

Estatística com Apoio Computacional

Análise Bidimensional

Universidade Estadual Vale do Acaraú – UVA

Paulo Regis Menezes Sousa

paulo_regis@uvanet.br

Análise Bidimensional

Variáveis Qualitativas

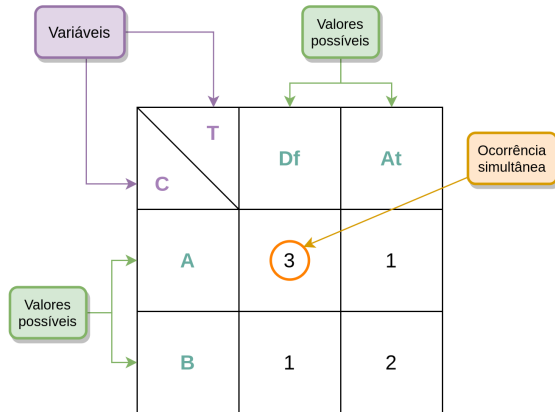
Variáveis Quantitativas

Qualitativa \times Quantitativa

- Frequentemente estamos interessados em analisar o **comportamento conjunto** de duas ou mais variáveis aleatórias.
- Os dados geralmente aparecem na forma de uma matriz, com as **colunas** indicando as variáveis e as **linhas** os indivíduos (ou elementos).
- O principal **objetivo** das análises nessa situação é explorar **relações** (similaridades) entre as colunas, ou algumas vezes entre as linhas.
- Quando consideramos **duas variáveis**, podemos ter três situações:
 - (a) duas variáveis qualitativas;
 - (b) duas variáveis quantitativas; ou
 - (c) uma variável qualitativa e outra quantitativa.

- Quando as variáveis são qualitativas, os dados são resumidos em tabelas de dupla entrada (ou de contingência).

P	C	E	T
1.0	A	Au	Df
6.2	B	Ag	At
0.4	A	Au	At
1.3	A	Pt	Df
2.9	B	Ag	Df
.	.	.	.
3.1	B	Pt	At
9.7	A	Au	Df



Exemplo 1

Suponha que queiramos analisar o comportamento conjunto das variáveis Y : grau de instrução e V : região de procedência, cujas observações estão contidas na Tabela 1. A distribuição de frequências é representada por uma tabela de dupla entrada e está na Tabela 2.

Nº	Estado civil	Grau de instrução	Nº de filhos	Salário (× sal. mín.)	Idade		Região de procedência
					anos	meses	
1	solteiro	ensino fundamental	—	4,00	26	03	interior
2	casado	ensino fundamental	1	4,56	32	10	capital
3	casado	ensino fundamental	2	5,25	36	05	capital
4	solteiro	ensino médio	—	5,73	20	10	outra
5	solteiro	ensino fundamental	—	6,26	40	07	outra
6	casado	ensino fundamental	0	6,66	28	00	interior

Tabela 1: Dados sobre os empregados da seção de orçamentos da Companhia MB.

$\begin{matrix} Y \\ \backslash \\ V \end{matrix}$	Ensino Fundamental	Ensino Médio	Superior	Total
Capital	4	5	2	11
Interior	3	7	2	12
Outra	5	6	2	13
Total	12	18	6	36

Tabela 2: Distribuição conjunta das frequências

- Para a visualização da tabela de contingência no R usamos os comandos `table()` ou `ftable()`.

```
1 emd <- read.csv("empregados-mb.csv", sep = ";", header = TRUE)
2
3 # retorna uma lista de nomes das colunas do data.frame
4 names(emd)
5
6 # cria uma tabela de contingência
7 tab.cont <- table( emd[ c("Procedência", "Grau.de.instrução") ] )
8 tab.cont
```

- Em vez de trabalharmos com as frequências absolutas, podemos construir tabelas com as frequências relativas (proporções).
- Existem três possibilidades de expressarmos a proporção, de acordo com o objetivo do problema em estudo:
 - (a) em relação ao total geral;
 - (b) em relação ao total de cada linha;
 - (c) ou em relação ao total de cada coluna.

