

Introdução à Análise de Dados I

Lupércio F. Bessegato & Marcel T. Vieira

UFJF – Departamento de Estatística
2013



Apresentação

- Lupércio França Bessegato
lupercio.bessegato@ufjf.edu.br
Departamento de Estatística
- Marcel de Toledo Vieira
marcel.vieira@ufjf.edu.br
Departamento de Estatística

Roteiro

1. Introdução
2. Tabelas de Frequência
3. Apresentação Gráfica
4. Medidas-resumo
5. Análise Exploratória Univariada

Análise Exploratória de Dados

O que é Análise Exploratória de Dados?

- Uma filosofia/abordagem para análise de dados
- Emprega uma variedade de técnicas (a maioria gráficas)...Neste curso, trabalhamos com alguns deles:
 - √ Diagrama de dispersão
 - √ **Ramo e folhas (p/ conhecer)**
 - √ **Boxplot**
 - √ Individual Plot

Técnicas que buscam:

- maximizar o “insight” do conjunto de dados;
- perceber a estrutura subjacente;
- extrair variáveis importantes;
- detectar valores atípicos (extremos) e anomalias;
- testar hipóteses fundamentais;
- desenvolver modelos parcimoniosos; e
- determinar conjunto ótimo de fatores

ideia Básica

- Modelo = Suave + Irregular (tosco)
- Técnicas visuais podem frequentemente separar mais o “suave” do “irregular” (“ruído”)

Clássica vs. Exploratória

- Sequencia Clássica:
 - √ Problema > Dados > Modelo > Análise > Conclusões
- Exploratória:
 - √ Problema > Dados > Análise > Modelo > Conclusões

Tratamento de Dados

- Clássica:
 - √ Média e desvio padrão = estimativas pontuais
 - √ Medida de variabilidade explicada – r de Pearson
- Exploratória
 - √ Resumo Numérico (5): Min, Q1, Median, Q3, Max
 - √ todos (maioria) dados=resumos visuais
 - √ Dispersão
 - √ Histograma
 - √ Boxplot

Análise Descritiva

- Inicia-se quase sempre pela verificação dos tipos disponíveis de variáveis
- Elas podem ser resumidas por tabelas, gráficos e/ou medidas

Objetivos

- Familiarização com os dados
- Detecção de estruturas interessantes
- Presença de valores atípicos (*outliers*)

- Todos estes aspectos foram tratados neste curso!

Estudo de Caso – SAEB

Estudo de Caso - SAEB

- Dados do Sistema Nacional de Avaliação da Educação Básica – SAEB

√ Ano de 1999

√ Variáveis incluem:

- Escores de proficiência
- Condições sócio-econômicas-demográficas

- Objetivo:

√ Manipulação, análise e interpretação de dados

- Banco de dados: *Saeb99_Mestrado.xlsx/banco*

• Variáveis e suas medidas

ufesc_c	Estado onde se localiza a Escola	31 - Minas Gerais 35 - São Paulo
codalu_c	Código do Aluno	
q12_3	Que idade você tinha quando entrou na 1ª série?	1 - Menos de 5 anos 2 - 5 anos 3 - 6 anos 4 - 7 anos 5 - 8 anos 6 - 9 anos ou mais
q13_3	Qual o seu sexo?	1 - Masculino 2 - Feminino
q14_3	Como você se considera?	1 - branco(a) 2 - pardo(a) / mulato(a) 3 - negro(a) 4 - amarelo(a) 5 - indígena
q17_3	Qual o nível de instrução do seu pai?	1 - Nunca frequentou a escola 2 - Ensino Fundamental (1o Grau) - 1a à 4a série 3 - Ensino Fundamental (1o Grau) - 5a à 8a série 4 - Ensino Médio (2o Grau) 5 - Superior 6 - Pós-Graduação 7 - Não Sei

