

FIAP

NBA

# MBA em DATA SCIENCE & ARTIFICIAL INTELLIGENCE

STATISTICS WITH R



## Dra. Regina Tomie Ivata Bernal

### Cientista de Dados na área da Saúde

#### Formação Acadêmica:

Estatístico - UFSCar

Mestre em Saúde Pública – FSP/USP

Doutor em Ciências – Epidemiologia - FSP/USP

#### Atividades Profissionais:

Professora de pós-graduação na FIAP

Consultora externa da SVS/MS

Cientista de Dados em Saúde

[profregina.bernal@fiap.com.br](mailto:profregina.bernal@fiap.com.br)  
[reginabernal@terra.com.br](mailto:reginabernal@terra.com.br)

# Avaliação da disciplina

Avaliação	Peso
Listas de exercícios	0.5
Projeto Integrado	0.5

# Objetivos da Disciplina

- Disseminar a cultura estatística quanto ao uso das técnicas descritivas, técnicas de associação e correlação tendo em vista a modelagem para previsão.
- Apresentar os conceitos básicos e metodologias para que seja extraído conhecimento de grandes bases de dados.
- Desenvolver conceitos de preparação de dados para fins estatísticos e informações para a geração de competitividade organizacional.
- Proporcionar o conhecimento necessário para reconhecer as seguintes técnicas Supervisionadas (Árvore de Decisão, Regressão Linear e Regressão Logística) e Não Supervisionadas como Componentes Principais e Análise de Cluster .

# Referências Bibliográficas

- BERRY, M.J.A.; LINOFF, G. Data Mining Techniques: for marketing, sales, and customer. Wiley Computer Publishing, 1997.
- BUSSAB, W.O.; MORETTIN, P. A., Estatística Básica, 5a. ed., São Paulo: Saraiva, 2006.
- KUHN, M. / JOHNSON, K. , Applied Predictive Modeling, 2013
- HAIR, J.F. / ANDERSON, R.E. / TATHAN, R.L. / BLACK, W.C. Análise multivariada de dados, 2009
- JAMES, G, / WITTEN, D. / HASTIE, T. / TIBSHIRANI, R. Na Introduction to Statistical Learning with Aplications in R, 2013
- LANTZ, B. Machine Learning with R. 2a. ed. Packt Publishing, 2015

# Referências Bibliográficas

- MOORE, S.D.; MCCABE, G.P.; DUCKWORTH, W.M.; SCLOVE, S.S.

**Estatística Empresarial como usar dados para tomar decisões.** Tradução Luis

Antonio Forjado. Rio de Janeiro: LTC, 2006.

- MORETIM, P.A.; TOLOI, C.M.C. **Análise de Séries Temporais**, 2ª ed., São Paulo: Edgard Blücher, 2006.
- SILVA, NN. **Amostragem Probabilística**. 2ª ed., São Paulo: Editora da Universidade de São Paulo, 2001.
- SOARES, J.; FARIAS, A. A.; CESAR, C. C., **Introdução a Estatística**, LTC, 2002.
- TORGO, L. Data Mining with R: Learning with Case Studies. 2.a ed. Chapman and Hall/CRC, 2007

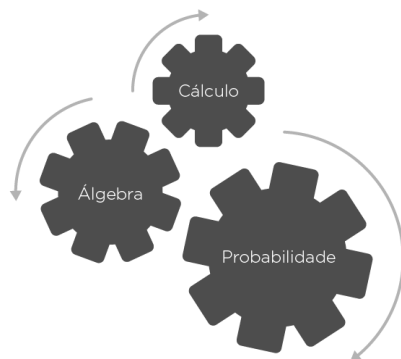


# DATA ANALYTICS

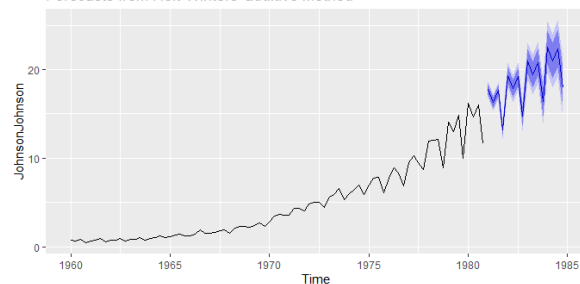


# DATA ANALYTICS

## Métodos



Forecasts from Holt-Winters' additive method



# UNIVERSO DE FORNECEDORES

MicroStrategy®

IBM®

SAP®

SAS®

Microsoft

alteryx

tableau

Qlik®

TIBCO®  
The Power of Now®

Information  
Builders

julia

ORACLE®

OPEN TEXT  
The Content Experts™

R

R Studio®



python



# SOFTWARE R





R is a [programming language](#) for [statistical computing](#) and graphics supported by the R Core Team and the R Foundation for Statistical Computing. Created by statisticians [Ross Ihaka](#) and [Robert Gentleman](#), R is used among [data miners](#), [bioinformaticians](#) and [statisticians](#) for [data analysis](#) and developing [statistical software](#). Users have created packages to augment the functions of the R language.