$\begin{matrix} Y \\ V \end{matrix}$	Fundamental	Médio	Superior	Total
Capital	11%	14%	6%	31%
Interior	8%	19%	6%	33%
Outra	14%	17%	5%	36%
Total	33%	50%	17%	100%

Tabela 3: Distribuição conjunta das proporções em relação ao total geral das variáveis Y e V.

$\begin{matrix} Y \\ V \end{matrix}$	Fundamental	Médio	Superior	Total
Capital	33%	28%	33%	31%
Interior	25%	39%	33%	33%
Outra	42%	33%	34%	36%
Total	100%	100%	100%	100%

Tabela 4: Distribuição conjunta das proporções em relação aos totais de cada coluna das variáveis Y e V.

```
1 # cria uma tabela de contingência exibida em porcentagens
2 tab.prop <- prop.table(tab.cont)
3 tab.prop
4
5 # tabela de proporções das linhas
6 tab.prop.lin <- prop.table(tab.cont, margin = 1)
7 tab.prop.lin
```

- A comparação entre as duas variáveis também pode ser feita utilizando-se representações gráficas.

```
1 # cores personalizadas
2 cores <- c("#B3E2CD", "#FDCDAC", "#CBD5E8")
3
4 # retângulos com áreas são proporcionais às frequências
5 plot(tab.cont, main = "Frequências", col = cores)
6
7 # gráficos de barra com a tabela de contingência
8 barplot(tab.cont, legend = TRUE, col = cores)
9 barplot(tab.cont, legend = TRUE, col = cores, beside = TRUE)
```

```
1 # gráficos de barra com a tabela de contingência em porcentagem
2 barplot(tab.prop, legend = TRUE, col = cores)
3 barplot(tab.prop, legend = TRUE, horiz = TRUE, col = cores)
4
5 # outros atributos do gráfico de barras
6 barplot(tab.prop,
7         legend = TRUE,
8         angle = c(90,45,0),      # ângulos das linhas de sombreamento
9         density = c(30,25,20)) # densidades das linhas de sombreamento
```

Exercício 1.

Usando os dados `empregados-mb.csv`:

- a) Construa a distribuição de frequência conjunta para as variáveis grau de instrução e região de procedência.
- b) Qual a porcentagem de funcionários que têm o ensino médio?
- c) Qual a porcentagem daqueles que têm o ensino médio e são do interior?
- d) Dentre os funcionários do interior, quantos por cento têm o ensino médio?

```
1 emd <- read.csv("empregados-mb.csv", sep = ";", header = TRUE)
2
3 # a)
4 names(emd)
5 tc <- table(emd[c("Grau.de.instrução", "Procedência")])
6 tp <- prop.table(tc)
7 tp
8
9 # b)
10 percentagem.ensino.medio <- sum(tp[2, ])
11 percentagem.ensino.medio
12
13 # c)
14 tp.lin <- prop.table(tc)
15 tp.lin
16 percentagem.ensino.medio.interior <- tp.lin[2,2]
17 percentagem.ensino.medio.interior
18
19 # d)
20 tp.col <- prop.table(tc, margin = 2)
```

Exercício 2.

Numa pesquisa sobre rotatividade de mão-de-obra, para uma amostra de 40 pessoas foram observadas duas variáveis: número de empregos nos últimos dois anos (X) e salário mais recente, em número de salários mínimos (Y). Os resultados estão no arquivo `rotatividade-mao-de-obra.csv`:

- Usando a mediana, classifique os indivíduos em dois níveis, alto e baixo, para cada uma das variáveis, e construa a distribuição de frequências conjunta das duas classificações.
- Qual a porcentagem das pessoas com baixa rotatividade e ganhando pouco?
- Qual a porcentagem das pessoas que ganham pouco?
- Entre as pessoas com baixa rotatividade, qual a porcentagem das que ganham pouco?
- A informação adicional dada em d) mudou muito a porcentagem observada em c)? O que isso significa?

```
1 # a)
2 rmo = read.csv("~/EAC/rotatividade-mao-de-obra.csv", sep=";", header=T)
3 rmo$XClass = rep(NA, length(rmo$X))
4 rmo$YClass = rep(NA, length(rmo$Y))
5
6 selx = rmo$X < median(rmo$X) # TRUE para todo X < mediana
7 sely = rmo$Y < median(rmo$Y) # TRUE para todo Y < mediana
8
9 rmo$XClass[selx] = "baixo"
10 rmo$XClass[!selx] = "alto" # ! operador de negação
11 rmo$YClass[sely] = "baixo"
12 rmo$YClass[!sely] = "alto"
13
14 tc = table(rmo[, c("XClass", "YClass")]) # tab. de contingência
15 tp = prop.table(tc) # tab. de proporções
16 tp
17
18 # b)
19 tp["baixo","baixo"] # baixa rotatividade e baixo salário
```



```
21 # c)
22 sum(tp[ , "baixo"]) # baixo salário total
23
24 # d)
25 tp.lin = prop.table(tc, margin = 1) # tab. prop. para linhas
26 tp.lin
27 tp.lin["baixo", "baixo"] # baixa rotatividade com baixo salário
28
29 # e) Bastante modificada;
30 #     maioria das pessoas que ganham pouco têm alta rotatividade.
```

- Um dispositivo bastante útil para se verificar a associação entre duas variáveis quantitativas, ou entre dois conjuntos de dados, é o *gráfico de dispersão*.

```
1 # agentes de vendas e núm. de clientes por ano de trabalho
2 aac = read.csv("agentes-anos-clientes.csv", sep = ";",
3               header = TRUE)
4 head(aac)
5
6 # gráfico de dispersão
7 plot(aac[, "Número.de.clientes"], aac[, "Anos.de.serviço"],
8      xlab = "Número de clientes",
9      ylab = "Anos de serviço",
10     pch = 19,
11     col = "#FF7F00")
```

```
1 # gasto com saúde por renda bruta familiar
2 fgs = read.csv("familia-gasto-saude.csv", sep = ";",
3               header = TRUE,
4               dec = ",") # contém números separados com vírgula
5 head(fgs)
6 plot(fgs[, "X"], fgs[, "Y"],
7      xlab = "Renda bruta mensal",
8      ylab = "Renda gasta em saúde (%)",
9      pch = 17,
10     col = "#00ccff")
11
12 # nota por tempo gasto no teste
13 tom = read.csv("teste-operacao-maquina.csv", sep = ";",
14               header = TRUE)
15 head(tom)
16 plot(tom[, "X"], tom[, "Y"],
17      xlab = "Resultado no teste (0-100)",
18      ylab = "Tempo para operar a máquina (min)",
19      pch = 15,
20     col = "magenta")
```

- A partir dos gráficos apresentados, verificamos que a representação gráfica das variáveis quantitativas ajuda muito a compreender o comportamento conjunto das duas variáveis quanto à existência ou não de associação entre elas.

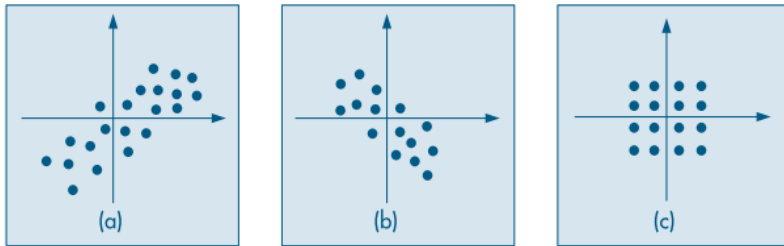
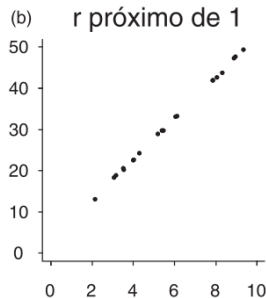
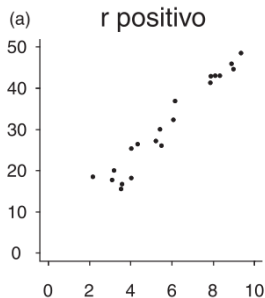


