

## Análise de Correlação e Associação

Charles Guimarães Cavalcante – RM 334409

Luan Nonato Figueiredo – RM 334325

Rodrigo Rossi de Lima Cano – RM 333927

## Insights a partir a geração de modelos preditivos em bases de dados

### Base de dados 1 – Dados sobre valores de imóveis na cidade de Boston, MA, EUA

**1.1) Considerando que a variável de interesse seja o preço ou valor do imóvel, quais são as variáveis que mais explicam o comportamento do preço dos imóveis? Comente e justifique seus insights.**

# leitura dos dados

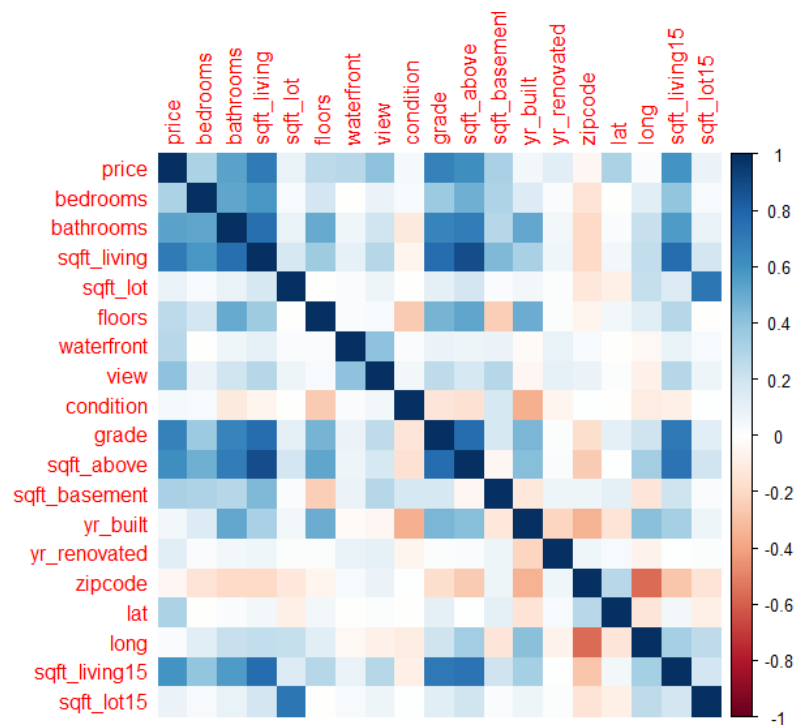
```
boston <- read.csv(file="../input/Boston_Housing_Data.csv")
```

# remoção dos campos id e data

```
boston <- boston[, c(-1,-2)]
```

# gráfico de correlação das variáveis

```
corrplot(round(cor(boston), 2), method = "color")
```



De acordo com o gráfico “price” tem boa correlação com: bedrooms, bathrooms, view, grade, sqft\_above, sqft\_basement, lat, sqft\_living, sqft\_living15.

**1.2) Existem variáveis redundantes presentes na base de dados? Se sim, quais as análises que você realizou para chegar a essa conclusão? Em caso de existência de redundância, quais foram as variáveis e como você endereçou o problema?**

Das variáveis selecionadas foram testadas agrupadas por semelhança.

**Grupo 1: bedrooms, bathrooms**

```
cor(boston$bedrooms, boston$bathrooms)
```

Resultado: 0.5158836 – correlação moderada

**Grupo 2: view, grade**

```
cor(boston$view, grade$bathrooms)
```

Resultado: 0.2513206 – correlação fraca

**Grupo 3: sqft\_living, sqft\_above, sqft\_basement, sqft\_living15**

```
cor(boston$sqft_living, boston$sqft_above)
```

Resultado: 0.8765966 – correlação muito forte

```
cor(boston$sqft_living, boston$sqft_basement)
```

