e Bootstrap.

Os testes de aleatorização são testes de significância nos quais são calculados os possíveis resultados estatísticos permutando os valores amostrais, tendo como princípio a hipótese nula como verdadeira. Já na validação cruzada, divide-se a amostra entre um grupo de treinamento e outro de teste, como por exemplo no estudo de regressão, onde o primeiro é usado para cálculo dos coeficientes da equação. O Jackknife, por sua vez, é usado na estimação da variância e a tendência de algum estimador. E, finalmente, o bootstrap também é usado para os mesmos propósitos do Jackknife, porém é considerado mais abrangente, pois permite mais replicações [8]. Em vista disso, optou-se pelo uso da técnica de bootstrap.

A ideia básica do bootstrap é que a inferência sobre uma população a partir de uma amostra possa ser modelada através de nova amostragem dos dados da amostra e realizando inferência sobre uma amostra a partir de dados novamente amostrados. O bootstrap funciona tratando a inferência da verdadeira distribuição de probabilidade  $X$  sobre os dados originais, como sendo análoga à inferência da distribuição empírica de  $\hat{X}$  sobre os dados reamostrados. A precisão das inferências sobre  $\hat{X}$  usando os dados reamostrados pode ser

estimada. Se  $\hat{X}$  é uma aproximação razoável de  $X$ , então a qualidade da inferência em  $X$  pode, por sua vez, ser inferida [3].

Em nossa análise, o método de bootstrap foi concebido a partir de nossa amostra única, como função estatística de interesse a variância e o número de replicações de 10000. Uma estimativa ideal da amostra usando bootstrap leva o número de replicações ao infinito [4], em regras gerais, usa-se um número relativamente grande, limitado ao poder de processamento do computador que irá executar o cálculo. Por fim, devemos encontrar o intervalo de confiança correspondente a reamostragem bootstrap, para tal utilizamos a função `boot.ci`, que retornará o intervalo de confiança da variância, visto que é a função estatística de interesse na reamostragem.

```
# Função estatística
statistic_function <- function(x, i){
  variance <- var(x[i])
  return(variance)
}

# Bootstrap com 10000 replicações
resampling <- boot(data = sample, statistic = statistic_function, R = 10000)

# Intervalo de confiança não-paramétrico
boot.ci(boot.out = resampling, conf = 0.95, type = 'basic')

## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 10000 bootstrap replicates
##
## CALL :
## boot.ci(boot.out = resampling, conf = 0.95, type = "basic")
##
## Intervals :
## Level      Basic
## 95%      (21.61, 50.90 )
## Calculations and Intervals on Original Scale
```

Considerando um nível de confiança de 95%, obtemos um intervalo de confiança para variância de 21.49 a 50.97, o intervalo foi calculado usando o método *basic bootstrap*. Baseado nestes resultados rejeitamos a hipótese nula, de que a variância é igual a 100.

## Conclusões

Com o objetivo de inferir sobre a melhora de desempenho da nova versão de um software em comparação com a versão atual, utilizamos de conceitos e ferramentas estatísticas para investigar e concluir sobre esta comparação.

Ficamos a cargo de investigar sobre dois pontos principais, custo médio de execução e sobre a variância no custo de execução. Sobre o custo médio de execução, concluímos que por não rejeitarmos a hipótese nula de que o custo médio de execução é igual a 50, neste caso, não podemos inferir sobre a melhora da nova versão sobre a versão atual. Vale lembrar que não rejeitar a hipótese nula não significa que ela é verdadeira, e sim que pelos testes que realizamos, com 95% de confiança, não conseguimos provar sua falsidade. Porém, sobre a variância do custo da execução, rejeitamos a hipótese nula, de que a variância é igual a 100, neste caso, com os testes realizados, com 95% de confiança, a hipótese alternativa foi aceita. Assim, a partir da nossa análise, concluímos que, por esse quesito, houve melhora da nova versão do software com relação a versão atual.

## Referências

- [1] Felipe Campelo. Modular Differential Evolution for Experimenting with Operators. <https://www.rdocumentation.org/packages/ExpDE/versions/0.1.2>, 2016. Version 0.1.2.



- [2] Felipe Campelo. Lecture Notes on Design and Analysis of Experiments. <http://git.io/v3Kh8>, 2018. Version 2.12; Creative Commons BY-NC-SA 4.0.
- [3] Bradley Efron. Bootstrap methods: another look at the jackknife. In *Breakthroughs in statistics*, pages 569–593. Springer, 1992.
- [4] Bradley Efron and Robert J Tibshirani. *An introduction to the bootstrap*. CRC press, 1994.
- [5] Andy Field. *Discovering statistics using SPSS:(and sex and drugs and rock’n’roll)*. Sage, 2009.
- [6] Ramanathan Gnanadesikan and Martin B Wilk. Probability Plotting Methods for the Analysis of Data. *Biometrika*, 55(1):1–17, 1968.
- [7] Douglas C Montgomery and George C Runger. *Applied Statistics and Probability for Engineers*. John Wiley & Sons, 2010.
- [8] Camilo Daleles Renno. Jackknife, bootstrap e outros metodos de reamostragem. [http://www.dpi.inpe.br/referata/arq/\\_2011/12\\_Camilo/Renno\\_2011\\_resampl.pdf](http://www.dpi.inpe.br/referata/arq/_2011/12_Camilo/Renno_2011_resampl.pdf), 2011. Acesso em 05 de agosto de 2020.
- [9] R Development Core Team. Power Calculations For One And Two Sample T Tests. <https://www.rdocumentation.org/packages/stats/versions/3.6.2/topics/power.t.test>, 2020. Documentation reproduced from package stats, version 3.6.2, License: Part of R 3.6.2.
- [10] Sharon Lawner Weinberg and Sarah Knapp Abramowitz. *Statistics using SPSS: An integrative approach*. Cambridge University Press, 2008.