Robert Gentleman, co-originator of R



Ross Ihaka, co-originator of R



O programa R é opensource, muitos pesquisadores criam e disponibilizam as funções para os usuários R no formato de pacotes. Esses pacotes são submetidos ao Comprehensive R Archive Network (CRAN) para avaliação, pois existe uma política para armazenar o pacote no repositório (CRAN package repository). Atualmente existem 19.024 pacotes disponíveis para os usuários R nesse repositório



DOWNLOAD

# RStudio Desktop

Used by millions of people weekly, the RStudio integrated development environment (IDE) is a set of tools built to help you be more productive with R and Python.

## 1: Install R

## 2: Install RStudio

RStudio requires R 3.3.0+. Choose a version of R that



+

+





# Friction free data science

Posit Cloud (formerly RStudio Cloud) lets you access Posit's powerful set of data science tools right in your browser – no installation or complex configuration required.

**GET STARTED**

**ALREADY A USER? LOG IN**

If you already have a shinyapps.io account, you can log in using your existing credentials.

Fonte: <https://posit.cloud>

Cloud  
Free

\$0 / forever

[Learn more](#)

SHARED SPACES	1 with member & project limits
PROJECTS	50
COMPUTE HOURS	25 hours per month
MAX RAM	1 GB
MAX CPU	1 CPU



## Ambiente Rstudio versão Windows

Versão do R

The screenshot shows the RStudio interface with the following components and annotations:

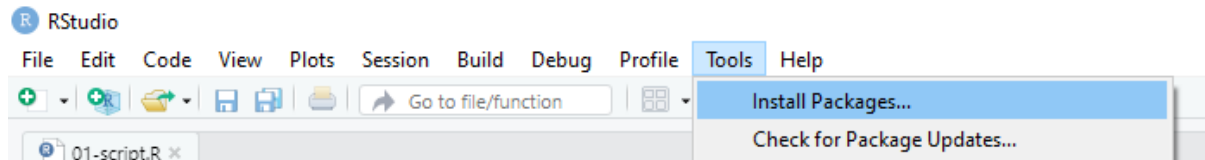
- (1) Editor de código:** The main editor window on the left, currently showing a blank file named 'Untitled1'.
- (2) Saída dos resultados:** The Console window at the bottom left, displaying the R version 4.1.2 (2021-11-01) and copyright information. A red box highlights the version number, with an arrow pointing to the 'Versão do R' label.
- (3) Saída dos gráficos:** The Plots window at the bottom right, which is currently empty.
- (4) Bases de dados:** The Environment window on the right, showing the 'Global Environment' with 132 MB of memory. A red box highlights the 'Import Dataset' button, with an arrow pointing to the 'Importar bases de dados' label.

Other visible elements include the top menu bar (File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help) and the toolbar with icons for running, saving, and other functions.

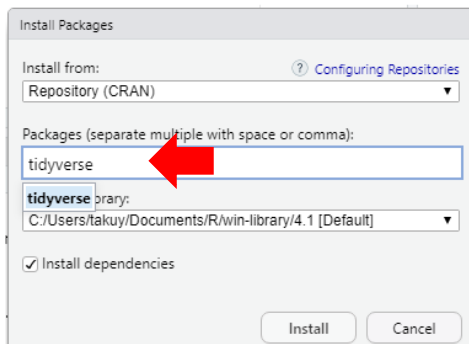




Os pacotes podem ser instalados usando a opção “Tools”:

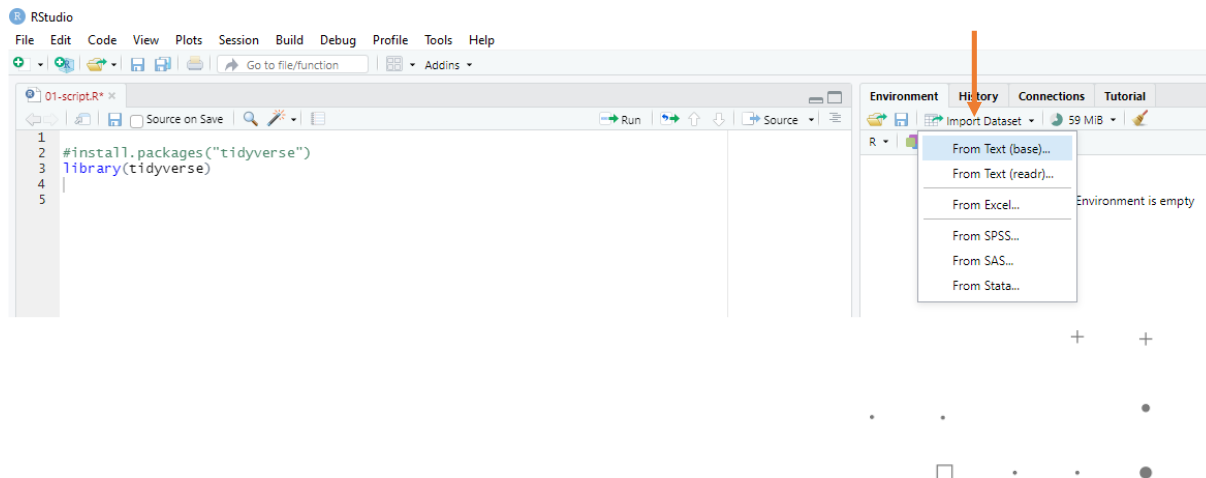


Digite o nome do pacote e selecione *Install*:





Na opção “*import Dataset*”, o usuário poderá importar a base de dados no formato **csv**, **Excel**, **SPSS**, **SAS** e **Stata**.





Matriz1

5 x 2

Coluna

1	2
3	4
5	6
7	8
9	10

Matriz2

5x 2

11	12
13	14
15	16
17	18
19	20

```
Matriz3 <- cbind(Matriz1, Matriz2)
```

Matriz3

5x 4

1	2	11	12
3	4	13	14
5	6	15	16
7	8	17	18
9	10	19	20

Linha



Matriz1  
5 x 2  
Coluna

1	2
3	4
5	6
7	8
9	10

Linha

Matriz1  
4 x 2

11	12
13	14
15	16
18	19

Matriz3 <- cbind(Matriz1, Matriz2)

1	2	11	12
3	4	13	14
5	6	15	16
7	8	18	19
9	10		



Matriz1  
5 x 2

1	2
3	4
5	6
7	8
9	10

Matriz2  
5 x 2

11	12
13	14
15	16
17	18
19	20

```
Matriz3 <- rbind(Matriz1, Matriz2)
```

Matriz2  
10 x 2

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20



Material sobre R

<https://www.r-tutor.com>

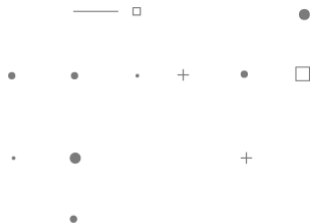
<https://www.r-bloggers.com>

# Exercitando!!!!



## Rstudio





# ESTATÍSTICA





## COMO TIRAR INFORMAÇÕES DELES

O QUE  
ACONTECEU?

DESCRITIVO



QUANTOS CANCELAMENTOS?  
QUANTOS CLIENTES NOVOS, QUANTOS ANTIGOS?  
QUAL REGIÃO?  
QUE TIPO DE CLIENTE?

POR QUE ISTO  
ACONTECEU?

DIAGNÓSTICO



QUAL A RELAÇÃO ENTRE CANCELAMENTO VOLUNTÁRIO POR TEMPO  
DE CONTA E TIPO DE CLIENTE?

O QUE  
ACONTECERÁ?

PREDITIVO



QUAL A PROBABILIDADE DE UM CLIENTE CANCELAR O  
SERVIÇO EM UMA CERTA REGIÃO NOS PRÓXIMOS 3  
MESES?

O QUE  
POSSO FAZER?

