

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

Ementa

- Introdução à análise de dados educacionais:
 - √ Variáveis e níveis de mensuração
 - √ distribuições de frequência, curva normal e escores padronizados
- Estatística descritiva:
 - √ Medidas de tendência central e de variabilidade
 - √ Representação e interpretação de dados em gráficos e tabelas

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

Questionário

Roteiro

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

Introdução

O que é *Estatística*?

- Segundo *Magalhães e Lima (2005)*, *Estatística* é um conjunto de técnicas que permite, de forma sistemática, organizar, descrever, analisar e interpretar dados oriundos de estudos ou experimentos, realizados em qualquer área do conhecimento.

Organização e Representação de Dados

- Uma das formas de organizar e resumir a informação contida em dados observados é por meio de tabela de frequências e gráficos.
- *Tabela de frequência*: relaciona categorias (ou classes) de valores, juntamente com contagem (ou frequências) do número de valores que se enquadram em cada categoria ou classe.
- *Elementos gráficos*: ajudam na visualização das principais características dos dados.
- *Medidas resumo*: Medidas de posição, dispersão, assimetria e curtose.

Variáveis

- Qualquer característica associada a um elemento pertencente a uma população ou uma amostra
- Classificação de variáveis:

Qualitativa { Nominal Sexo, cor dos olhos
Ordinal Classe social, grau de instrução

Quantitativa { Discreta Número de filhos, nº de carros
Contínua Peso, altura, salário

Dados Brutos

- Obtidos diretamente de pesquisa
 - √ Ainda sem qualquer processo de síntese ou análise
- Incluídos em tabelas
 - √ Porém, não incluídos em publicações

Atividade nº 1

Id	Turma	Sexo	Q158	Alt	Peso	Fibros	Fuma	Tslm	Elevs	One	OpOne	TV	OpTV
1	A	F	17	1,60	60,5	2	NÃO	P	0	1	B	16	R
2	A	F	18	1,69	68,0	1	NÃO	M	0	1	B	7	R
3	A	M	18	1,88	72,8	2	NÃO	P	6	2	M	16	R
4	A	M	20	1,86	69,8	2	NÃO	P	6	2	B	20	R
5	A	F	19	1,59	55,0	1	NÃO	M	2	2	B	5	R
6	A	M	19	1,76	60,0	3	NÃO	M	2	1	B	2	R
7	A	F	20	1,60	65,0	1	NÃO	P	3	1	B	7	R
8	A	F	18	1,64	47,0	1	SIM	I	2	2	M	10	R
9	A	F	18	1,62	67,8	3	NÃO	M	3	3	M	12	R
10	A	F	17	1,54	50,0	2	NÃO	M	2	2	M	10	R
11	A	F	18	1,72	70,0	1	SIM	I	10	2	B	8	N
12	A	F	18	1,00	54,0	3	NÃO	M	0	2	B	0	R
13	A	F	21	1,70	68,0	2	NÃO	M	5	1	M	30	R
14	A	M	19	1,79	69,8	1	SIM	I	6	1	M	2	N
15	A	F	18	1,66	63,8	1	NÃO	I	4	1	R	10	R
16	A	F	19	1,61	47,4	3	NÃO	B	0	1	B	18	R
17	A	F	17	1,62	60,0	1	NÃO	P	3	1	B	10	N
18	A	M	18	1,80	35,2	2	NÃO	P	3	4	B	10	R
19	A	F	20	1,60	64,5	1	NÃO	P	3	2	B	5	R
20	A	F	19	1,66	62,5	3	NÃO	M	7	2	B	14	M
21	A	F	21	1,70	60,0	2	NÃO	P	5	2	B	5	R
22	A	F	19	1,66	68,6	1	NÃO	M	0	3	B	5	R
23	A	F	19	1,57	48,2	1	SIM	I	5	4	B	10	R
24	A	F	20	1,55	49,0	1	SIM	I	0	1	M	28	R
25	A	F	20	1,69	51,5	2	NÃO	P	9	5	M	4	N
26	A	F	19	1,54	67,0	2	NÃO	I	8	2	B	8	R
27	B	F	23	1,62	63,0	2	NÃO	M	3	2	M	5	R
28	B	F	19	1,62	62,0	1	NÃO	P	1	1	M	10	R
29	B	F	19	1,57	10,0	2	NÃO	P	2	1	B	12	R

Tabelas de Frequência

- ### Tabelas de Frequências
- Uso:
 - √ Variáveis Qualitativas ou Quantitativas Discretas.
 - Contém valores da variável e suas respectivas contagens (frequências absolutas e relativas)
 - √ Frequência absoluta (n_i): contagem das ocorrências de cada valor da variável; seu total é n (o total da amostra);
 - √ Frequência relativa (f_i): proporção de ocorrência de cada valor ($f_i = n_i/n$); seu total é 1 (útil para fazer comparações entre grupos).

Tabelas de Frequência - Exemplo

Tabela de Frequências		
Sexo	Freq. Absolutas	Freq. Relativas
F	37	0,74
M	13	0,26
Total	50	1

- Classe: contém, na base de dados, quantos alunos são do sexo Masculino e quantos são do sexo Feminino.

Tabelas de Frequência para Variáveis Ordenadas

- Quando existe uma ordenação das categorias de uma variável (qualitativa ordinal ou quantitativa), faz sentido inserirmos na tabela uma outra coluna, a da frequência acumulada (f_{ac}), que é a soma das frequências relativas, do menor valor até o atual.

Exemplo: Tabela de Frequência para a Variável 'Tolerância'

Toler	Frequência	
	Absoluta (n_i)	Relativa (f_i)
M	19	38,0%
P	21	42,0%
I	10	20,0%
Total	50	100,0%

Exemplo: Tabela de Frequência para a Variável 'Nº de filhos'

Tabela de Frequências para a variável "Nº de filhos"				
Filhos	Freq. Absolutas	Freq. Acumuladas	Freq. Relativas	Freq. Acumuladas relativas
1	28	28	0,56	0,56
2	14	42	0,28	0,84
3	6	48	0,12	0,96
4	1	49	0,02	0,98
5	0	49	0	0,98
6	0	49	0	0,98
7	1	50	0,02	1
Total	50			1

- % famílias que não têm filho único?
- % famílias com pelo menos 2 filhos?
- % famílias com mais de 3 filhos?

Atividade nº 2

Nº	Estado Civil	Grau de Instrução	No de filhos	Salário (R\$ mil/mês)	Idade anos	Região de procedência
1	Solteiro	1º grau	-	4,00	26 03	Interior
2	Casado	1º grau	1	4,50	32 10	Capital
3	Casado	1º grau	2	4,25	36 05	Capital
4	Solteiro	2º grau	-	2,73	20 10	Outro
5	Solteiro	1º grau	-	6,26	40 07	Outro
6	Casado	1º grau	0	6,66	28 00	Interior
7	Solteiro	1º grau	-	6,66	41 00	Interior
8	Solteiro	1º grau	-	7,39	43 04	Capital
9	Casado	2º grau	1	4,50	32 10	Capital
10	Solteiro	2º grau	-	7,44	23 06	Outro
11	Casado	2º grau	2	8,12	23 06	Interior
12	Solteiro	1º grau	-	6,46	27 11	Capital
13	Solteiro	2º grau	-	6,74	27 05	Outro
14	Casado	2º grau	3	6,97	44 02	Outro
15	Casado	2º grau	0	8,13	30 05	Interior
16	Solteiro	2º grau	-	6,35	38 08	Outro
17	Casado	2º grau	1	9,77	31 07	Capital
18	Casado	1º grau	2	5,50	35 07	Outro
19	Solteiro	Superior	-	10,53	25 08	Interior
20	Solteiro	2º grau	-	12,76	37 04	Interior
21	Casado	2º grau	1	11,00	30 09	Outro
22	Solteiro	2º grau	-	11,49	34 02	Capital
23	Solteiro	1º grau	-	12,60	41 06	Outro
24	Casado	Superior	0	12,79	26 01	Outro
25	Casado	2º grau	2	13,23	32 05	Interior
26	Casado	2º grau	2	13,60	35 00	Outro
27	Solteiro	1º grau	-	13,85	46 07	Outro
28	Casado	2º grau	0	14,69	29 08	Interior
29	Casado	2º grau	4	14,71	40 06	Interior
30	Casado	2º grau	2	15,99	35 10	Capital
31	Solteiro	Superior	-	16,22	31 05	Outro
32	Casado	2º grau	1	16,61	36 04	Interior
33	Casado	Superior	3	17,26	43 07	Capital
34	Solteiro	Superior	-	18,70	33 07	Capital
35	Casado	2º grau	4	19,00	48 11	Capital
36	Casado	Superior	3	23,30	42 02	Interior

Apresentação Gráfica

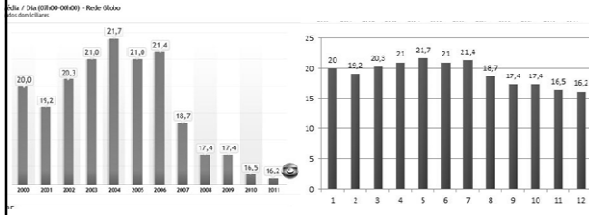
Gráficos

- Objetivo:
 - √ Identificação da forma do conjunto de dados
 - √ Resumo e identificação
 - √ Padrão dos dados
- Em geral, facilita a visualização de informações contida em tabelas
- Construção simplificada atualmente por programas computacionais