Resultado: 0.435043 – correlação moderada

```
cor(boston$sqft_living, boston$sqft_living15)
```

Resultado: 0.7564203 – correlação forte

```
cor(boston$sqft_above, boston$sqft_basement)
```

Resultado: -0.05194331 – correlação muito fraca

```
cor(boston$sqft_above, boston$sqft_living15)
```

Resultado: 0.7318703 – correlação forte

```
cor(boston$sqft_basement, boston$sqft_living15)
```

Resultado: 0.200355 – correlação fraca

As variáveis **sqft\_living** e **sqft\_above** tem correlação muito forte. As variáveis

`sqft_living` e `sqft_living15`, e as variáveis `sqft_living15` e `sqft_above` tem correlação forte.

**1.3) Construa um modelo preditivo que explique o preço dos imóveis. Quais variáveis entraram no modelo final? Exiba a matriz de parâmetros e interprete os resultados. Qual o nível de acurácia do modelo? Justifique a métrica utilizada e interprete o resultado.**

Primeiro selecionamos a matrix somente com as variáveis selecionadas: `price`, `bedrooms`, `bathrooms`, `sqft_living`, `view`, `grade`, `sqft`, `basement`.

```
data = boston[,c(1, 2, 3, 4, 8, 10, 12)]
```

A seguir, separamos os dados em dados de treino (80%) e teste (20%):

```
set.seed(41) # semente para reproduzir os mesmos dados
sample = sample.split(data, SplitRatio=0.8) # separação em 80%
train_data = subset(data, sample==TRUE) # dados de treino
test_data = subset(data, sample==FALSE) # dados de teste
```

Criação do modelo de predição com os dados de treino:

```
modelo <- lm(price ~ ., data=train_data)
```

Predição e teste de acurácia:

```
predicao <- predict(newdata=test_data, modelo)
teste <- data.frame(actual=test_data$price, predicted=predicao)
media <- mean(abs(teste$actual-teste$predicted) / teste$actual)
acuracia <- 1 - media
```

Resultado: com o modelo proposto chegamos a uma acurácia para predição do preço do imóvel de **68,1%**.

## Base de dados 2 – Dados sobre rotatividade de funcionários de uma empresa

Considere a tabela `HR_Analytics.xlsx`, que traz dados de quase 15.000 empregados de uma empresa, incluindo dados sobre satisfação, desempenho, participação em projetos, registro de promoção, entre outros elementos. A ideia é buscar compreender os fatores que mais influenciam na saída do funcionário da empresa.

**2.1) Considerando que a variável de interesse seja a variável “left”, que indica se o funcionário saiu (left=1) ou não (left=0), qual metodologia ou técnica poderia ser mais adequada para entender o perfil de rotatividade dos funcionários?**

```
# leitura dos dados
```

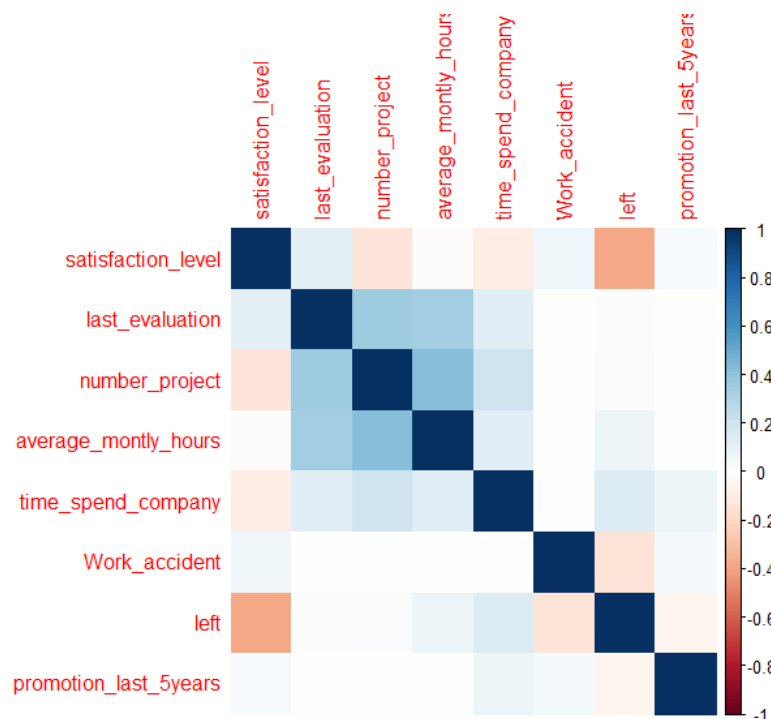
```
hr <- read_excel("HR_Analytics.xlsx")
```