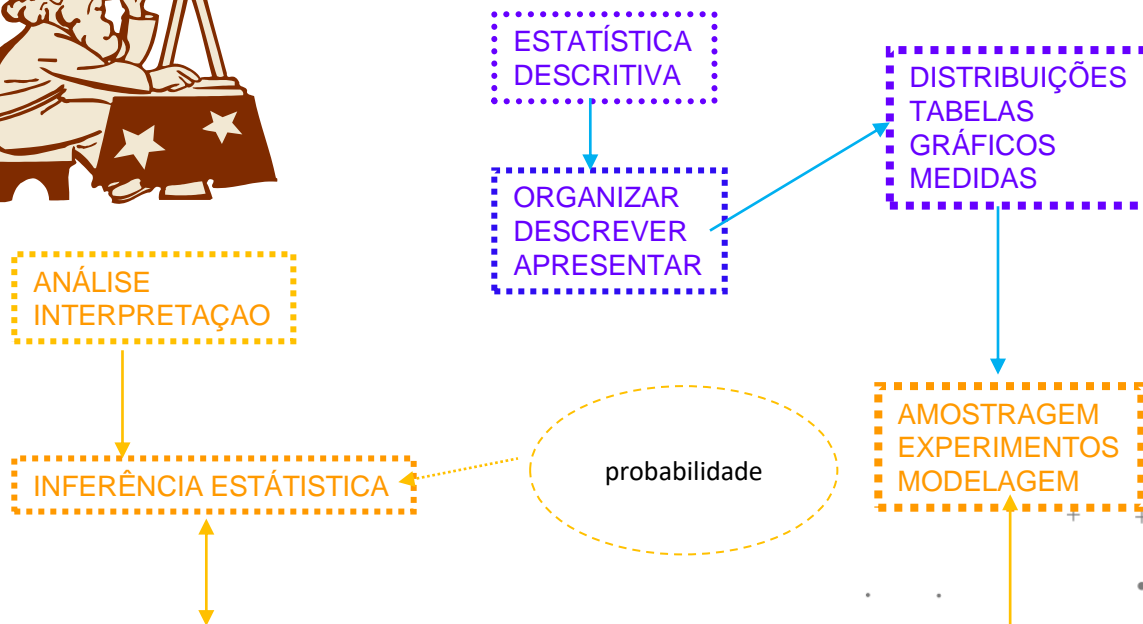
PRESCRITIVO



LISTA DE AÇÕES PARA RETER OS CLIENTES, SEGUNDO SEU  
VALOR?



# ESTATÍSTICA



# Estatística

É a ciência que trata dados numéricos provenientes de mensuração em grupos de indivíduos.

Trata da organização, descrição, apresentação análise e interpretação de dados resultantes da observação de fenômenos coletivos. Produz métodos para inferência estatística.

## ✓ Propriedades

Estuda as variações:

- entre indivíduos;
- em um mesmo indivíduo.



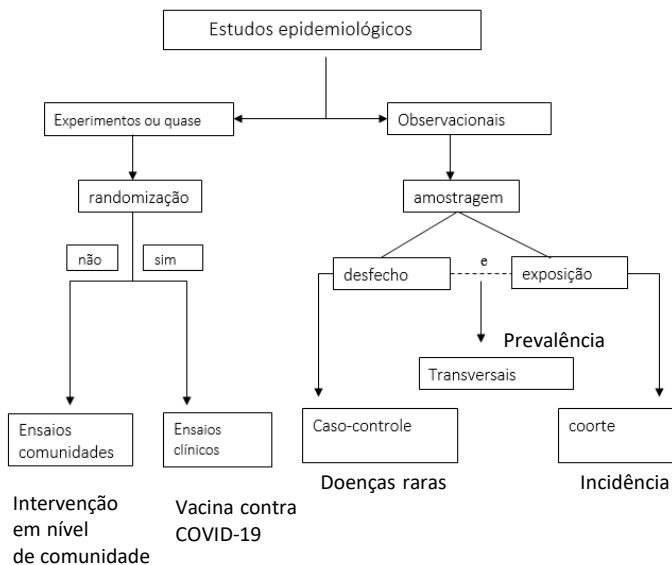
## ESTATÍSTICA

ÁREA	ESTATÍSTICA
Saúde	Bioestatística
Economia	Econometria
População	Demografia
Jurídica	Jurimetria
Biologia	Biometria
Contabilidade	Contabilometria

# BIOESTATÍSTICA

“Bioestatística é a Estatística aplicada às ciências médica e biológica.”

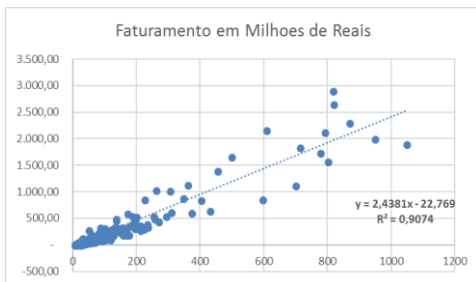
Fonte: Vieira, Sonia. Introdução a Bioestatística. Rio de Janeiro: Elsevier, 1980



## ECONOMETRIA

A Econometria consiste em uma série de ferramentas estatísticas que visam obter relações relevantes entre as variáveis econômicas a partir da aplicação de modelos matemáticos.