Cuidados

- Gráfico com medidas desproporcionais pode:
 - √ Dar falsa impressão de desempenho
 - √ Conduzir a conclusões equivocadas

Exemplo: Audiência



Tipos Básicos

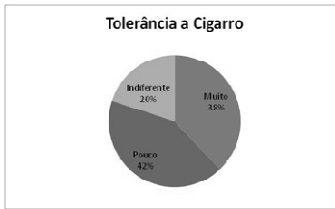
- Gráfico de setores (disco, pizza)
√ adapta-se muito bem às variáveis qualitativas nominais
- Gráfico de barras
√ adapta-se melhor às variáveis quantitativas discretas ou às variáveis qualitativas ordinais
- Histograma
√ utilizado com variáveis quantitativas contínuas

Gráfico de Setores

- Adapta-se muito bem às variáveis qualitativas nominais
- Repartição de disco em setores circulares correspondentes às frequências relativas de cada valor da variável

Exemplo: Tolerância a Cigarro

Toler	n_i	f_i
M	19	38,0%
P	21	42,0%
I	10	20,0%
Total	50	100,0%



- Importante:
 - ✓ Use com variáveis com até no máximo 6 níveis
 - ✓ Os valores não devem ser muito próximos

Gráfico de Setores – Comentários

- O gráfico de setores não é uma forma boa de visualizar informações!
 - ✓ O olho é bom para julgar medidas lineares e ruim em julgar áreas relativas.
- Um gráfico de barras ou um diagrama de pontos são formas preferíveis de dispor este tipo de dado.

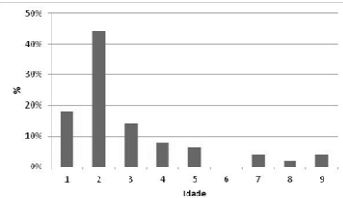
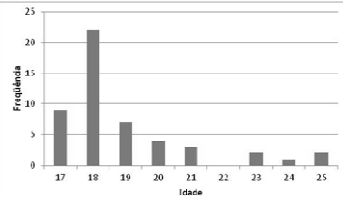
Cleveland (1985): "Dados que podem ser mostrados por um gráfico de setores sempre podem ser mostrados por um gráfico de barras ou um diagrama de pontos. Isto significa que julgamentos da posição em meio a uma escala comum podem ser feitos em vez de julgamentos menos acurados via ângulos dos setores."

Gráfico de Barras

- Para cada valor da variável desenha-se uma barra com altura correspondente à sua frequência (absoluta ou relativa)
 - ✓ Eixo das abscissas (x): valores da variável
 - ✓ Eixo das ordenadas (y): frequências absolutas ou relativas
- Adapta-se melhor às variáveis quantitativas discretas ou qualitativas ordinais

Exemplo: Idade de Alunos

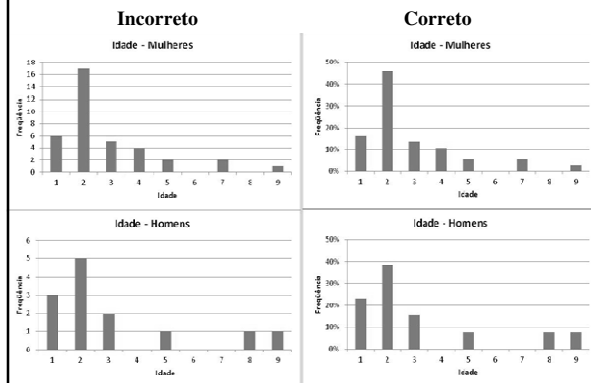
Idade	n_i	f_i
17	9	18,0%
18	22	44,0%
19	7	14,0%
20	4	8,0%
21	3	6,0%
22	0	0,0%
23	2	4,0%
24	1	2,0%
25	2	4,0%
Total	50	100,0%



Recomendações

- Colunas sempre com mesma largura
- Distância entre colunas deve ser constante
- Para comparar diferentes amostras:
 - √ Utilizar frequências relativas
 - √ Uniformizar as escalas de ambos os eixos

Comparação Idade vs. Sexo



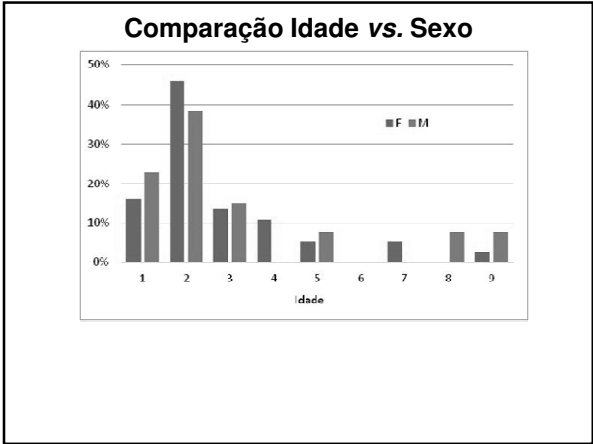
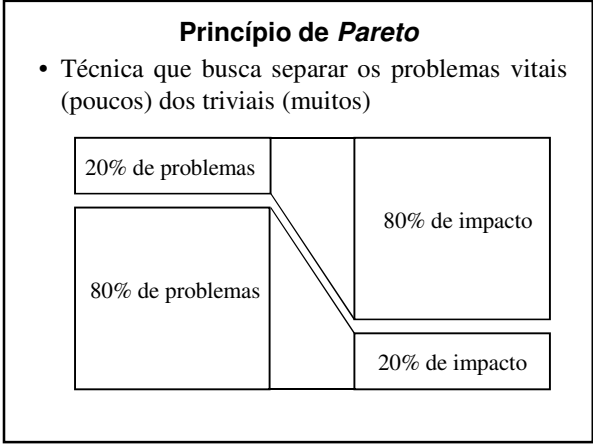


Gráfico de Pareto

- É essencialmente um gráfico de barras com os itens ordenados por tamanho
- Objetivo:
 - √ Ordenar tipo de problemas por tamanho
 - √ Foco na gestão dos problemas mais importantes



Problemas

- “Poucos e vitais”:
 - √ Representam um **pequeno número de problemas** que, no entanto, resultam em **grandes perdas**.
- “Muitos e triviais”:
 - √ São um **grande número de problemas** que resultam em **perdas pouco significativas**.

Objetivo

- Identificar as causas dos “poucos problemas vitais”;
- Focar na solução dessas causas;
- Eliminar uma parcela importante dos problemas com um pequeno número de ações.

Diagrama de Pareto

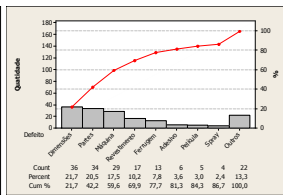
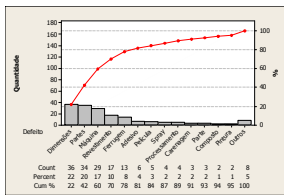
- Distribuição de frequências de dados organizados por categorias:
 - √ Marca-se a frequência total de ocorrência de cada defeito vs. o tipo de defeito
 - √ Uma escala para frequência absoluta e outra para a frequência relativa acumulada.

Diagrama de Pareto

- Identifica-se rapidamente os problemas que ocorrem com maior frequência
- Os problemas mais frequentes não são necessariamente os mais importantes.

Exemplo

- Gráfico Pareto



Outros: 5%

Outros: 15%

Procedimento

- Categorizar os quesitos (problemas) do processo
- Coletar a frequência de cada um deles durante um período
- Ordenar do mais frequente para o menos frequente
- Construir um gráfico de barras
- Adicionar um gráfico de frequências acumuladas

Exemplo

√ Problemas em empréstimos de livros em biblioteca escolar

Problema	n _i	f _i	f _{ac}
Empréstados	135	32,3%	32,3%
Em uso no recinto	103	26,3%	59,3%
Pedido não localizado	57	13,3%	73,2%
Na encadernação	43	10,5%	83,7%
Fora de lugar	17	2,3%	86,6%
Em processamento	12	2,3%	89,5%
Classificação errada	7	1,7%	91,2%
Outros	35	8,3%	100,0%
Total	410	100,0%	

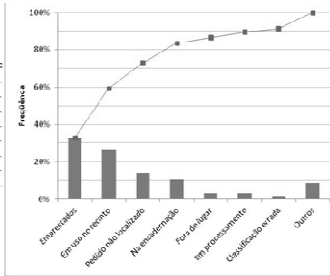
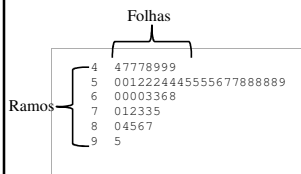


Gráfico Ramo-e-Folhas

- Dados são agrupados preservando quase toda a informação numérica
- Adequado para representação de conjunto de dados de 15 a 150 valores, aproximadamente

Exemplo: Peso



4	4
4	7778999
5	001222444
5	5555677888889
6	000033
6	68
7	01233
7	5
8	04
8	567
9	
9	5

cada linha: folhas 0, 1, 2, ..., 9

1ª. linha: folhas 0, 1, 2, 3, 4

2ª. linha: folhas 5, 6, 7, 8, 9

- folha representa um único dígito

√ 60,5 kg → 6 1 0

- Representar os valores:
220 214 222 218 223 210 223 210 227 225 212
- Suponha que queremos dividir cada número após o 2º. dígito:
220 = 22 | 0
- Procedimento:
 - √ Ramos do gráfico
 - √ Adicione 220 ao gráfico
 - √ Adicione 214 ao gráfico
 - √ Adicione demais números
 - √ Ordene as folhas
- No exemplo: intervalo de classes = 10