## # remoção dos campos irrelevantes

```
hr <- hr[, c(-1)]
```

```
# gráfico de correlação das variáveis
```

```
corrplot(round(cor(hr[, c(-9, -10)]), 2), method = "color")
```



Podemos notar que a variável “left” tem uma relação muito negativa com o índice de satisfação (satisfacion\_level).

**2.2) Faça uma análise bivariada de cada uma das variáveis que potencialmente influenciam a saída do funcionário. Calcule o “IV” (Information Value) para cada uma das análises, interprete as análises de cada variável e ordene, por grau de importância, as variáveis que mais explicam a rotatividade dos empregados.**

Primeiro fizemos a categorização das variáveis contínuas:  
satisfaction\_level, last\_evaluation e average\_monthly\_hours.

```
satisfaction_level_cl <- discretize(hr$satisfaction_level,  
"frequency", breaks=4)  
last_evaluation_cl <- discretize(hr$last_evaluation,  
"frequency", breaks=4)  
average_monthly_hours_cl <- discretize(hr$average_monthly_hours,  
"frequency", breaks=4)  
hr <- data.frame(hr, satisfaction_level_cl, last_evaluation_cl,  
average_monthly_hours_cl)
```

Cálculo de IV (excluindo os campos que foram categorizados):

```
IV <- create_infotables(data = hr[, c(-1,-2,-4)], y = "left")  
print(head(IV$Summary, 100), row.names = FALSE)
```

Resultado:

Variable	IV
number_project	1.97240680
satisfaction_level_cl	1.07808144
time_spend_company	0.92658691
average_monthly_hours_cl	0.57291526
last_evaluation_cl	0.44550034
work_accident	0.18535538
salary	0.17904981
sales	0.03561297
promotion_last_5years	0.03385306

Vamos realizar a análise bivariada com as seguintes variáveis:

- number\_project
- satisfaction\_level\_cl
- time\_spend\_company
- average\_monthly\_hours\_cl
- last\_evaluation\_cl

```
CrossTable(hr$number_project, hr$left, prop.r=FALSE,
prop.t=FALSE, prop.chisq=FALSE)
```

hr\$number_project	hr\$left		Total
	0	1	
2	821 0.072	1567 0.439	2388
3	3983 0.349	72 0.020	4055
4	3956 0.346	409 0.115	4365
5	2149 0.188	612 0.171	2761
6	519 0.045	655 0.183	1174
7	0 0.000	256 0.072	256
Total	11428 0.762	3571 0.238	14999

A variável **number\_project** demonstra que 43,9% das pessoas que deixaram a empresa trabalhavam em apenas dois projetos.

```
CrossTable(hr$satisfaction_level_cl, hr$left, prop.r=FALSE,
prop.t=FALSE, prop.chisq=FALSE)
```

hr\$satisfaction_level_cl	hr\$left		Total
	0	1	
[0.09,0.44)	1521 0.133	2153 0.603	3674
[0.44,0.64)	3201 0.280	459 0.129	3660
[0.64,0.82)	3422 0.299	461 0.129	3883
[0.82,1]	3284 0.287	498 0.139	3782
Total	11428 0.762	3571 0.238	14999