Fonte: <http://www.suno.com.br/artigos/econometria/>



## Exemplos de modelos econométricos

### Um modelo econométrico para previsão de impostos no Brasil.

<https://doi.org/10.1590/S1413-80502013000200006>

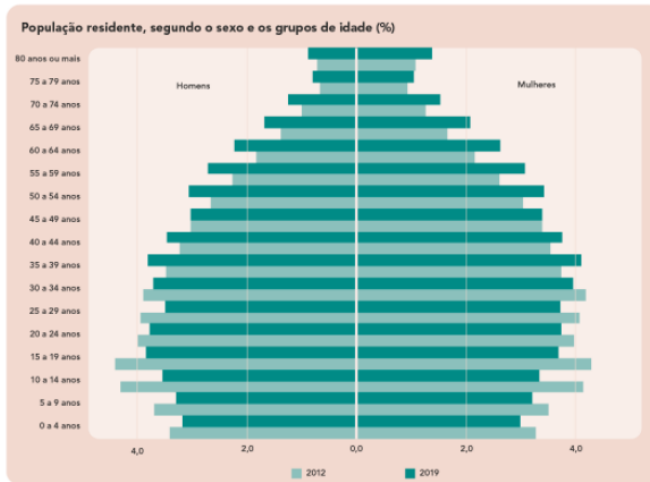
### Análise Econométrica do Comportamento dos Preços do Minério de Ferro no Mercado Mundial

[http://www.repositorio.ufop.br/bitstream/123456789/3608/1/DISSERTAÇÃO\\_AnáliseEconométricaComportamento.pdf](http://www.repositorio.ufop.br/bitstream/123456789/3608/1/DISSERTAÇÃO_AnáliseEconométricaComportamento.pdf)

## DEMOGRAFIA

“A Demografia é uma área do conhecimento que estuda a dinâmica das populações sejam elas humanas ou não.”

Fonte: <http://www.todamateria.com.br/demografia/#:::text=A%20demografia%20baseia-se%20em,dos%20seres%20vivos%20pelo%20planeta.>



Fonte: IBGE, Diretoria de Pesquisas, Coordenação de Trabalho e Rendimento, Pesquisa Nacional por Amostra de Domicílios Contínua 2012/2019.



“É um ramo da ciência que estuda a mensuração dos seres vivos. É a parte da Estatística que investiga atributos biológicos quantitativos, pertinentes a uma população de seres vivos.”

[/www.youtube.com/watch?v=DiUejcMFsjY](http://www.youtube.com/watch?v=DiUejcMFsjY)

Fonte: <http://ufpa.br/biome/bioconba.htm>

GARCIA, I. A. A segurança na identificação: a biometria da íris e da retina. 2009. 129 f. Dissertação (Mestrado em Direito Penal)— Faculdade de Direito, Universidade de São Paulo, São Paulo, 2009.

MORAES, A. F. Método para avaliação da tecnologia biométrica na segurança de aeroportos. 2006. 203 f. Dissertação (Mestrado em Engenharia Elétrica)— Escola Politécnica, Universidade de São Paulo, São Paulo, 2006.

NAKASHIRO, M. M. Biometria no Brasil e o registro de identidade civil: novos rumos para identificação. 2011. 126 f. Tese (Doutorado em Sociologia)— Departamento de Pós-Graduação em Sociologia, Universidade do Estado de São Paulo, São Paulo, 2011.

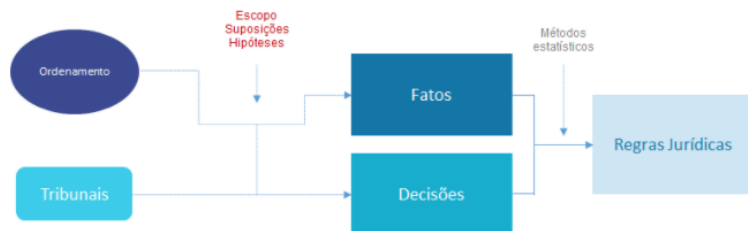




## JURIMETRIA

“Os avanços da computação possibilitaram uma nova forma de encarar as normas e a sua aplicação que baseia-se em dados e, conseqüentemente, em estatísticas. Por isso, ela pode ser genericamente definida como “a estatística aplicada do Direito”.

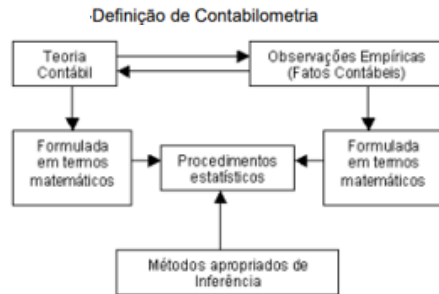
Fonte: [http:// abj.org.br/conteudo/jurimetria/](http://abj.org.br/conteudo/jurimetria/)