21	0	0	2	4	8
22	0	2	3	3	5

Expandindo o Gráfico

- Folhas 0, 1, 2, 3, 4 em uma linha
- Folhas 5, 6, 7, 8, 9 na seguinte
- Valores:
220 214 222 218 223 210 223 210 227 225 212

21	0	0	2	4	
21	8				
22	0	2	3	3	
22	5	7			

Informe de Unidades

- Unidades 8 | 3 = 83.000
9 | 7 = 97.000
- Unidades 8 | 3 = 0,083
9 | 7 = 0,097

Comentários

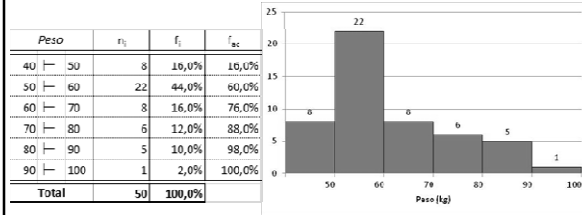
- Um gráfico ramo-e-folhas com menos de 5 ramos ativos é altamente não informativo
- Em geral, não se usa mais que 10 a 15 ramos ativos
- Regras práticas e definitivas são improdutivas
 - √ gráficos de comprimentos diferentes podem transmitir informações diferentes

Atividade nº 3

Histograma

- Características da forma do histograma:
 - √ número, largura e altura dos retângulos
- Retângulos contíguos:
 - √ eixo abscissas (x): base correspondente ao intervalo de classe
 - √ eixo das ordenadas (y): altura correspondente à frequência (ou porcentagem) do intervalo de classe
- Usado para representação gráfica da distribuição de variáveis contínuas
 - √ São parecidos com os gráficos de ramo-e-folhas

Exemplo: Peso



- Em geral, utilizam-se de 5 a 8 faixas com mesma amplitude (preferencialmente)

Histograma – Construção

- Determinam-se o máximo e o mínimo dos dados
- Divide-se a amplitude dos dados em um número conveniente de intervalos de classe de tamanhos iguais
- Contam-se a quantidade de observações que caem em cada um desses intervalos (frequência)
- Altura do retângulo acima de um intervalo de classe é igual à frequência

ESTATURA DAS MENINAS DESTA SALA - 2009

CLASSE	ESTATURAS (cm)	FREQUÊNCIA
Li	150 154	4
Ls	154 158	9
	158 162	11
	162 166	8
h = Ls-Li	166 170	5
	170 174	3
	TOTAL	40

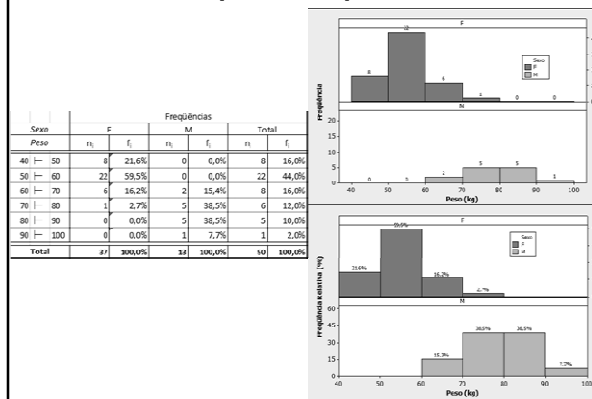
FONTE: Novaes, 2009.

AT = Ls max-Li min Ponto médio = (Ls - Li)/2

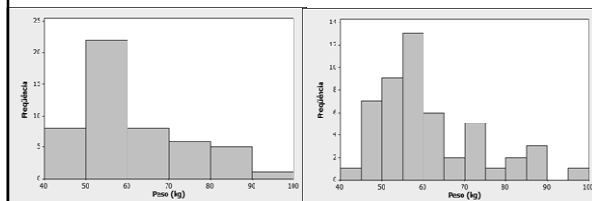
Histograma – Comparações

- Histograma de frequência relativa:
 - √ Altura do retângulo = frequência relativa do intervalo
 - √ Conveniente para comparar histogramas baseados em amostras de tamanhos diferentes
- Motivo: aspectos principais captados no histograma: formato geral e área dos retângulos
 - √ Se intervalos de classe são iguais essas áreas são proporcionais às frequências

Exemplo: Peso por Sexo

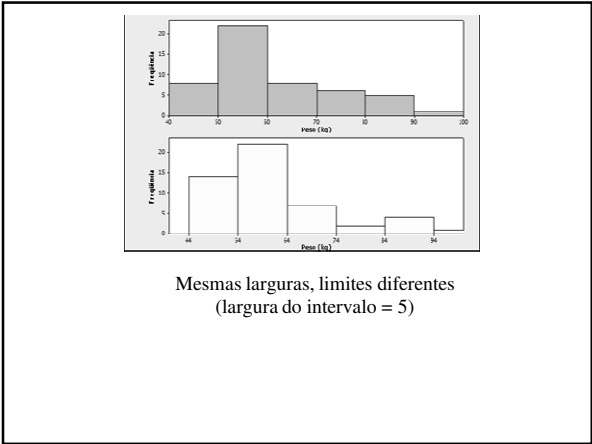


- Formato do histograma depende:
 - √ largura escolhida para os intervalos de classe
 - √ posicionamento dos extremos dos intervalos de classe



Histograma original
(largura do intervalo = 10)

Largura de intervalo modificada
(largura do intervalo = 5)

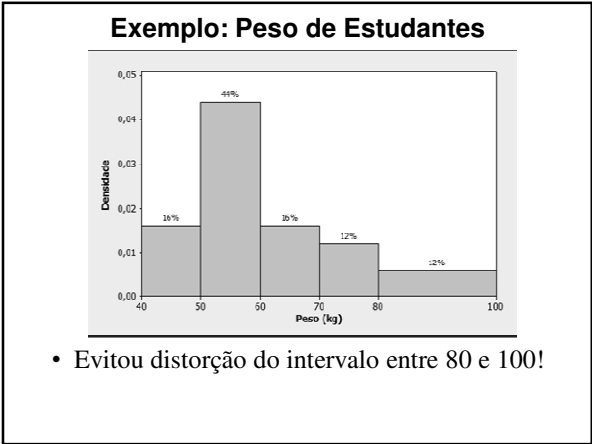


Histograma de Densidade

- Área de cada retângulo representa a frequência relativa do intervalo de classe correspondente
 \checkmark Soma das áreas de todos os retângulos = 1 (100%)
- Densidade de frequência: altura do retângulo

$$\text{densidade} = \frac{\text{área retângulo}}{\text{amplitude intervalo}}$$

- O histograma de densidade não fica distorcido quando ele é construído com intervalos de amplitudes diferente

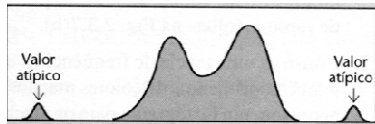


Interpretação de Gráficos de Ramo-e-Folhas & Histograma

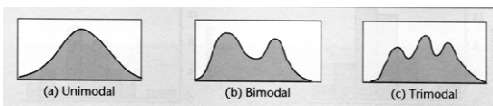
- Em uma análise gráfica procuramos identificar:
 - √ PADRÃO GLOBAL nos dados
 - √ Desvios acentuados em relação ao mesmo
- Importante:
 - √ Não perceberemos padrões nos dados se houver um número muito pequeno ou muito grande de intervalos de classe
- Procuramos uma impressão geral suavizada (não reagimos a pequenas subidas ou descidas)

Valores Atípicos (*Outliers*)

- Procuramos por observações que estejam bem afastadas da maioria dos dados
 - √ Observações discrepantes (*outliers*)
- Analisar estas observações com mais cuidado
 - √ Porque razão são tão diferentes?
 - √ Está ocorrendo algo incomum ou interessante?
 - √ São erros?



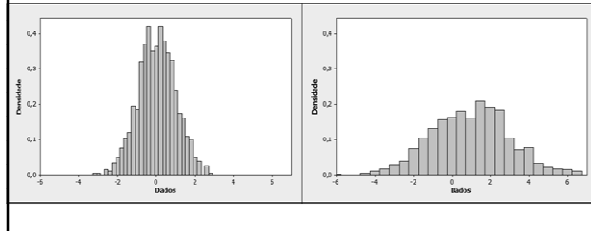
Existência de Mais de Um Pico



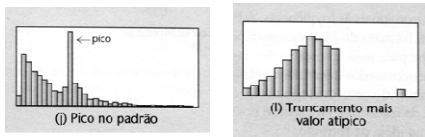
- √ Picos são chamados Modas
- √ Quando há apenas um pico, a moda representa o valor mais popular (ou classe)
- √ Presença de diversas modas é indicador de diversos grupos distintos de dados
- √ Em geral, deve-se investigar os motivos de multimodalidade

Valores Centrais e Dispersão

- Observar:
 - √ Onde os dados parecem estar centrados
 - √ Quão espalhados estão os dados
 - √ Posição das modas (caso de multimodalidade)



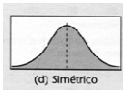
Mudanças Abruptas



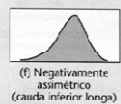
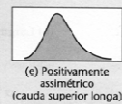
- √ Suspeite de mudanças abruptas
- √ Tente estabelecer suas causas

Forma da Distribuição

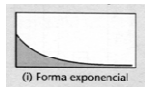
- O gráfico parece ser aproximadamente simétrico?