q18_3	Qual o nível de instrução de sua mãe? (mesma codificação da variável anterior)	
q111_3	Como é a sua situação econômica?	1 - Eu não trabalho e dependo economicamente da minha família 2 - Eu trabalho, mas dependo economicamente da minha família 3 - Eu trabalho e não dependo da minha família para me manter 4 - Eu trabalho e outras pessoas dependem de meu trabalho para viver
q114_3	Você já deixou de frequentar a escola por algum período?	1 - Não 2 - Sim, por 1 ano 3 - Sim, por 2 anos 4 - Sim, por 3 anos 5 - Sim, por 4 anos 6 - Sim, por mais de 4 anos
q116_3	Você já repetiu de ano?	1 - Não 2 - Sim, 1 vez 3 - Sim, 2 vezes 4 - Sim, 3 vezes 5 - Sim, 4 vezes 6 - Sim, mais de 4 vezes
q121_3	Quanto tempo você leva para chegar de sua casa até a escola?	1 - Menos de 15 minutos 2 - De 15 a 30 minutos 3 - De 31 minutos a 1 hora 4 - Mais de 1 hora
q26_3	Você gosta de Física?	1 - Não gosta 2 - Gosto mais ou menos 3 - Gosto muito
profic99	Proficiência do Aluno em Física	

Análise Exploratória – Objetivos

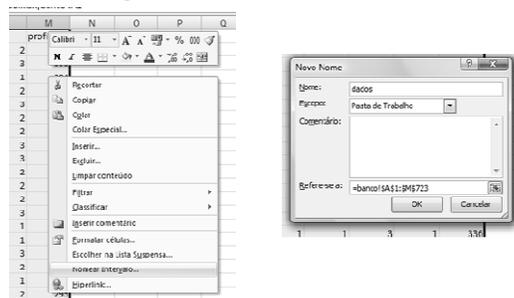
- ✓ Análise e comparação das variáveis relevantes da base de dados
(resumo por meio de gráficos, tabelas e medida-resumo)
- ✓ Identificação de estratificações importantes ao estudo de variáveis relevantes do banco de dados;
- ✓ Identificação de padrões de variáveis relevantes do banco de dados;
- ✓ Estabelecimento de hipóteses para estudo posterior mais detalhado

• 1º Passo:

- ✓ Abrir a planilha *Saeb99_Mestrado.xlsx*
- ✓ Copiar o banco de dados para uma nova planilha
- ✓ Nomear a guia com o banco com um nome especial (*banco*, por exemplo)
- ✓ Salvar a nova planilha para que você possa trabalhar nela, preservando os dados originais

• 2º Passo:

- ✓ Selecione o banco em sua planilha de trabalho
- ✓ Nomeie o intervalo selecionado (denomine-o, por exemplo, *dados*). Use o botão direito do mouse



Item 1 – Desempenho por Estado

√ Análise da proficiência de todos os aluno:

- Tendência central (média)
- Dispersão (desvio padrão)
- Assimetria (histograma)

√ Análise da proficiência por estado

- Repetir o procedimento anterior

√ Pergunta importante durante o exercício:

- Há diferença na proficiência por estado?

Bancos de Dados

Funções de Banco de Dados

- BDMÉDIA
- BDDESVP (ou BDDEST)
- BDMÁX
- BDMÍN
- BDCONTAR
- BDSOMA

Funcionamento Funções de Banco de Dados

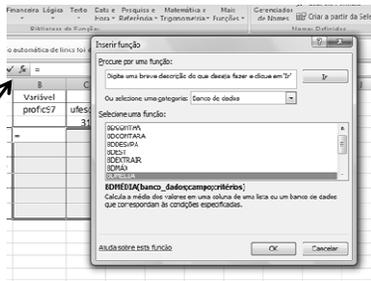
- Preparar uma tabela para cálculo de medidas:

	A	B	C	D
1		variável	critério	
2		profic97	ufasc_c	ufesc_c
3			31	35
4	Média			
5	Desvio-padrão			
6	Coefficiente de variação			
7	Mínimo			
8	Máximo			
9	Quantidade (n)			
10				

- Acesso à funções Banco de Dados:

√ Serão usadas as funções:

BDMÉDIA; BDESVPA; BDMÍN; BDMÁX; BDCONTARA



- Use o tutorial do Excel para montar as fórmulas!