## CONTABILOMETRIA

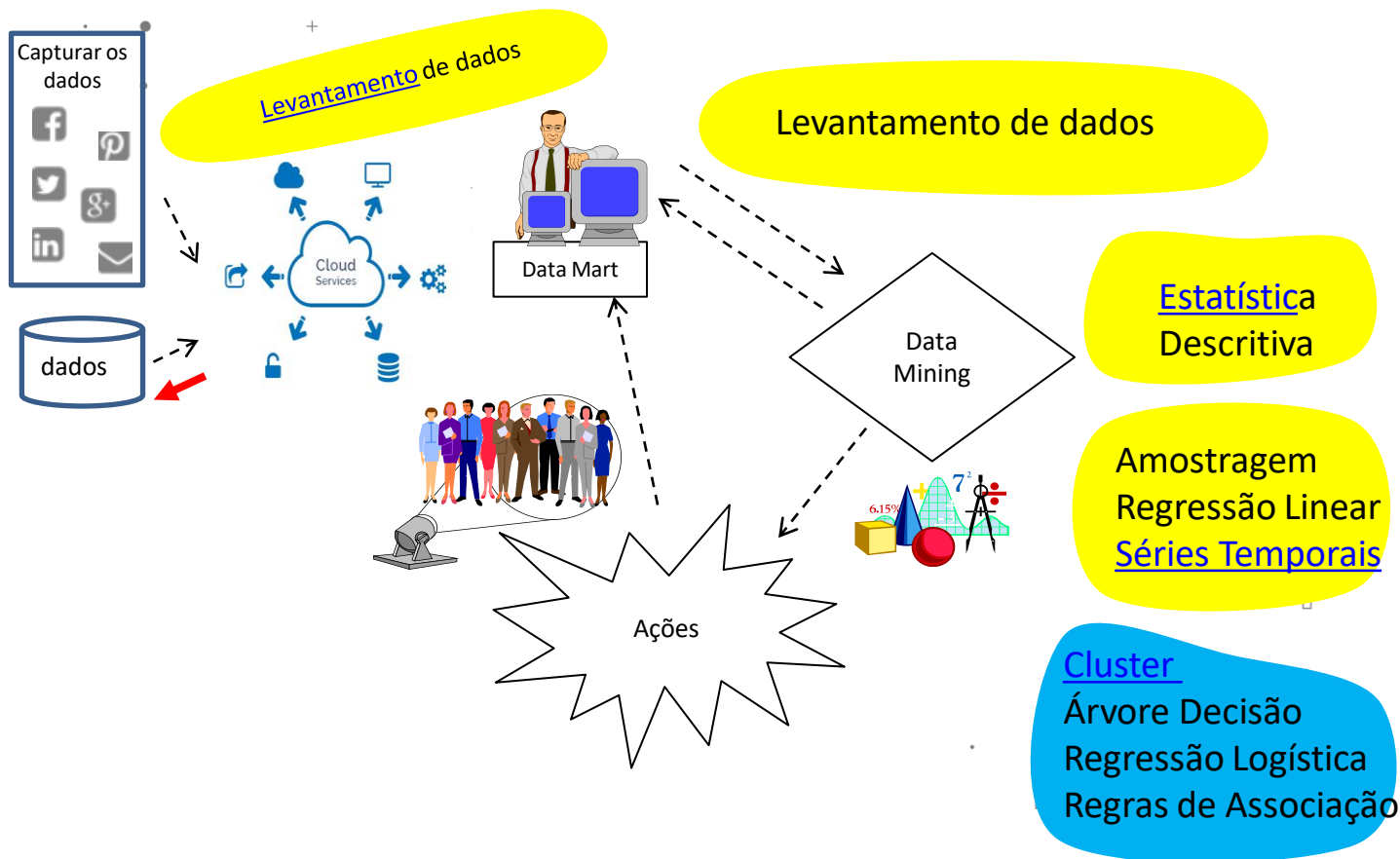
“A Contabilometria tem como característica fundamental a sua capacidade preditiva, ou seja, através da Contabilometria é possível projetar modelos de decisão eficazes que possam antecipar, prever ou estimar, de alguma forma o que irá ocorrer no futuro, provocando assim, uma melhor utilização que se pode fazer dos dados contábeis, como instrumento informativo projetado para o futuro, tornando a Contabilidade uma disciplina mais forte e mais útil.”

Fonte: <http://congressosp.fipecafi.org/anais/artigos32006/255.pdf>



Fonte: (MARION; SILVA, 1986)

# DATA ANALYTICS



# Levantamento de Dados

## Data Cleaning:

- Padronização
- Transformação de Dados
- Adoção de De-Para de Atributos

Atributo      Descrição

	De	Para	
Sexo: Masculino	2	0	
Feminino	4	1	
Idade:Criação de Faixa Etária			
	0-10		1
	11-18		2
	19-25		3
	26-30		4
	31-35		5
	36-40		6
	41-45		7
	46	8	
	sem informação	0	

+

+

.

□

.

.

□

.

.