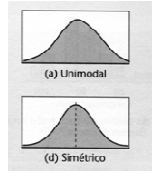
- O gráfico apresenta assimetria moderada?



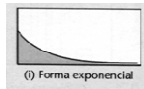
- O gráfico apresenta assimetria extrema?



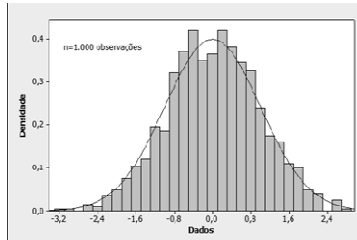
- A envoltória do gráfico tem aproximadamente forma de sino?



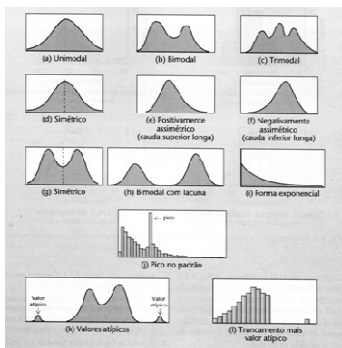
- ou tem forma exponencial?



- Usualmente, técnicas estatísticas formais preferem trabalhar com um histograma simétrico com forma de sino
- A forma do histograma pode sugerir uma função matemática cuja curva se ajusta bem ao histograma



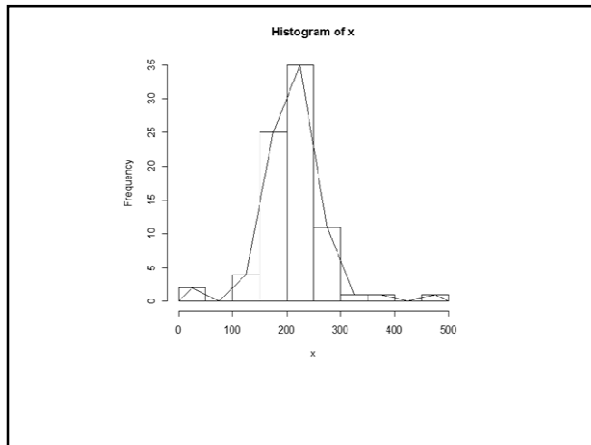
- Características a serem procuradas nos histogramas:



Fonte: Wild, C.J & Seber, G.A *Encontros com o Acaso, LTC, 2000*

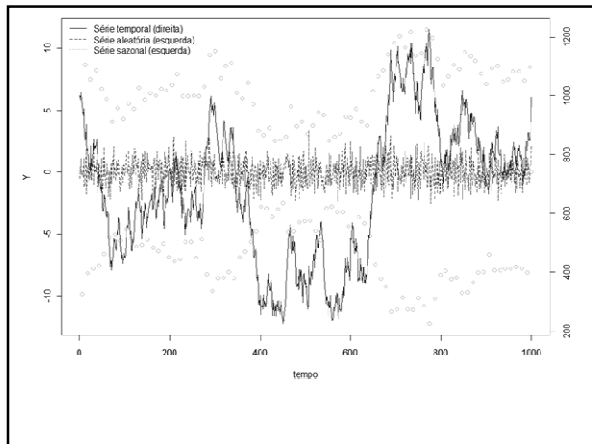
Polígono de Frequências

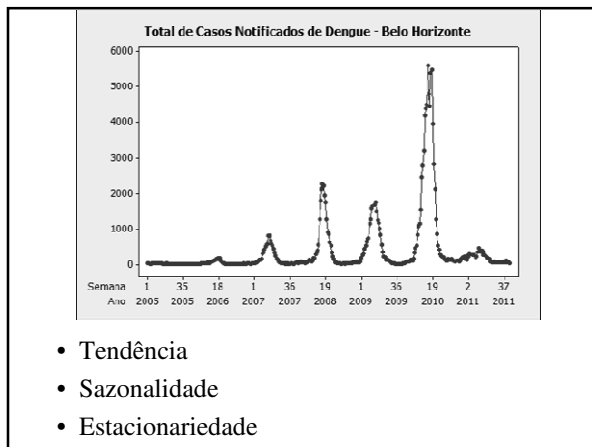
- Construído a partir do histograma
- Segmentos de retas unindo as ordenadas dos pontos médios de cada classe
- Assim como o histograma, serve para visualização da forma da distribuição de frequências da variável



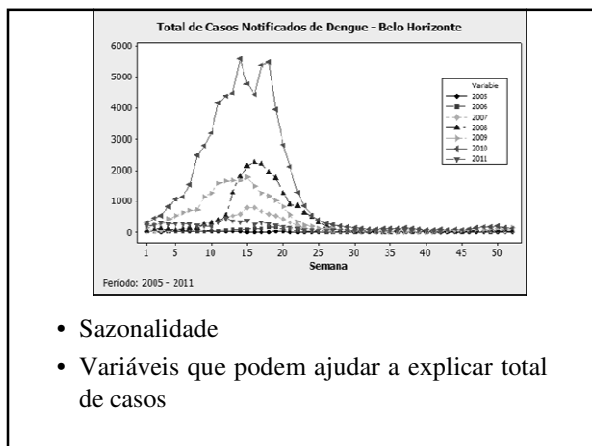
Séries Temporais

- Coleção de observações feitas sequencialmente ao longo do tempo
 - √ Em séries temporais a ordem dos dados é fundamental.
- Característica importante:
 - √ Observações vizinhas são dependentes
- Interesse: analisar e modelar esta dependência



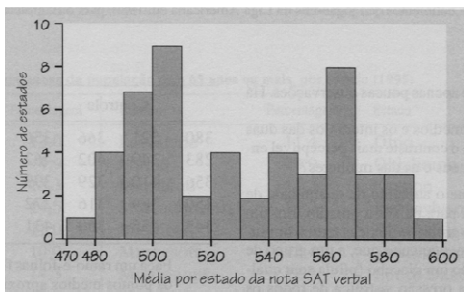


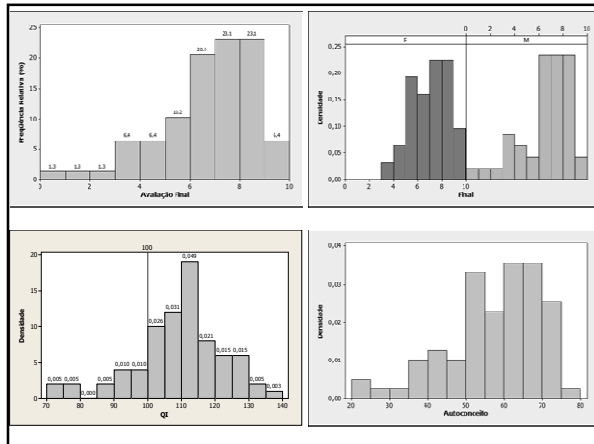
- Tendência
- Sazonalidade
- Estacionariedade



- Sazonalidade
- Variáveis que podem ajudar a explicar total de casos

Atividade nº 4





Medidas Resumo

- Medidas Resumo**
- Medidas que sintetizam informações contidas nas variáveis em um único número
 - Tipos:
 - √ Medidas de tendência central
 - √ Medidas de dispersão
 - √ Quartis, Decis e Percentis
 - √ Medidas de assimetria
 - √ Medidas de curtose

Medidas de Tendência Central

Medidas de Tendência Central

- Em geral, podem ser interpretadas como o ponto ao redor do qual os dados são distribuídos
- Algumas medidas de posição (tendência central):
 - √ Média
 - √ Mediana
 - √ Moda

Média

- Tendência central dos dados caracterizada pela média aritmética simples;
 - √ Média amostral
 - √ Média populacional

Média Amostral

- Os dados em geral são provenientes de uma amostra de observações selecionada de uma população
- Definição:

Se n observações em uma amostra forem denotadas por x_1, x_2, \dots, x_n , a média amostral será:

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

Exemplo – Peso

- Peso (kg)
- $n = 50$ indivíduos
- Média amostral

$$\bar{x} = \frac{3.046,4}{50} = 60,93 \text{ kg}$$

Média Populacional

- Valor médio de todas as observações em uma população:

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i$$

- A média amostral é um '*bom*' estimador da média populacional

Mediana

- Valor que divide a distribuição dos dados em duas partes de igual tamanho



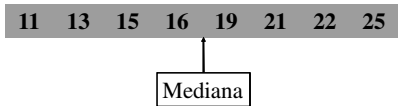
- 50% das observações ficam acima da mediana e 50%, abaixo

- Determinação da mediana:

√ Quantidade ímpar de observações:



√ Quantidade par de observações



Procedimento

- Ordenar os dados
- Se n for ímpar:
 - √ A mediana é o valor do elemento central
 - √ Elemento de ordem $\frac{n+1}{2}$
- Se n for par:
 - √ A mediana é o valor médio entre os dois elementos centrais
 - √ Elementos de ordem $\frac{n}{2}$ e $\frac{n}{2} + 1$

Exemplo – Peso (kg)