- Cálculo da Média por Estado:

critério

campo

Argumentos da função

Banco_dados: dados = {"ufasc_c"; "ufesc_c"; "q11..."}
 Campo: ufasc_c = "ufasc_c"
 Critérios: \$C\$2:\$C\$3 = \$C\$2:\$C\$3
 Resultado da fórmula = 337,762
 Ajuda sobre esta função

- Procedimento similar para calcular as próximas medidas
 - √ Montar as fórmulas de uma coluna (C, por exemplo) e copiar para outra (D).
- Fórmulas da montagem da tabela

	variavel	critério	
	profici77	ufrec_c	ufrec_c
1			
2			
3			
4 Média	=BOMEDIA(dados\$B\$2:\$B\$28)	=BOMEDIA(dados\$B\$2:\$B\$3)	=BOMEDIA(dados\$B\$2:\$B\$3)
5 Desvio-padrão	=DESPAD(dados\$B\$2:\$B\$28)	=BODISP(dados\$B\$2:\$B\$3)	=BODISP(dados\$B\$2:\$B\$3)
6 Coeficiente de variação	=BIVAR(dados\$B\$2:\$B\$28)	=BIVAR(dados\$B\$2:\$B\$3)	=BIVAR(dados\$B\$2:\$B\$3)
7 Mínimo	=BOMIN(dados\$B\$2:\$B\$28)	=BOMIN(dados\$B\$2:\$B\$3)	=BOMIN(dados\$B\$2:\$B\$3)
8 Máximo	=BOMAX(dados\$B\$2:\$B\$28)	=BOMAX(dados\$B\$2:\$B\$3)	=BOMAX(dados\$B\$2:\$B\$3)
9 Quantidade (n)	=BCONTAR(dados\$B\$2:\$B\$28)	=BCONTAR(dados\$B\$2:\$B\$3)	=BCONTAR(dados\$B\$2:\$B\$3)
10			
11			

• Resultados:

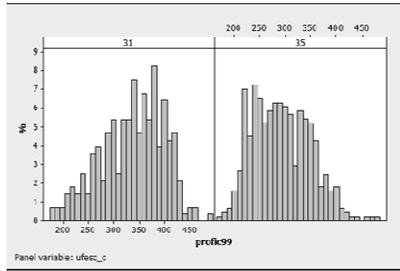
	Variável	Critério	
		ufrec_c	ufrec_c
	profici77	31	35
Média	309,807	337,762	292,201
Desvio-padrão	62,414	61,518	56,166
Coefficiente de variação	20,146	18,214	19,222
Mínimo	168,280	177,200	168,280
Máximo	485,940	465,940	475,910
Quantidade (n)	722	279	443

- Há diferença na proficiência entre os estados? É significativa?
 - √ Em qual estado o desempenho é melhor?
 - √ Onde há maior homogeneidade?

Item 1 – Histograma

- Montagem da tabela de frequência:
 - √ (Máximo; Mínimo) = (168,28; 485,94)
 - √ Intervalo de classes: 10
 - √ Limite de classe inferior (1º): 165
 - √ Limite de classe superior (último): 495
 - √ Quantidade de intervalos: 33
- Tabelas de frequência:
 - √ *case_sae_b_sol.xlsx/tabela_a_estado*

• Uma outra visão



- √ É mais 'limpa'?
- √ Facilitou a análise?
- √ As conclusões são as mesmas?

Item 2 – Influência Gênero e Etnia

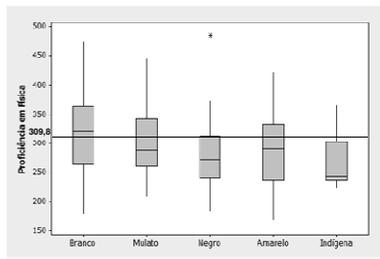
- Questão:
 - √ Há evidências empíricas que suportem a hipótese de que gênero e etnia influenciem o desempenho?
- Objetivo:
 - √ Visualizar graficamente possíveis diferenças (ou igualdades → diferença = 0)

• Tabela de medidas-resumo - Etnias:

	Global	Etnia				
		Branco	Mulato	Negro	Amarelo	Indígena
Média	309,807	315,526	302,057	281,773	288,606	267,673
Desvio-padrão	62,414	63,118	57,198	62,768	57,989	46,278
Coefficiente de variação	20,146	20,004	18,936	22,276	20,093	17,289
Mínimo	158,280	177,200	208,520	183,760	168,280	224,210
Máximo	435,040	475,010	446,330	435,040	422,050	365,440
Quantidade (n)	722	509	145	23	36	9

- √ Há evidências que indiquem diferenças de desempenho entre alunos de diferentes etnias? As diferenças são significativas?
- √ Quais suas conclusões?
- √ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

• Desempenho e Etnia



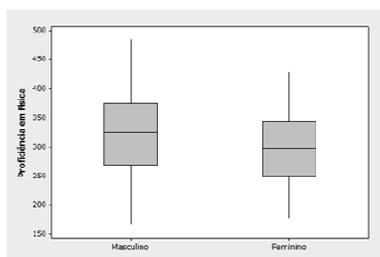
- √ Há evidências de diferenças no desempenho médio?
- √ Há evidências de diferenças na dispersão?
 - Qual seu significado?
- √ Há outliers? Quais?

• Tabela de medidas-resumo - Sexo:

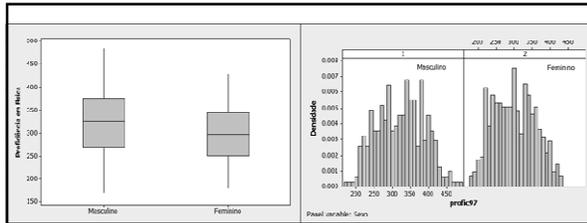
	Global	Sexo	
		Masculino	Feminino
Média	309,807	322,492	300,316
Desvio-padrão	67,414	65,767	58,536
Coefficiente de variação	20,146	20,238	19,425
Mínimo	168,280	168,280	177,200
Máximo	485,940	485,940	429,820
Quantidade (n)	722	309	413

- √ Há evidências que indiquem diferenças de desempenho entre alunos do sexo masculino e feminino? Elas são significativas?
- √ Quais suas conclusões?
- √ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

• Desempenho e Gênero



- √ Há evidências de diferenças no desempenho médio?
- √ Há evidências de diferenças na dispersão?
 - Qual seu significado?
- √ Há outliers? Quais?

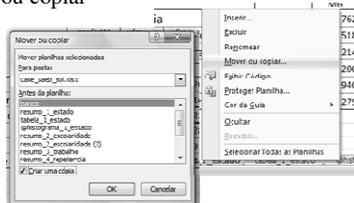


- Neste caso, o histograma facilitou sua análise?
 ✓ Suas conclusões são as mesmas?

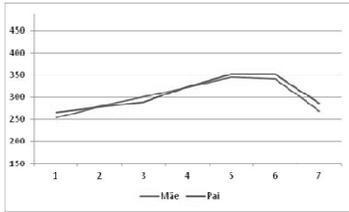
Item 3 – Influência Escolaridade Pais

- Questão:
 ✓ O nível de escolaridade dos pais ou das mães influenciam mais no desempenho? Ou é indiferente?
- Objetivo:
 ✓ Utilizar médias (e desvio-padrão) para observar evidências de possíveis diferenças (ou indiferente → diferença = 0)

- Montar uma tabela para cálculo de:
 ✓ média
 ✓ desvio padrão
 ✓ coeficiente de variação
- Faça uma cópia da guia *resumo_1_estado*:
 ✓ Com o cursor na guia da planilha que deseja copiar, clique em Mover ou copiar
 ✓ Marque copiar

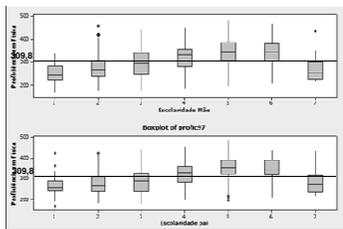


- Gráfico das médias de desempenho por nível de escolaridade (pai/mãe)