A variável **satisfaction\_level** demonstra que 60% das pessoas que deixaram a empresa tinham índice baixo de satisfação.

```
CrossTable(hr$time_spend_company, hr$left, prop.r=FALSE,
prop.t=FALSE, prop.chisq=FALSE)
```

hr\$time_spend_company	hr\$left		Total
	0	1	
2	3191 0.279	53 0.015	3244
3	4857 0.425	1586 0.444	6443
4	1667 0.146	890 0.249	2557
5	640 0.056	833 0.233	1473
6	509 0.045	209 0.059	718
7	188 0.016	0 0.000	188
8	162 0.014	0 0.000	162
10	214 0.019	0 0.000	214
Total	11428 0.762	3571 0.238	14999

A variável **time\_spend\_company** demonstra que a maioria das pessoas que deixaram a empresa trabalhavam de 3 a 5 horas por dia.

```
CrossTable(hr$average_monthly_hours_cl, hr$left, prop.r=FALSE,
prop.t=FALSE, prop.chisq=FALSE)
```

hr\$average_monthly_hours_cl	hr\$left		Total
	0	1	
[96,156)	2354 0.206	1326 0.371	3680
[156,200)	3456 0.302	330 0.092	3786
[200,245)	3227 0.282	491 0.137	3718
[245,310]	2391 0.209	1424 0.399	3815
Total	11428 0.762	3571 0.238	14999

A variável **average\_monthly\_hours** demonstra que 37,1% das pessoas que deixam a empresa trabalhavam menos do que 156 horas por mês deixaram a empresa. Porém demonstra também que 39,9% com mais de 245 horas também deixaram a empresa.

```
CrossTable(hr$last_evaluation_cl, hr$left, prop.r=FALSE,
prop.t=FALSE, prop.chisq=FALSE)
```

hr\$last_evaluation_cl	hr\$left		Total
	0	1	
[0.36,0.56)	2259 0.198	1348 0.377	3607
[0.56,0.72)	3456 0.302	330 0.092	3786
[0.72,0.87)	3099 0.271	662 0.185	3761
[0.87,1]	2614 0.229	1231 0.345	3845
Total	11428 0.762	3571 0.238	14999

A variável **last\_evaluation** demonstra que 37,7% das pessoas que deixam a empresa tiveram índice baixo na última avaliação. Porém demonstra também que 34,5% com índice muito elevado também deixaram a empresa.

**2.3) Construa um modelo preditivo que explique a rotatividade dos empregados. Quais variáveis entraram no modelo final? Exiba a matriz de parâmetros e interprete os resultados. Qual o nível de acurácia do modelo? Justifique a métrica utilizada e interprete o resultado.**

Separação dos dados que serão utilizados:

```
data = hr[, c(3, 5, 7, 10, 11, 12)]
```

A seguir, separamos os dados em dados de treino (80%) e teste (20%):

```
set.seed(41) # semente para reproduzir os mesmos dados
sample = sample.split(data, SplitRatio=0.8) # separação em 80%
train_data = subset(data, sample==TRUE) # dados de treino
test_data = subset(data, sample==FALSE) # dados de teste
```

Criação do modelo de predição de regressão logística com os dados de treino:

```
modelo <- glm(left ~ .,
              family=binomial(link='logit'),
              data=train_data)
```



Predição e teste de acurácia:

```
predicao <- predict(modelo, newdata=test_data, type="response")  
ks.test(predicao[test_data$left==0],predicao[test_data$left==1])
```

Two-sample Kolmogorov-Smirnov test

```
data: predicao[test_data$left == 0] and predicao[test_data$left == 1]  
D = 0.59795, p-value < 0.000000000000000022  
alternative hypothesis: two-sided
```

Resultado: com o modelo proposto chegamos a uma acurácia para predição de saída do funcionário de **59,7%**.