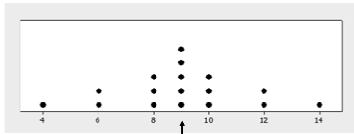
- Peso (kg)
- $n = 50$ indivíduos
- Valor médio entre o 25º e o 26º valores ordenados

$$x_{(25)} = 58; x_{(26)} = 58$$

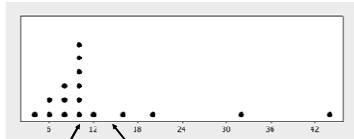
- Mediana

$$\hat{x} = \frac{58 + 58}{2} = 58 \text{ kg}$$

Média & Mediana



$\bar{x} = \tilde{x} = 9,0$



$\tilde{x} = 9,0$ $\bar{x} = 12,8$

Média e Mediana

- Valores atípicos (muito grandes ou muito pequenos) causam grandes variações na média
- Em geral, a mediana não é afetada da mesma forma que a média
- A mediana é uma medida mais robusta (menos afetada pro valores atípicos)

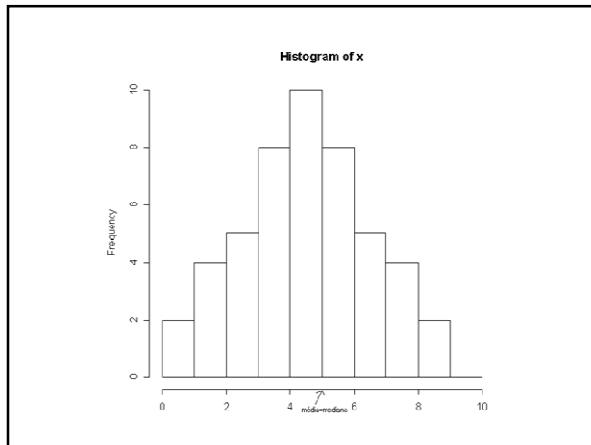
Média vs. Mediana

Média	Mediana
<ul style="list-style-type: none"> • fácil de ser manipulada algebricamente; • representa o “centro de massa” dos dados (ponto de equilíbrio no histograma). • afetada grandemente por valores extremos . 	<ul style="list-style-type: none"> • difícil de ser manipulada algebricamente; • valor da posição central dos dados ordenados; • não é afetada por valores extremos.

Média vs. Mediana (2)

- Para distribuições muito assimétricas, a mediana é uma medida mais apropriada para caracterizar um conjunto de dados.
- Se a distribuição é aproximadamente simétrica, então média e mediana são aproximadamente iguais.

√ Em distribuições perfeitamente simétricas média = mediana.



Média – Dados em Tabelas de Frequência

- Para dados disponíveis apenas em tabela de frequências
- Para calcular a média em tabela com k classes:

Ponto Médio	Frequência
x_1	f_1
x_2	f_2
\vdots	\vdots
x_k	f_k

$$n = \sum_{i=1}^k f_i$$

$$\bar{x} = \frac{x_1 f_1 + x_2 f_2 + \dots + x_k f_k}{n} = \frac{1}{n} \sum_{i=1}^k x_i f_i$$

• Exemplo - Tabela de Frequências – Peso

Peso (kg)	Ponto Médio (x_i)	Freq. Absoluta (f_i)	$x_i \cdot f_i$
40 50	45	8	360
50 60	55	22	1210
60 70	65	8	520
70 80	75	6	450
80 90	85	5	425
90 100	95	1	95
Total			3060

$$\bar{x}_{tab.} = \frac{3060}{50} = 61,20 \text{ kg}$$

$$\bar{x}_{exata} = 60,93 \text{ kg}$$

Moda

- É o valor mais frequente da distribuição.
- No histograma, ou na tabela de frequências, a classe modal é a classe de maior frequência e a moda são aproximadas pelo ponto médio da classe.

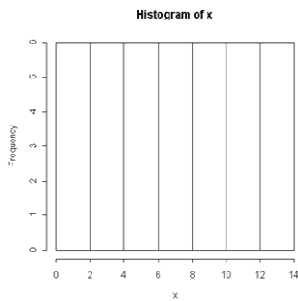
Exemplo: Peso

- Classe Modal: [50; 60)
 $\sqrt{\text{Maior frequência}} = 22$ observações
- Moda: 55 kg

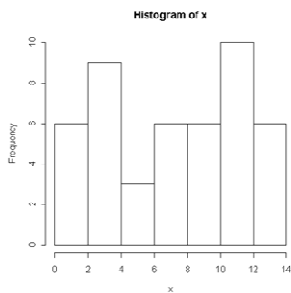
Moda (2)

- Uma distribuição pode não possuir moda (amodal – distribuição “achatada”).
- Uma distribuição pode possuir mais de uma moda (multimodal).
- Uma distribuição pode possuir apenas uma moda (unimodal).

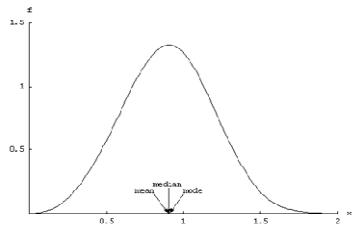
Distribuição “Achatada”



Distribuição Multimodal

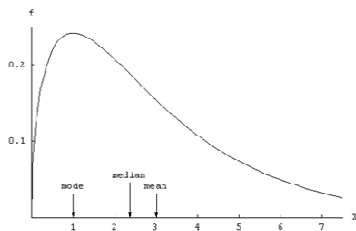


Medidas de Posição – Distribuições Simétricas



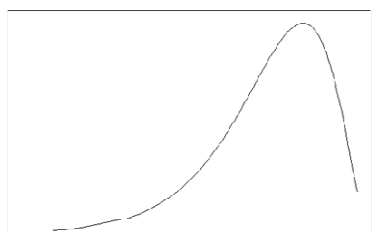
média = mediana = moda

Medidas de Posição – Distribuições Assimétricas à Direita



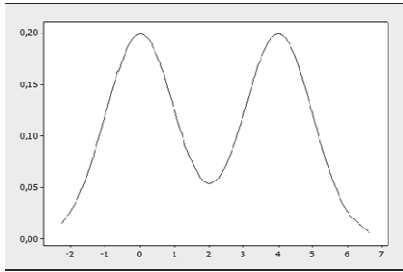
média > mediana > moda

Medidas de Posição – Distribuições Assimétricas à Esquerda



média < mediana < moda

Distribuições Bimodais



média = mediana \neq moda

Medidas de Dispersão

Comparação entre Grupos de Dados

Stem-and-Leaf Display: grupo_1

Stem-and-leaf of grupo_1 N = 10
Leaf Unit = 0,10

```
(10) 5 0000000000
```

Stem-and-Leaf Display: grupo_2

Stem-and-leaf of grupo_2 N = 10
Leaf Unit = 0,10

```
4 2 0000  
5 3 0  
5 4  
5 5  
5 6  
5 7 0  
4 8 0000
```

Stem-and-Leaf Display: grupo_3

Stem-and-leaf of grupo_3 N = 10
Leaf Unit = 0,10

```
3 4 000  
(4) 5 0000  
3 6 000
```

Stem-and-Leaf Display: grupo_4

Stem-and-leaf of grupo_4 N = 10
Leaf Unit = 0,10

```
1 1 0  
2 2 0  
3 3 0  
4 4 0  
(2) 5 00  
4 6 0  
3 7 0  
2 8 0  
1 9 0
```

Stem-and-Leaf Display: grupo_5

Stem-and-leaf of grupo_5 N = 10
Leaf Unit = 0,10

```
1 3 0  
3 4 00  
(4) 5 0000  
3 6 00  
1 7 0
```

Média e Mediana

- Todos os conjuntos têm média e mediana iguais a 5
- Podemos afirmar que a distribuição dos dados é a mesma?

Comentários

- Há grandes diferenças entre os grupos;
 - √ Grupo 1: Todos os valores são iguais a 5.
 - √ Grupo 2: Nenhum valor igual a 5;
 - √ Grupo 3: Valores concentrados entre 4 e 6.
 - √ Grupo 4: Valores espalhados entre 1 e 9
 - √ Grupo 5: Valores dispersos entre 3 e 7
- Além da média e da mediana, é necessário outro tipo de medida para caracterizar os grupos!

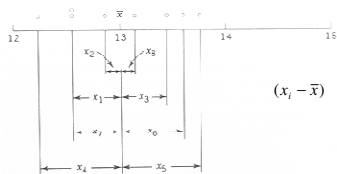
Medidas de Dispersão

- Informações importantes sobre os dados:
 - √ Valor em torno do qual os dados se **concentram**
 - √ Valor do grau de dispersão dos dados
- Medidas de dispersão mais comuns:
 - √ Amplitude amostral
 - √ Variância amostral (Desvio-padrão amostral)
 - √ Distância interquartílica (ou desvio interquartílico)

Amplitude Amostral - r

- É a mais simples das medidas de dispersão.
- É definida como: $r = \max(x_i) - \min(x_i)$
- Desvantagem:
 - ✓ Omite toda a informação entre o mínimo e o máximo
 - ✓ Em geral, quando $n < 10$, esta perda de informações não será muito séria

Construção de uma Medida de Dispersão



- Quanto maior a variabilidade dos dados, maior o valor absoluto de alguns desvios
- Valor absoluto complica o tratamento matemático
- A soma dos desvios é zero
- Uma solução: considerar o quadrado dos desvios

Variância Amostral

- É a média dos desvios quadráticos em relação à média.

