



**Universidade de Brasília**

DEPARTAMENTO DE ESTATÍSTICA

18 de janeiro de 2023

**Resolução - Marcos Augusto D. Barbosa (220006024)**

### **Lista 4: Desafio de velocidade**

Computação em Estatística para dados e cálculos massivos

Tópicos especiais em Estatística 2

Prof. Guilherme Rodrigues

César Augusto Fernandes Galvão (aluno colaborador)

Gabriel Jose dos Reis Carvalho (aluno colaborador)

1. As questões deverão ser respondidas em um único relatório *PDF* ou *html*, produzido usando as funcionalidades do *Quarto* ou outra ferramenta equivalente.
2. O aluno poderá consultar materiais relevantes disponíveis na internet, tais como livros, *blogs* e artigos.
3. O trabalho poderá ser feito individualmente ou em dupla. Suspeitas de plágio e compartilhamento de soluções serão tratadas com rigor.
4. Os códigos *R* utilizados devem ser disponibilizados na íntegra, seja no corpo do texto ou como anexo.
5. O aluno deverá enviar o trabalho até a data especificada na plataforma Microsoft Teams.
6. O trabalho será avaliado considerando o nível de qualidade do relatório, o que inclui a precisão das respostas, a pertinência das soluções encontradas, a formatação adotada, dentre outros aspectos correlatos.
7. Escreva seu código com esmero, evitando operações redundantes, visando eficiência computacional, otimizando o uso de memória, comentando os resultados e usando as melhores práticas em programação.

Simulação computacional ([https://en.wikipedia.org/wiki/Monte\\_Carlo\\_method](https://en.wikipedia.org/wiki/Monte_Carlo_method)) é uma poderosa ferramenta amplamente adotada em estudos de sistemas complexos. Aqui, para fins meramente didáticos, simularemos os resultados dos jogos da Copa do Mundo Fifa 2022, sediada no Catar, para responder questões de possível interesse prático.

Consideraremos um modelo probabilístico notavelmente rudimentar e de baixa precisão. Especificamente, assumamos que o resultado do jogo entre os times  $i$  e  $j$ , com  $i \neq j$ , segue a distribuição Poisson bivariada definida a seguir.

$$\begin{aligned}(X_i, X_j) &\sim \text{Poisson}(\lambda_{ij}, \lambda_{ji}), \quad \text{com} \\ P(X_i = x_i, X_j = x_j) &= P(X_i = x_i) P(X_j = x_j) \\ &= \frac{\lambda_{ij}^{x_i}}{x_i!} \exp(-\lambda_{ij}) \frac{\lambda_{ji}^{x_j}}{x_j!} \exp(-\lambda_{ji}),\end{aligned}$$

onde  $X_i$  e  $X_j$  representam o número de gols marcados pelas seleções  $i$  e  $j$ , respectivamente,  $P(X_i, X_j)$  denota a densidade conjunta do vetor  $(X_i, X_j)$  e  $\lambda_{ij}$  e  $\lambda_{ji}$  indicam, respectivamente, as médias (esperanças matemáticas) de  $X_i$  e  $X_j$ . Considere ainda que  $\lambda_{ij}$  é calculado, deterministicamente, como a média entre  $GF_i$  e  $GS_j$ , onde  $GF_i$  e  $GS_j$  representam, respectivamente, a média de gols feitos pelo time  $i$  nos últimos 10 jogos e a média de gols sofridos pelo time  $j$  nos últimos 10 jogos.

As estatísticas dos times classificados para o torneio estão disponíveis em <https://footystats.org/world-cup> e na pasta da tarefa no Teams. A tabela de jogos e o regulamento da Copa estão disponíveis em <https://ge.globo.com/futebol/copa-do-mundo/2022/>.

## Questão 1: Simulando a Copa do mundo

Para responder os itens a seguir, use os conhecimentos adquiridos no curso para acelerar o máximo possível os cálculos. Uma lista não exaustiva de opções inclui:

1. Usar uma lógica que evite realizar cálculos desnecessários;
2. Investigar os gargalos do código (*profiling*);
3. Criar parte do código em C++ usando o pacote Rcpp;
4. Executar as operações em paralelo usando um cluster (com múltiplas *cores*) na nuvem.

a) Sob o modelo assumido, qual era a probabilidade do Brasil vencer na estreia por 5x0? Compare o resultado exato com uma aproximação de Monte Carlo baseada em uma amostra de tamanho 1 milhão.