- √ Facilitou sua análise?
- √ Suas conclusões são as mesmas?
- √ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

- Box-plot do desempenho por nível de escolaridade pais



- √ Há algum novo aspecto que você observa?
- √ Suas conclusões são as mesmas?
- √ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

Item 4 – Influência do Trabalho

- Questão:
 - √ Há evidências na influência do trabalho do aluno em seu desempenho?
- Objetivo:
 - √ Visualizar possível diferença por meio de medidas resumo (ou graficamente)
 - igualdade desempenho \rightarrow diferença = 0 \rightarrow não influencia

- Montar uma tabela para cálculo de:
 - √ média
 - √ desvio padrão
 - √ coeficiente de variação
- Faça uma cópia da guia *resumo_1_estado*:
 1. Mude o critério de estratificação

	A	B	C	D	E
1		Variável		Critério	
2		profic:99	q111_3	q111_3	
3			1	>1	
4		Global			

1

- Tabela de medidas-resumo:

	Global	Trabalho	
		Não	Sim
Média	309,807	331,473	277,717
Desvio-padrão	62,414	60,801	49,513
Coeficiente de variação	20,146	18,345	17,829
Mínimo	168,280	183,760	168,280
Máximo	485,940	485,940	431,460
Quantidade (n)	722	431	291

- √ Considerado situação econômica, há diferença significativa nos desempenhos?
- √ Quais suas conclusões?
- √ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

- Box-plot do desempenho por situação trabalhista

1. Não trabalha
2. Trabalha e depende família
3. Trabalha e não depende família
4. Trabalha e família depende

- √ Há algum novo aspecto que você observa?
- √ Suas conclusões são as mesmas?
- √ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

Item 5 – Influência Repetência

- Questão:
 - √ Há evidências da influência da repetência no desempenho do aluno em Física?
- Objetivo:
 - √ Visualizar possível diferença por meio de medidas resumo (ou graficamente)
 - igualdade desempenho → diferença = 0 → não influencia

- Montar uma tabela para cálculo de:

- √ média
- √ desvio padrão
- √ coeficiente de variação

- Faça uma cópia da guia *resumo_1_estado*:

1. Mude o critério de estratificação

	Variável	Critério	
		q116_3	q116_3
	profic99	1	>1
	Global	Não	Sim
Média	309,807	329,476	279,119
Desvio-padrão	62,414	66,930	51,205

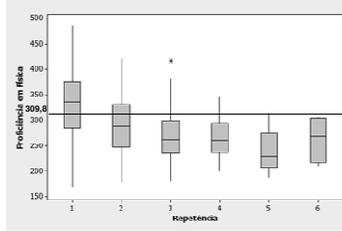
- Tabela de medidas-resumo:

	Global	Repetência	
		Não	Sim
Média	309,807	329,476	279,119
Desvio-padrão	62,414	66,930	51,205
Coeficiente de variação	20,146	18,493	18,345
Mínimo	168,280	158,280	177,200
Máximo	485,940	485,940	420,830
Quantidade (n)	722	440	282

- √ Considerada repetências no passado, há evidências de influência no desempenho? Ela é significativa?
- √ Quais suas conclusões?
- √ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

• **Box-plot do desempenho por situação trabalhista**

1. Não
2. Sim, 1 vez
3. Sim, 2 vezes
4. Sim, 3 vezes
5. Sim, 4 vezes
6. Sim, mais de 4 vezes



- ✓ Há algum novo aspecto que você observa?
- ✓ Suas conclusões são as mesmas?
- ✓ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

Item 6 – Local de Residência

- **Questão:**
 - ✓ Há evidências de que os alunos que moram mais próximos da escola tenham melhor desempenho em Física?
 - Para essa particular amostra de alunos
- **Objetivo:**
 - ✓ Visualizar possível diferença por meio de medidas resumo (ou graficamente)
 - igualdade desempenho → diferença = 0 → não influencia