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

- Tem unidade diferente dos dados.
- Por questões técnicas (Inferência), adota-se $n-1$ no denominador da média.
 - ✓ Torna-se o 'melhor' estimador

Desvio-padrão Amostral (s)

- É a raiz quadrada da variância amostral
 $\sqrt{\quad}$ A unidade de medida é a mesma dos dados!

- Conjunto de dados:

5 2 3 4 8

$$s^2 = \frac{\sum(x_i - \bar{x})^2}{n-1} = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_5 - \bar{x})^2}{5-1}$$

$$= \frac{(5-4,4)^2 + (2-4,4)^2 + (3-4,4)^2 + (4-4,4)^2 + (8-4,4)^2}{4}$$

$$= \frac{21,2}{4} = 5,3$$

$$s = \sqrt{s^2} = \sqrt{5,3} = 2,30$$

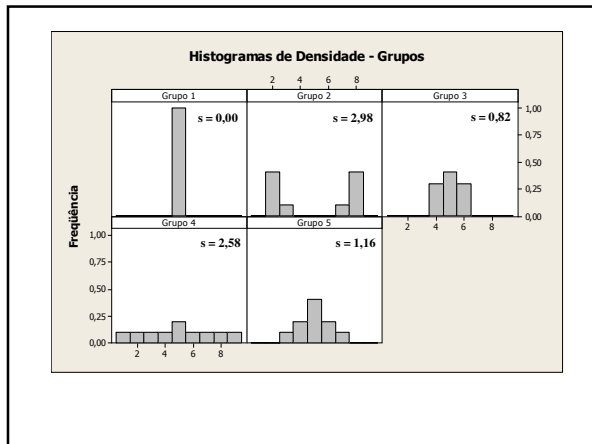
Cálculo Alternativo

- Variância: $s^2 = \frac{1}{n-1} [\sum x_i^2 - n(\bar{x})^2]$

x_i	x_i^2
5	25
2	4
3	9
4	16
8	64
22	118

$$s^2 = \frac{1}{5-1} [118 - 5(4,4)^2] = \frac{21,2}{4} = 5,3$$

$$s = \sqrt{5,3} = 2,30$$



Coeficiente de Variação

- Medida relativa de dispersão: $cv = \frac{s}{\bar{x}} \cdot 100$
- Medida adimensional
- Fornece medida de homogeneidade dos dados
 - √ Quanto menor o cv , maior a homogeneidade
- Utilidades:
 - √ Comparação grau de concentração (dispersão) em torno da média
 - √ Comparação entre variáveis (ou grupos)

Exemplo – Peso

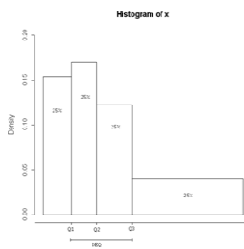
- Peso (kg)
- $n = 50$ indivíduos
- Variância: $s^2 = 148,33$
- Desvio-padrão: $s = \sqrt{148,33} = 12,18$
- Média: $\bar{x} = 60,93$
- Coeficiente de variação: $cv = \frac{s}{\bar{x}} = \frac{12,18}{60,93} = 19,99\%$

Atividade nº 5

Quartis e Percentis

Quartis

- Dividem o conjunto de dados em 4 partes iguais



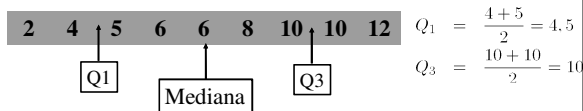
- 1º Quartil (Q_1):
25% dos dados estão abaixo (75% acima)
- 3º Quartil (Q_3):
75% dos dados estão abaixo (25% acima)
- 2º Quartil:
É a mediana!

Procedimento para Determinação dos Quartis

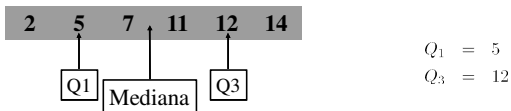
- Várias definições são usadas na literatura e por diferentes pacotes computacionais
 - √ As diferentes definições dão respostas muito parecidas
- Regra que adotaremos:
 - √ O primeiro quartil (Q_1) é a mediana de todas as observações com posição estritamente abaixo da posição da mediana
 - √ O terceiro quartil (Q_3) é a mediana das observações que estão estritamente acima da posição da mediana.

• Determinação da mediana:

$$\sqrt{n} = 9$$



$$\sqrt{n} = 6$$



Exemplo – Peso

- Peso (kg)

Mínimo	44,0	$x_{(1)} = 44,0$
Q_1	52,0	$x_{(13)} = 52,0$
$Q_2 = \text{Mediana}$	58,0	$x_{(25)} = 58,0$ $x_{(26)} = 58,0$
Q_3	68,5	$x_{(38)} = 68,5$
Máximo	95,0	$x_{(50)} = 95,0$

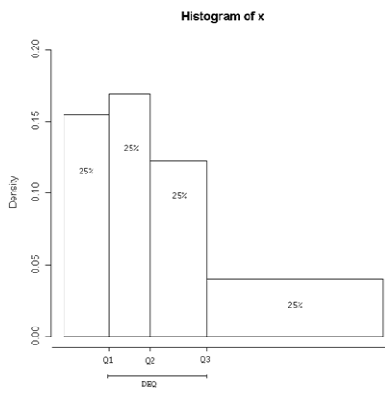
Distância Interquartílica

- Medida de variabilidade dada por .

$$DI = Q_3 - Q_1$$

- Menos sensível a valores extremos que a amplitude e a variância (desvio-padrão)
- É uma medida um pouco mais refinada que a amplitude amostral.

FB45



Exemplo: Peso

- Peso (kg)

Q_1	52,0	$x_{(13)} = 52,0$
$Q_2 = \text{Mediana}$	58,0	$x_{(25)} = 58,0$ $x_{(26)} = 58,0$
Q_3	68,5	$x_{(38)} = 68,5$
Distância Interquartílica	16,50	$Q_3 - Q_1$

Slide 133

LFB45 Refazer o gráfico em Minitab
Lupercio Bessegato; 25/09/2007

Box-plot

Esquema dos 5 Números

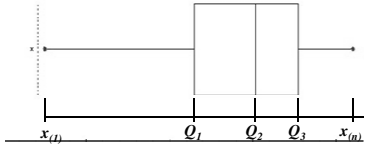
- São os cinco valores importantes para se ter uma boa ideia da assimetria dos dados.
- São as seguintes medidas da distribuição:
- $x_{(1)}$, Q_1 , Q_2 , Q_3 e $x_{(n)}$.

Esquema dos 5 Números (2)

- Para uma distribuição aproximadamente simétrica, tem-se:
 - $\sqrt{Q_2 - x_{(1)}} \cong \sqrt{x_{(n)} - Q_2}$;
 - $\sqrt{Q_2 - Q_1} \cong \sqrt{Q_3 - Q_2}$;
 - $\sqrt{Q_1 - x_{(1)}} \cong \sqrt{x_{(n)} - Q_3}$;
 - √ distâncias entre mediana e Q1, mediana e Q3 menores do que distâncias entre os extremos e Q1 e Q3.

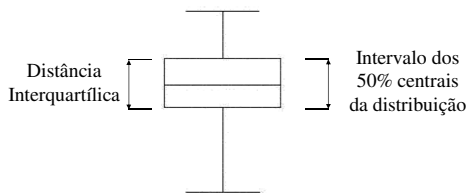
Box Plot

- A informação do esquema dos cinco números pode ser expressa num diagrama, conhecido como *box plot* (*gráfico-caixa*).
- Descreve várias características dos dados:
√ Centro, dispersão, simetria e valores atípicos



Box Plot (2)

- O retângulo é traçado de maneira que suas bases têm alturas correspondentes Q_1 e Q_3 .
- Corta-se o retângulo por segmento paralelo às bases, na altura correspondente Q_2 .
- O retângulo do *boxplot* corresponde aos 50% valores centrais da distribuição.

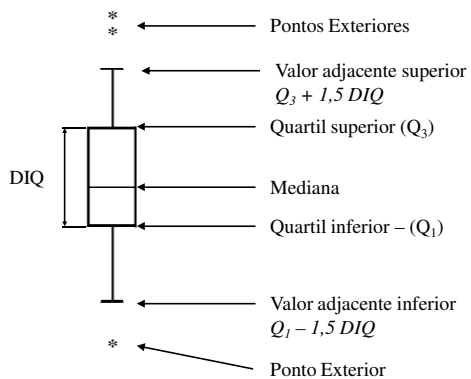


Região de Observações Típicas

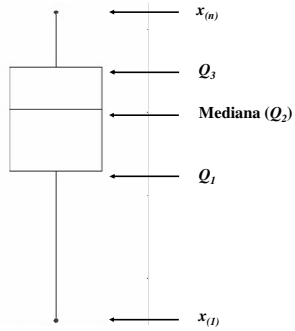
- Delimita-se a região que vai da base superior do retângulo até o maior valor observado que NÃO supere o valor de $Q_3 + 1,5 \times DIQ$.
- Procedimento similar para delimitar a região que vai da base inferior do retângulo, até o menor valor que NÃO é menor do que $Q_1 - 1,5 \times DIQ$.