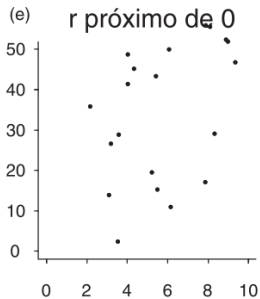
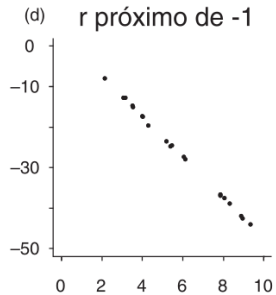
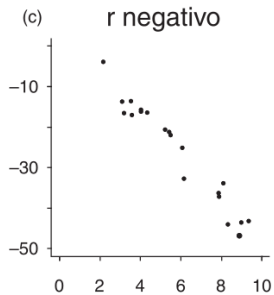
Figura 1: Tipos de associações entre duas variáveis.

Coeficiente de correlação (linear)

É uma medida do grau de associação entre elas e também da proximidade dos dados a uma reta.

- O coeficiente de correlação r_{xy} que pode assumir qualquer valor real entre -1 e 1.





- A função `cor(x,y)` calcula o coeficiente de correlação entre as variáveis x e y .

[illegible]

- É comum também analisar o que acontece com a variável **quantitativa** dentro de cada categoria da variável **qualitativa**.

Exemplo

Retomemos os dados da Tabela 1, para os quais desejamos analisar agora o comportamento dos salários (S) dentro de cada categoria de grau de instrução (Y). Começemos a análise construindo a Tabela 5, que contém medidas-resumo da variável S para cada categoria de Y .

Grau de instrução	n	\bar{s}	$dp(S)$	$var(S)$	$s_{(1)}$	q_1	q_2	q_3	$s_{(n)}$
Fundamental	12	7,84	2,79	7,77	4,00	6,01	7,13	9,16	13,65
Médio	18	11,54	3,62	13,10	5,73	8,84	10,91	14,48	19,40
Superior	6	16,48	4,11	16,89	10,53	13,65	16,74	18,38	23,30
Todos	36	11,12	4,52	20,46	4,00	7,55	10,17	14,06	23,30

Tabela 5: Medidas-resumo para a variável salário, segundo o grau de instrução, na Companhia MB.

- Criação da tabela 5 no R:
 - Na primeira linha os dados são carregados.
 - Nas linhas 3, 5 e 6, todos os salários de pessoas com ensino fundamental, médio e superior são selecionados respectivamente.

```
1 emd = read.csv("~/EAC/empregados-mb.csv", sep = ";", header = TRUE,  
2           dec = ",")  
3 fund = emd[emd[, "Grau.de.instrução"] == "Ensino fundamental", "Salário"]  
4 med  = emd[emd[, "Grau.de.instrução"] == "Ensino médio", "Salário"]  
5 sup  = emd[emd[, "Grau.de.instrução"] == "Superior", "Salário"]  
6  
7 tab = data.frame( N      = c(length(fund), length(med), length(sup)),
```

```
8 Média = c(mean(fund) , mean(med) , mean(sup) ),
9 DP    = c(sd(fund)   , sd(med)   , sd(sup)   ),
10 Var   = c(var(fund)  , var(med)  , var(sup)  ),
11 S1     = c(min(fund)  , min(med)  , min(sup)  ),
12 Q1     = c(quantile(fund, 0.25, names = FALSE),
13             quantile(med , 0.25, names = FALSE),
14             quantile(sup , 0.25, names = FALSE)),
15 Q2     = c(quantile(fund, 0.50, names = FALSE),
16             quantile(med , 0.50, names = FALSE),
17             quantile(sup , 0.50, names = FALSE)),
18 Q3     = c(quantile(fund, 0.75, names = FALSE),
19             quantile(med , 0.75, names = FALSE),
20             quantile(sup , 0.75, names = FALSE)),
21 Sn     = c(max(fund), max(med), max(sup)) )
22
23 row.names(tab) = c("Fundamental", "Médio", "Superior")
24 todos = c(fund, med, sup)
```

```
25 tab = rbind(tab, Todos = c(length(todos),
26                             mean(todos),
27                             sd(todos),
28                             var(todos),
29                             min(todos),
30                             quantile(todos, 0.25, names = FALSE),
31                             quantile(todos, 0.50, names = FALSE),
32                             quantile(todos, 0.75, names = FALSE),
33                             max(todos)) )
34
35 tab
```