- Montar uma tabela para cálculo de:
 - ✓ média
 - ✓ desvio padrão
 - ✓ coeficiente de variação
- Faça uma cópia da guia *resumo_1_estado*:
 1. Mude o critério de estratificação
 2. Copie o grupo de colunas 1 vez e substitua valor pelos níveis de tempo de deslocamento desejados 3 e 4

	vartierres		q121_3		q121_3	
profic:99	1	2	3	4	3	4
Global	< 15'	15' a 30'	31' a 1h	> 1h	< 15'	15' a 30'
Média	309,807	318,440	305,584	304,630	294,141	294,141

• Tabela de medidas-resumo:

	Global	Tempo de Deslocamento			
		< 15'	15' a 30'	31' a 1h	> 1h
Média	309,807	318,400	305,584	304,630	254,141
Desvio-padrão	62,414	65,428	61,538	54,498	65,719
Coefficiente de variação	20,146	20,549	20,138	17,890	22,343
Mínimo	168,280	186,930	177,200	202,880	168,280
Máximo	485,940	485,940	458,520	431,460	413,350
Quantidade (n)	722	270	296	131	25

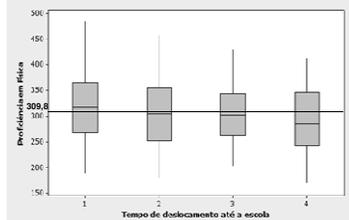
√ Considerado tempo de deslocamento, você encontra evidências que suportem sua influência no desempenho do aluno? A influência é significativa?

√ Quais suas conclusões?

√ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

• Box-plot do desempenho por tempo de deslocamento até a escola:

1. Menos de 15 minutos
2. De 15 a 30 minutos
3. De 31 minutos a 1 hora
4. Mais de 1 hora



√ Há algum novo aspecto que você observa?

√ Suas conclusões são as mesmas?

√ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

Item 7 – Influência Interesse Pessoal

• Questão:

√ Alunos que gostam de Física parecem ter um desempenho melhor do que aqueles que não gostam?

• Objetivo:

√ Visualizar possível diferença por meio de medidas resumo (ou graficamente)

- igualdade desempenho → diferença = 0 → não influencia

- Montar uma tabela para cálculo de:
 - √ média
 - √ desvio padrão
 - √ coeficiente de variação
- Faça uma cópia da guia *resumo_1_estado*:
 1. Mude o critério de estratificação
 2. Copie uma coluna substitua o valor pelo nível de interesse em Física desejado (3)

	A	B	C	D	E
1		Variável		Critério	
2		profis99	q26_3	q26_3	q26_3
3			1	2	3
4			Interesse em Física		
5		Global			
6	Média	309,807	294,037	309,739	339,528

1 2

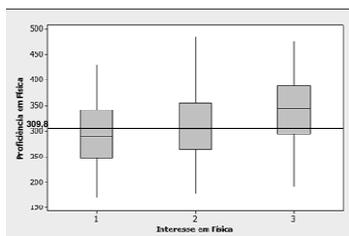
- Tabela de medidas-resumo:

	Global	Interesse em Física		
		Não	Médio	Muito
Média	309,807	254,037	309,739	339,528
Desvio-padrão	62,414	57,238	60,160	66,912
Coefficiente de variação	20,146	19,466	19,423	19,707
Mínimo	168,280	168,280	177,200	189,670
Máximo	485,940	429,820	485,940	475,910
Quantidade (n)	722	234	363	125

- √ Há evidências que indiquem diferenças nos desempenhos entre os três grupos de alunos? Elas são significativas?
- √ Quais suas conclusões?
- √ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

- Box-plot do desempenho por interesse em Física:

1. Não gosto
2. Gosto mais ou menos
3. Gosto muito



- √ Há algum novo aspecto que você observa?
- √ Suas conclusões são as mesmas?
- √ Que hipótese(s) você gostaria de testar de uma maneira mais formal?

Referências

Bibliografia

- Magalhães, M.N. e Lima, A.C.P.L. (Edusp)
Noções de Probabilidade e Estatística
- Wild, C.J. e Seber, G.A.F. (LTC)
Encontros com o Acaso: um Primeiro Curso de Análise de Dados e Inferência
- Agresti, A. e Agresti, B.F. (Dellen Pub.)
Statistical Methods for the Social Sciences