### Solução:

Primeiro, importamos todos os pacotes necessários.

```
library(readxl)
library(dplyr)
library(tictoc)
library(furrr)
```

Primeiro, verificamos como são os dados e, em seguida, calculamos algumas médias.

```
n=10**6
fifa_stats = read_excel("estatisticas-times.xlsx")
head(fifa_stats)
```

```
## # A tibble: 6 x 4
##   country      P    GF    GS
##   <chr>      <dbl> <dbl> <dbl>
## 1 Argentina    15    28    4
## 2 Belgium     15    32   19
## 3 Brazil      15    31    5
```

```
## 4 Cameroon      15    27    12
## 5 Canada         15    27     8
## 6 Croatia        15    29    15
```

```
fifa_stats = fifa_stats %>%
  mutate(GF_mean = GF/P, GS_mean = GS/P)
head(fifa_stats)
```

```
## # A tibble: 6 x 6
##   country      P    GF    GS GF_mean GS_mean
##   <chr>      <dbl> <dbl> <dbl>   <dbl>   <dbl>
## 1 Argentina    15    28     4    1.87    0.267
## 2 Belgium      15    32    19    2.13    1.27
## 3 Brazil        15    31     5    2.07    0.333
## 4 Cameroon     15    27    12    1.8     0.8
## 5 Canada        15    27     8    1.8    0.533
## 6 Croatia       15    29    15    1.93     1
```

Definimos a função que computa o parâmetro lambda, conforme a definição do enunciado. Depois, por meio dela, calculamos os parâmetros associados ao jogo Brasil versus Sérvia.

```
get_lambda = function(country_i, country_j, data){
  gf_mean_i = data %>% filter(country==country_i) %>% select(GF_mean) %>% pull()
  gs_mean_j = data %>% filter(country==country_j) %>% select(GS_mean) %>% pull()

  return((gf_mean_i + gs_mean_j)/2)
}

lambda_ij = get_lambda(country_i='Brazil', country_j='Serbia', data=fifa_stats)
lambda_ji = get_lambda(country_i='Serbia', country_j='Brazil', data=fifa_stats)
```

Finalmente, computamos a probabilidade de probabilidade do Brasil vencer, na estreia, por 5x0 da Sérvia, conforme (1) modelo probabilístico e (2) probabilidade empírica via simulação.

```
(prob_model = dpois(x=5, lambda=lambda_ij) * dpois(x=0, lambda=lambda_ji))
```

```
## [1] 0.004342969
```

```
score = tibble(brasil = rpois(n=n, lambda=lambda_ij),
               serbia = rpois(n=n, lambda=lambda_ji))
(prob_empirical = score %>%
  count(outcome=(brasil==5 & serbia==0)) %>%
  mutate(prob=n/sum(n)) %>%
  filter(outcome==TRUE) %>%
  pull())
```

```
## [1] 0.004404
```

Verificamos que as probabilidades são bastante próximas.

b) Qual era o jogo mais decisivo do Brasil na fase de grupos? Isso é, aquele que, se vencido, levaria à maior probabilidade de classificação da seleção para a segunda fase. Responda simulando os resultados do grupo do Brasil.

**Observação:** Esse tipo de análise é usado para definir questões comercialmente estratégicas como o calendário de competições, preço de comercialização do produto, entre outras.

## Solução:

Primeiro, armazenamos alguns dados.

```
group_G = c("Brazil", "Serbia", "Switzerland", "Cameroon")
# create all six games combinations
games = combn(x=group_G, m=2) %>%
  t() %>%
  as_tibble() %>%
  setNames(c('home', 'guest'))

lambdas = games %>%
  left_join(fifa_stats, by=join_by(home==country)) %>%
  left_join(fifa_stats, by=join_by(guest==country), suffix=c('_i', '_j')) %>%
  mutate(lambda_ij = (GF_mean_i + GS_mean_j)/2,
         lambda_ji = (GF_mean_j + GS_mean_i)/2)

board = lambdas %>% select(home, guest, lambda_ij, lambda_ji)
n_games = 6
```