●

●

# Levantamento de Dados

Atributo Descrição

Internet: Acesso últimos 3 meses

De	Para
1	1
3	0

Anos

de Estudo: Criação da Faixa Grau de Instrução

0-4	1
5-8	2
9-11	3
12	4

Renda: Criação de Faixa Salarial baseando em salários mínimos (Valor atual R\$380,00)

0-380	1
381-760	2
761-1900	3
1901-3800	4
3801	5
sem informação	0

Área: Agrupamento da área de residência em 1-Urbana e 2-Rural

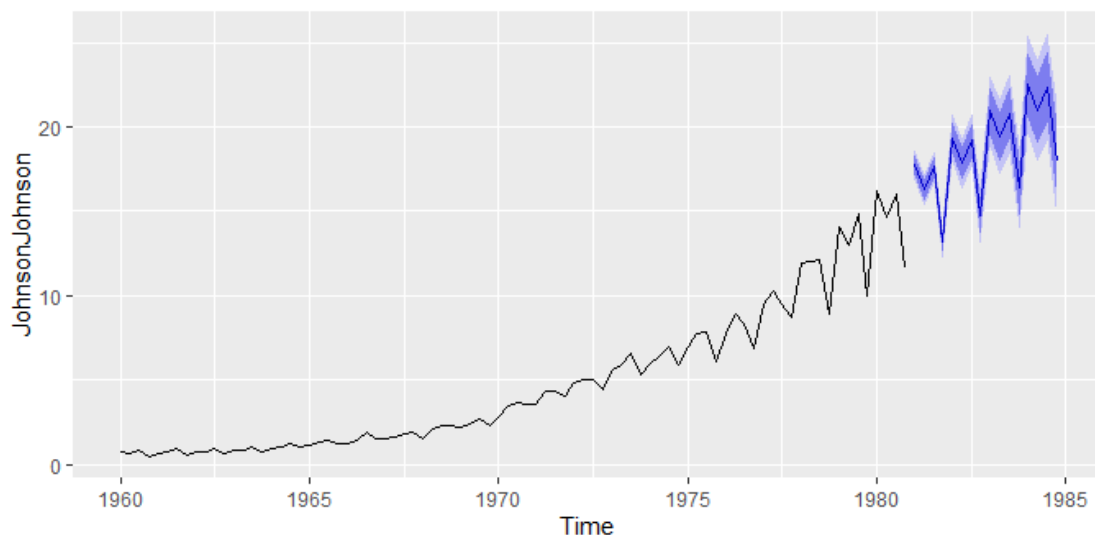
1-3	1
4-8	2

+ + .