- Na Figura 2, apresentamos uma visualização gráfica por meio de *box plots*.

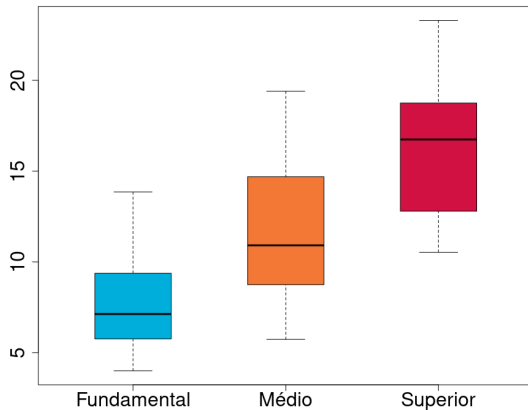


Figura 2: *Box plots* de salário segundo grau de instrução.

- Criação dos Box plots no R.

```
1 b = boxplot(fund, med, sup,  
2           names = c("Fundamental", "Médio", "Superior"),  
3           col    = c("#00aedb", "#f37735", "#d11141"),  
4           boxwex = 0.5,  
5           cex.axis = 2)  
6 b
```

- Os atributos `border` e `boxwex` são, respectivamente, um vetor com as cores das bordas dos desenhos de cada uma das caixas e a proporção de largura para todas as caixas.
- A função `boxplot` possui um retorno, neste caso o objeto `b`. O retorno da função exibe uma lista dos atributos do gráfico especificando dados que podem ser difíceis de visualizar apenas no gráfico.

A leitura desses resultados sugere uma dependência dos salários em relação ao grau de instrução: o salário aumenta conforme aumenta o nível de educação do indivíduo.

- Funcionários com o ensino fundamental completo recebem, em média, 7,84.
- O salário médio de um funcionário é 11,12 (salários mínimos).
- Para um funcionário com curso superior o salário médio passa a ser 16,48.

- Na Tabela 6 e Figura 3 temos os resultados da análise dos salários em função da região de procedência (V).

Região de procedência	n	\bar{s}	$dp(S)$	$var(S)$	$s_{(1)}$	q_1	q_2	q_3	$s_{(n)}$
Capital	11	11,46	5,22	27,27	4,56	7,49	9,77	16,63	19,40
Interior	12	11,55	5,07	25,71	4,00	7,81	10,64	14,70	23,30
Outra	13	10,45	3,02	9,13	5,73	8,74	9,80	12,79	16,22
Todos	36	11,12	4,52	20,46	4,00	7,55	10,17	14,06	23,30

Tabela 6: Medidas-resumo para a variável salário segundo a região de procedência, na Companhia MB.

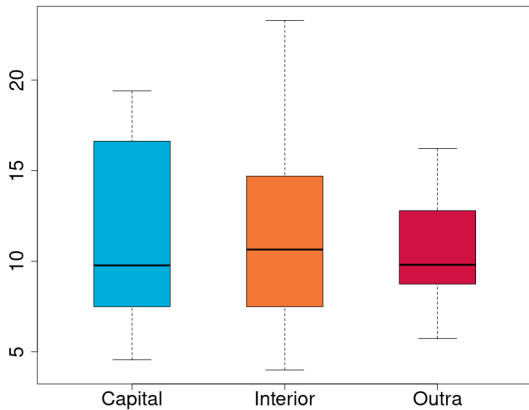


Figura 3: Medidas-resumo para a variável salário segundo a região de procedência, na Companhia MB.