Em seguida, definimos a função que gera apenas uma simulação que será iterada, posteriormente.

```
generate_simulation = function(i){

  # generate scored goals according to home team lambda
  board$GF_i = rpois(n=n_games, lambda=lambdas$lambda_ij)
  # generate scored goals according to visitor team lambda
  board$GF_j = rpois(n=n_games, lambda=lambdas$lambda_ji)

  # tag results
  board = board %>%
    mutate(result_i = case_when(GF_i > GF_j ~ 'VICTORY',
                                GF_i < GF_j ~ 'LOSS',
                                TRUE ~ 'DRAW'),
           result_j = case_when(GF_j > GF_i ~ 'VICTORY',
                                GF_j < GF_i ~ 'LOSS',
                                TRUE ~ 'DRAW'))

  # get outcomes according to home teams
  outcome_i = board %>%
    mutate(GS_i = GF_j) %>% select(home, GF_i, GS_i, result_i)
  # get outcomes according to guest teams
  outcome_j = board %>%
    mutate(GS_j = GF_i) %>% select(guest, GF_j, GS_j, result_j)
  # bind everything to long format
  colnames(outcome_j) = colnames(outcome_i)
  outcome = outcome_i %>% bind_rows(outcome_j)

  # compute points by result
  outcome = outcome %>%
    mutate(points = case_when(result_i == 'VICTORY' ~ 3,
                              result_i == 'DRAW' ~ 1,
                              TRUE ~ 0))

  # create classification board
  classification = outcome %>%
    group_by(home) %>%
    summarise(points=sum(points),
```

```

        GF=sum(GF_i),
        GS=sum(GS_i),
        DIFF=GF-GS) %>%
arrange(desc(points), desc(DIFF), desc(GF)) %>%
mutate(pos = 1:n())

# create final outcome for Brazil, tagging result for each rival and adding
# flag to check if Brazil classified
output =
  # slice for Brazil results
  board %>%
  slice(1:3) %>%
  select(result_i) %>%
  # convert results to wide format
  t() %>%
  # flag if Brazil classified, i.e, final position 1 or 2
  cbind(classification %>%
        filter(home == 'Brazil') %>%
        select(pos) %>% pull() <= 2) %>%
  as.data.frame() %>%
  setNames(c('serbia', 'switzerland', 'cameroon', 'is_classified')) %>%
  mutate(is_classified = as.logical(is_classified))

  rownames(output) = i
  return(output)
}

```

Apresentamos um exemplo de output da simulação.

```

# simulation example
generate_simulation(i=1)

```

```

##   serbia switzerland cameroon is_classified
## 1   DRAW      VICTORY      DRAW          TRUE

```

Finalmente, calculamos as probabilidades condicionais empíricas de o Brasil se classificar dado o time vencido.

```

tic()
plan(multisession, workers=4)
result = future_map(1:10000, ~ generate_simulation(.x))
toc()

```

```
## 376.75 sec elapsed
```

```

final = do.call(bind_rows, result)
final %>%
  filter(serbia=='VICTORY') %>% select(is_classified) %>% pull() %>% mean()

```

```
## [1] 0.8902972
```

```

final %>%
  filter(switzerland=='VICTORY') %>% select(is_classified) %>% pull() %>% mean()

```

```
## [1] 0.8449698
```

```
final %>%  
  filter(cameroon=='VICTORY') %>% select(is_classified) %>% pull() %>% mean()
```

```
## [1] 0.8805537
```

Portanto, o jogo contra Sérvia é o mais decisivo do Brasil na fase de grupos.

c) Qual era a probabilidade do Brasil ser campeão, em uma final contra a Argentina, tendo se classificado em primeiro do grupo? Para responder ao item, gere 10 milhões de amostra de Monte Carlo usando um cluster na nuvem!

**Atenção:** Nas fases eliminatórias, em caso de empate, sorteie o classificado considerando probabilidade de 50% para cada time (como dizem - equivocadamente -, *penalty* é loteria).

**Solução:**

## Considerações finais

Aqui consideramos um exemplo lúdico, mas o mesmo procedimento é útil para resolver problemas em genética, engenharia, finanças, energia, etc.

Há uma vasta literatura na área de modelagem preditiva de resultados esportivos (via modelos probabilísticos e de aprendizagem de máquina - algorítmicos). Entretanto, por não ser esse o foco do curso, optamos por não modelar o número esperado de gols marcados por equipe. Com base em resultados passados, seria possível ajustar modelos bem mais sofisticados, que levassem em consideração, por exemplo, contra quem os últimos resultados foram alcançados. Decidimos também modelar a incerteza usando distribuições Poisson independentes. Essa é obviamente uma suposição equivocada. Alternativas mais flexíveis podem ser adotadas para melhorar a capacidade preditiva do processo.