Região de Observações Atípicas

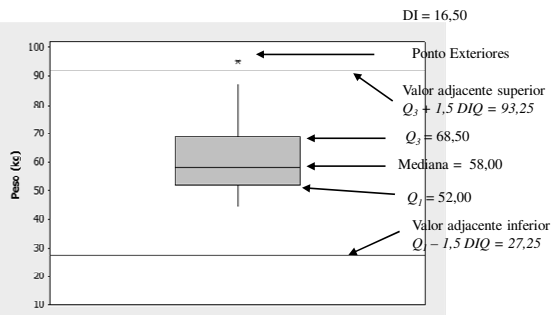
- Observações são representadas por asteriscos e situam-se:
 - √ ou, acima do Valor adjacente superior ($Q_3 + 1,5 \times DIQ$)
 - √ ou, abaixo do Valor adjacente inferior ($Q_1 - 1,5 \times DIQ$)
- Estes pontos exteriores são denominados *outliers* ou valores atípicos.



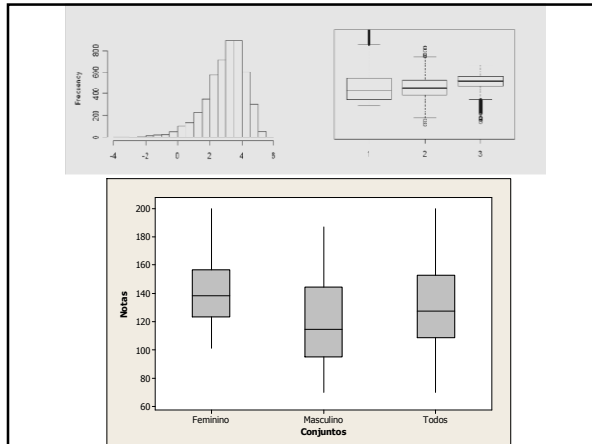
- Se não houver pontos exteriores:



Exemplo - Peso

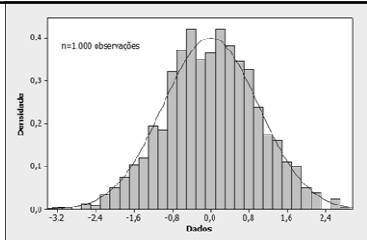


Atividade nº 6

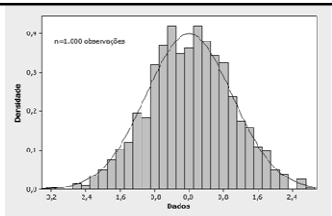


Distribuição Normal

- ### Exploração de Dados Univariados
- Faça sempre um gráfico de seus dados
 - √ Em geral, ramo-e-folhas ou um histograma
 - Procure um padrão global e desvios acentuados
 - √ *Outliers*
 - Calcule um resumo numérico para descrever o centro e a dispersão
 - Às vezes, o padrão global de um grande número de observações é tão regular que pode ser descrito por uma **curva suave**



- Curva descreve toda a distribuição em uma única expressão
 - √ Mais fácil para trabalhar
- A curva é um modelo matemático
 - √ descrição matemática idealizada



- Áreas das barras em um histograma representam contagens (ou proporções)
- Área sob a curva é exatamente 1
- Área sob a curva representa proporção de observações

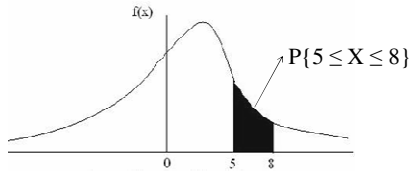
área = frequência relativa

Curva de Densidade

- A curva é denominada curva de densidade
- Propriedades:
 - √ Está sempre sobre ou acima do eixo horizontal
 - √ Tem área exatamente igual a 1 entre ela e o eixo horizontal

[FB1]

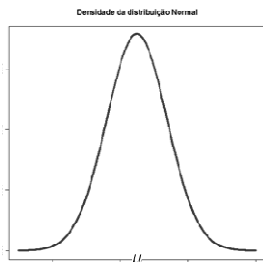
- A proporção entre 2 pontos é igual à área sob a curva, entre os dois pontos e o eixo x



Curvas Normais

- É uma classe importante de curvas de densidade
- Características:
 - √ São simétricas, unimodais e tem forma de sino
 - √ Descrevem distribuições normais (gaussianas)

Função de Densidade



$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

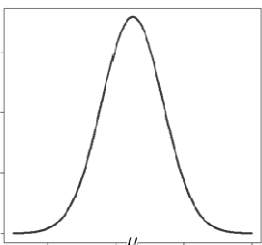
- √ O gráfico tem o formato de sino
- √ Parâmetros da distribuição normal:
 - Média (μ)
 - Desvio-padrão (σ) ou variância (σ^2)

Slide 153

LFB1 Trocar proporção por probabilidade
Lupércio Bessegato; 08/01/2012

Características

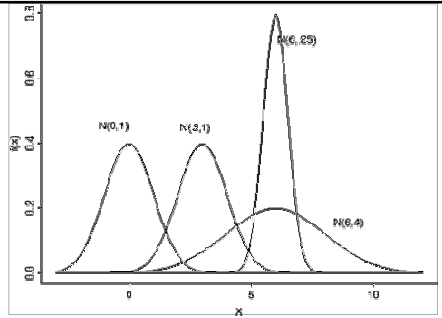
Densidade da distribuição Normal



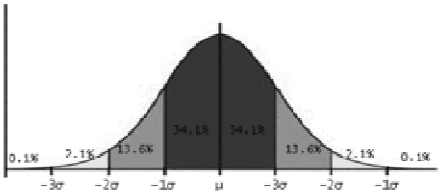
√ Simétrica em torno da média (μ)

- área antes de μ = área depois de $\mu = 0,5$
- média = mediana = moda

√ Varia de $-\infty$ a $+\infty$

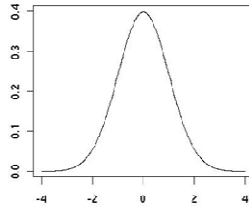


- Parâmetro de locação: μ
- Parâmetro de escala: σ^2



- Áreas de intervalos
- √ $\mu \pm \sigma \approx 68\%$
- √ $\mu \pm 2\sigma \approx 95\%$
- √ $\mu \pm 3\sigma \approx 99,7\%$

Distribuição Normal Padrão



- $Z \sim N(0, 1)$
- Média (μ) = 0
- Desvio-padrão (σ) = 1
- Valores de área tabelados

Tabela A3

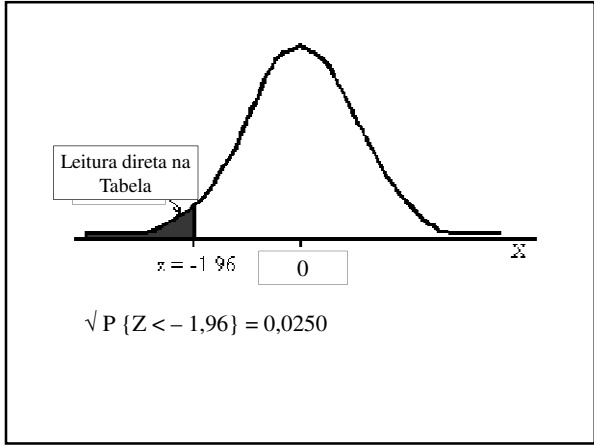
z	0,00	0,01	0,02	0,03	0,04	0,05	0,06	0,07	0,08	0,09
-4,0	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000
-3,5	0,0001	0,0002	0,0003	0,0004	0,0005	0,0006	0,0007	0,0008	0,0009	0,0010
...
0,0	0,5000	0,5040	0,5080	0,5120	0,5160	0,5199	0,5239	0,5279	0,5319	0,5359
...
4,0	0,9999	0,9998	0,9997	0,9996	0,9995	0,9994	0,9993	0,9992	0,9991	0,9990

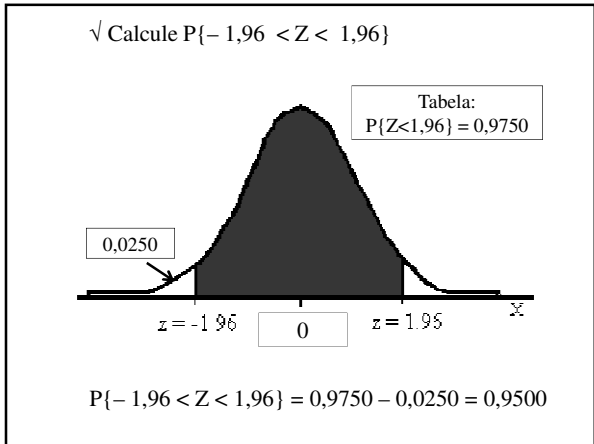
Distribuição Normal – Cálculo de Probabilidades

- Seja a variável aleatória $Z \sim N(0, 1)$
- Calcule $P\{Z < -1,96\}$
- Roteiro:
 - Esboce a curva normal
 - Trace uma linha para $z = -1,96$
 - Verifique a área que se deseja calcular
 - Determine a área a partir da tabela