- Sem usar a informação da variável categorizada, a variância calculada para a variável quantitativa para todos os dados mede a dispersão dos dados globalmente

Se a variância dentro de cada categoria for pequena e menor do que a global, significa que a variável qualitativa melhora a capacidade de previsão da quantitativa e portanto existe uma relação entre as duas variáveis.

- Para termos uma medida-resumo da variância entre as categorias da variável qualitativa, podemos usar a média das variâncias ponderada pelo número de observações em cada categoria.

$$\overline{\text{var}(S)} = \frac{\sum_{i=1}^k n_i \text{var}_i(S)}{\sum_{i=1}^k n_i}, \quad (1)$$

- No qual k é o número de categorias ($k = 3$ nos dois exemplos) e $\text{var}_i(S)$ denota a variância de S dentro da categoria $i, i = 1, 2, \dots, k$.
- Pode-se mostrar que $\overline{\text{var}(S)} \leq \text{var}(S)$, de modo que podemos definir o grau de associação entre as duas variáveis como o ganho relativo na variância, obtido pela introdução da variável qualitativa.

$$R^2 = \frac{\text{var}(S) - \overline{\text{var}(S)}}{\text{var}(S)} = 1 - \frac{\overline{\text{var}(S)}}{\text{var}(S)}. \quad (2)$$

Exemplo

Voltando aos dados do Exemplo anterior, vemos que para a variável S na presença de grau de instrução, tem-se

$$\overline{\text{var}(S)} = \frac{12(7,77) + 18(13,10) + 6(16,89)}{12 + 18 + 6} = 11,96, \quad (3)$$

$$\text{var}(S) = 20,46 \quad (4)$$

de modo que

$$R^2 = 1 - \frac{11,96}{20,46} = 0,415, \quad (5)$$

e dizemos que 41,5% da variação total do salário é explicada pela variável grau de instrução. Fazendo o mesmo cálculo para S e região de procedência temos $R^2 = 0,013$ de modo que apenas 1,3% da variabilidade dos salários é explicada pela região de procedência.

Exercício 3.

Uma amostra de 200 habitantes de uma cidade foi escolhida para declarar sua opinião sobre um certo projeto governamental. O resultado você pode encontrar no arquivo `pesquisa-opiniao.csv`:

- a) Crie uma tabela de proporções em relação ao total das colunas.
- b) Você diria que a opinião independe do local de residência?

Exercício 4.

Uma pesquisa sobre a participação em atividades esportivas de adultos moradores nas proximidades de centros esportivos construídos pelo estado de São Paulo mostrou os resultados da tabela abaixo. Baseado nesses resultados você diria que a participação em atividades esportivas depende da cidade?

Participam	Cidade			
	São Paulo	Campinas	Rib. Preto	Santos
Sim	50	65	105	120
Não	150	185	195	180

Exercício 5.

Uma amostra de dez casais e seus respectivos salários anuais (em s.m.) foi colhida num certo bairro conforme vemos na tabela abaixo.

Salário	Casal nº	1	2	3	4	5	6	7	8	9	10
	Homem (X)	10	10	10	15	15	15	15	20	20	20
	Mulher (Y)	5	10	10	5	10	10	15	10	10	15

- Encontre o salário anual médio dos homens e o seu desvio padrão.
- Encontre o salário anual médio das mulheres e o seu desvio padrão.
- Construa o diagrama de dispersão.
- Qual o salário médio familiar? E a variância do salário familiar?
- Se o homem é descontado em 8% e a mulher em 6%, qual o salário líquido médio familiar? E a variância?