√ Área sob a curva para $Z < -1,96$:

z	-0.09	-0.08	-0.07	-0.06	-0.05	-0.04	-0.03	-0.02	-0.01	0.00	z
-3.00	.0041	.0041	.0041	.0041	.0041	.0041	.0041	.0041	.0041	.0041	-3.80
-3.00	.0041	.0041	.0041	.0041	.0041	.0041	.0041	.0041	.0041	.0041	-3.70
-3.60	.0001	.0001	.0001	.0001	.0001	.0001	.0001	.0001	.0001	.0002	-3.60
-3.50	.0002	.0002	.0002	.0002	.0002	.0002	.0002	.0002	.0002	.0002	-3.50
-3.40	.0002	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	-3.40
-3.30	.0003	.0004	.0004	.0004	.0004	.0004	.0004	.0005	.0005	.0005	-3.30
-3.20	.0005	.0005	.0005	.0006	.0006	.0006	.0006	.0006	.0007	.0007	-3.20
-3.10	.0007	.0007	.0008	.0008	.0008	.0008	.0009	.0009	.0009	.0010	-3.10
-3.00	.0010	.0010	.0011	.0011	.0011	.0012	.0012	.0013	.0013	.0013	-3.00
-2.90	.0014	.0014	.0015	.0015	.0016	.0016	.0017	.0018	.0018	.0019	-2.90
-2.80	.0019	.0020	.0021	.0021	.0022	.0023	.0023	.0024	.0025	.0026	-2.80
-2.70	.0026	.0027	.0028	.0029	.0029	.0031	.0032	.0033	.0034	.0035	-2.70
-2.60	.0036	.0037	.0038	.0039	.0040	.0041	.0042	.0044	.0045	.0047	-2.60
-2.50	.0048	.0049	.0051	.0052	.0054	.0055	.0057	.0059	.0060	.0062	-2.50
-2.40	.0064	.0066	.0068	.0069	.0071	.0073	.0075	.0078	.0080	.0082	-2.40
-2.30	.0084	.0087	.0089	.0091	.0094	.0096	.0098	.0102	.0104	.0107	-2.30
-2.20	.0110	.0113	.0116	.0119	.0122	.0125	.0128	.0132	.0136	.0139	-2.20
-2.10	.0143	.0146	.0150	.0154	.0158	.0162	.0166	.0170	.0174	.0179	-2.10
-2.00	.0183	.0188	.0192	.0197	.0202	.0207	.0212	.0217	.0222	.0228	-2.00
-1.90	.0233	.0239	.0244	.0250	.0256	.0262	.0268	.0274	.0281	.0287	-1.90
-1.80	.0294	.0301	.0307	.0314	.0322	.0329	.0336	.0344	.0351	.0359	-1.80
-1.70	.0367	.0375	.0384	.0392	.0401	.0409	.0418	.0427	.0436	.0446	-1.70
-1.60	.0455	.0465	.0475	.0485	.0495	.0505	.0516	.0526	.0537	.0548	-1.60

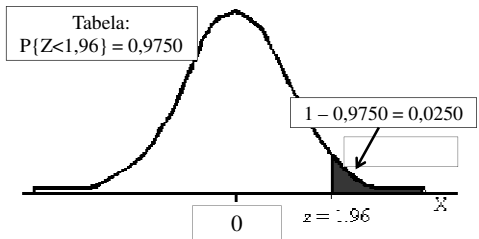




- Área sob a curva para $Z < 1,96$

z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09	z
0.00	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359	0.00
0.10	.5398	.5438	.5478	.5517	.5557	.5596	.5636	.5675	.5714	.5753	0.10
0.20	.5793	.5832	.5871	.5910	.5948	.5987	.6026	.6064	.6103	.6141	0.20
0.30	.6179	.6217	.6255	.6293	.6331	.6368	.6406	.6443	.6480	.6517	0.30
0.40	.6554	.6591	.6628	.6664	.6700	.6736	.6772	.6808	.6844	.6879	0.40
0.50	.6915	.6950	.6985	.7019	.7054	.7088	.7123	.7157	.7190	.7224	0.50
0.60	.7257	.7291	.7324	.7357	.7389	.7422	.7454	.7486	.7517	.7549	0.60
0.70	.7580	.7611	.7642	.7673	.7704	.7734	.7764	.7794	.7823	.7852	0.70
0.80	.7881	.7910	.7939	.7967	.7995	.8023	.8051	.8078	.8106	.8133	0.80
0.90	.8150	.8176	.8212	.8239	.8264	.8289	.8315	.8340	.8365	.8389	0.90
1.00	.8413	.8438	.8461	.8485	.8508	.8531	.8554	.8577	.8599	.8621	1.00
1.10	.8643	.8665	.8686	.8708	.8729	.8749	.8770	.8790	.8810	.8830	1.10
1.20	.8849	.8869	.8888	.8907	.8925	.8944	.8962	.8980	.8997	.9015	1.20
1.30	.9032	.9049	.9066	.9082	.9099	.9115	.9131	.9147	.9162	.9177	1.30
1.40	.9192	.9207	.9222	.9236	.9251	.9265	.9279	.9292	.9306	.9319	1.40
1.50	.9332	.9345	.9357	.9370	.9382	.9394	.9406	.9418	.9429	.9441	1.50
1.60	.9452	.9463	.9474	.9484	.9495	.9505	.9515	.9525	.9535	.9545	1.60
1.70	.9554	.9564	.9573	.9582	.9591	.9599	.9608	.9616	.9625	.9633	1.70
1.80	.9641	.9649	.9656	.9664	.9671	.9678	.9686	.9693	.9699	.9706	1.80
1.90	.9713	.9719	.9726	.9732	.9738	.9744	.9750	.9756	.9761	.9767	1.90
2.00	.9772	.9778	.9783	.9788	.9793	.9798	.9803	.9808	.9812	.9817	2.00
2.10	.9821	.9826	.9830	.9834	.9838	.9842	.9846	.9850	.9854	.9857	2.10
2.20	.9861	.9864	.9868	.9871	.9875	.9878	.9881	.9884	.9887	.9890	2.20
2.30	.9893	.9896	.9898	.9901	.9904	.9906	.9909	.9911	.9913	.9916	2.30
2.40	.9918	.9920	.9922	.9925	.9927	.9929	.9931	.9932	.9934	.9936	2.40

- Calcule $P\{Z > 1,96\}$

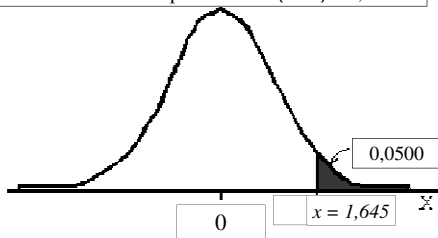


- Probabilidade contida em alguns intervalos

Intervalo	Probabilidade
$-1 < Z < 1$	$0,8413 - 0,1587 = 0,6826$
$-2 < Z < 2$	$0,9772 - 0,0228 = 0,9544$
$-3 < Z < 3$	$0,9987 - 0,0013 = 0,9974$

- Determinar x , tal que $P\{Z > x\} = 0,05$

Tabela: Valor mais próximo de $P\{Z < x\} = 0,9500$



$$\checkmark P\{Z < 1,65\} = 0,9505$$

$$\checkmark P\{Z < 1,64\} = 0,9495$$

Intervalos Simétricos em Torno de Zero

Probabilidade	Intervalo
90%	$1,645 < Z < 1,645$
95%	$-1,96 < Z < 1,96$
99%	$-2,58 < Z < 2,58$

Outras Distribuições Normais

- Caso Geral:
 - ✓ Média: μ
 - ✓ Desvio-padrão: σ

- Transformação:

$$Z = \frac{(X - \mu)}{\sigma}$$

- Mesmos procedimentos após transformação (tabela Normal Padrão)

Conversão na Normal Padrão

$$z = \frac{X - \mu}{\sigma}$$

μ x X

0 z Z

• $P\{\mu < X < x\} = P\{0 < Z < z\}$

Exemplo

- As alturas de mulheres com 18 a 24 anos de idade é aproximadamente normal com média 164 cm e desvio-padrão 6,4 cm.
- ✓ X: altura de mulheres entre 18 e 24 anos (cm)
- $X \sim N(164, 6,4)$

✓ Encontre a proporção de mulheres com altura inferior a 172 cm

√ Padronização

$$z = \frac{x - \mu}{\sigma}$$

$$z = \frac{172 - 164}{6,4} = 1,25$$

√ Pela tabela

$$P \{ Z < 1,25 \} = 0,8944$$

$$P \{ X < 100 \} = 0,8944 = 89,44\%$$

- Qual o valor de altura que delimita 5% das mulheres mais altas?

$$X = \mu + Z\sigma = 164 + 1,645(6,4) = 174,5\text{cm}$$

Aplicações da Distribuição Normal

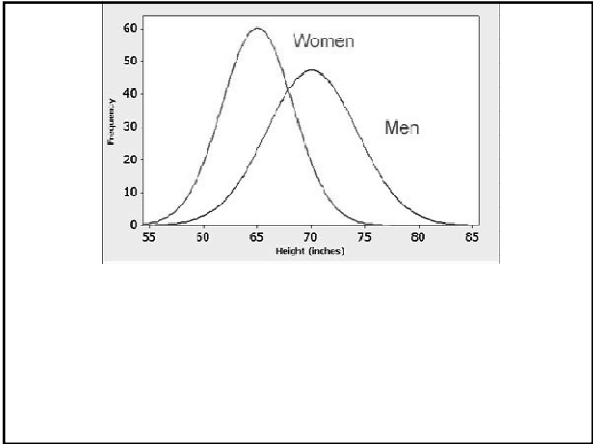
- Usada como um modelo para estudar uma grande variedade de variáveis

√ Objetivo: responder questões sobre probabilidades relacionadas com essas variáveis

- Exemplos:

√ Altura humana

√ Inteligência



Atividade nº 7

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!

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