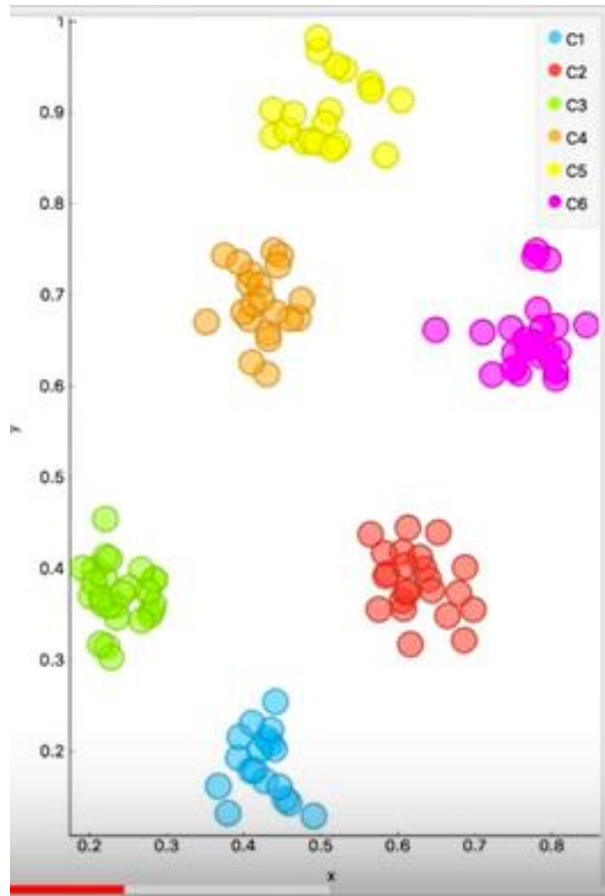
# Estatística Descritiva



Forecasts from Holt-Winters' additive method

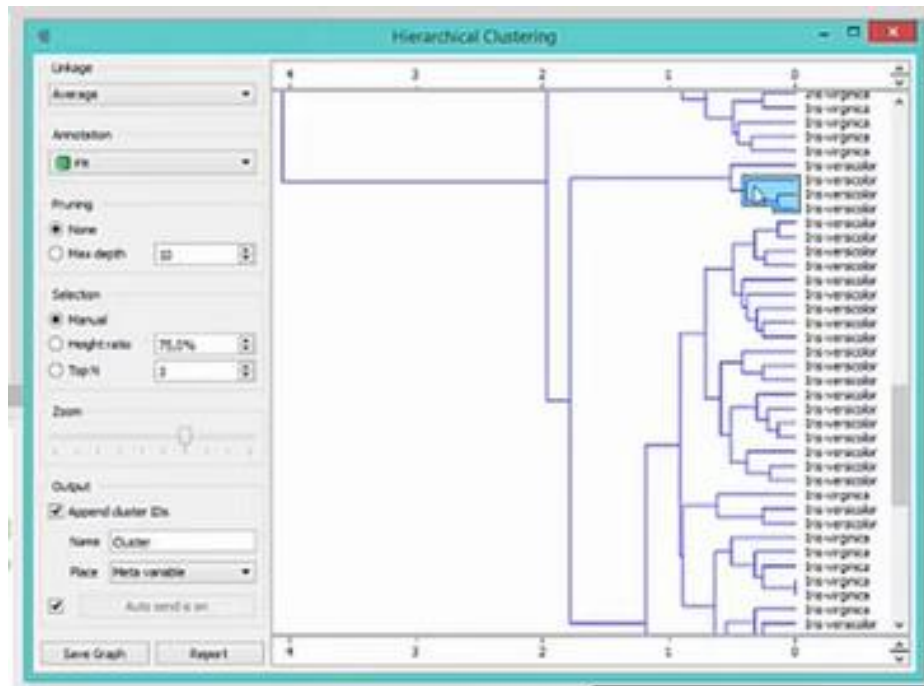


# Exemplo de Cluster Não Hierárquico



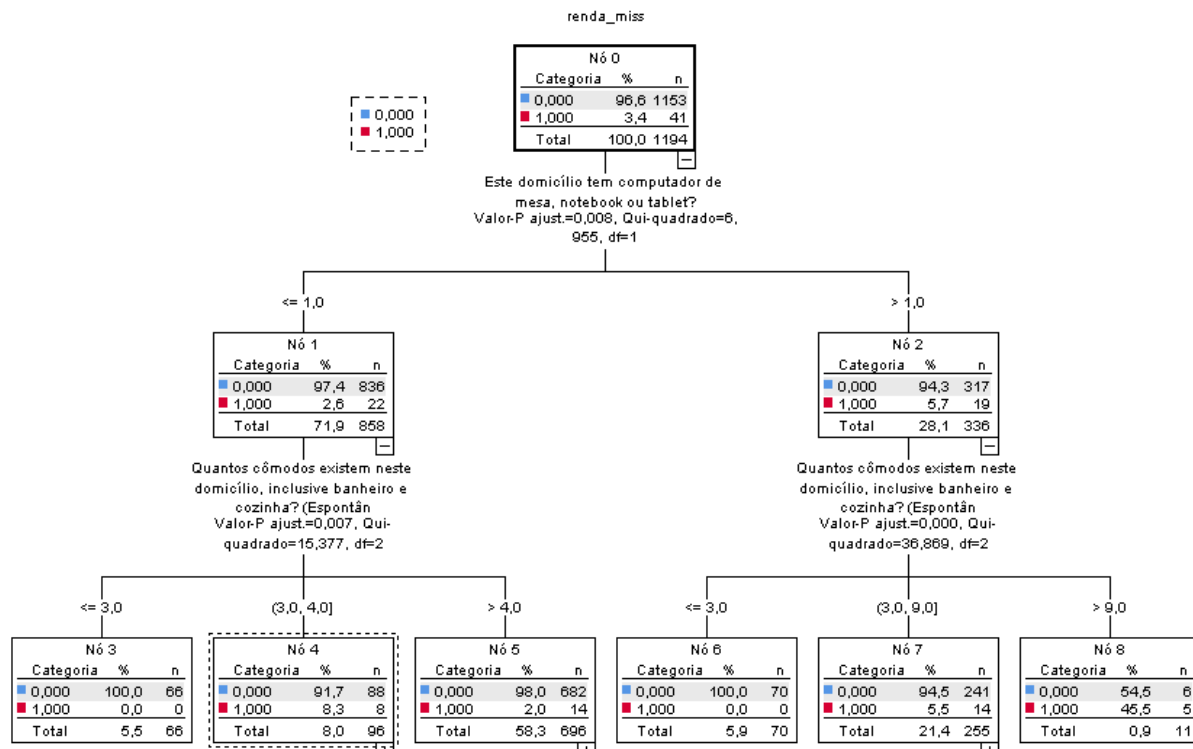


## Exemplo de Cluster Hierárquico

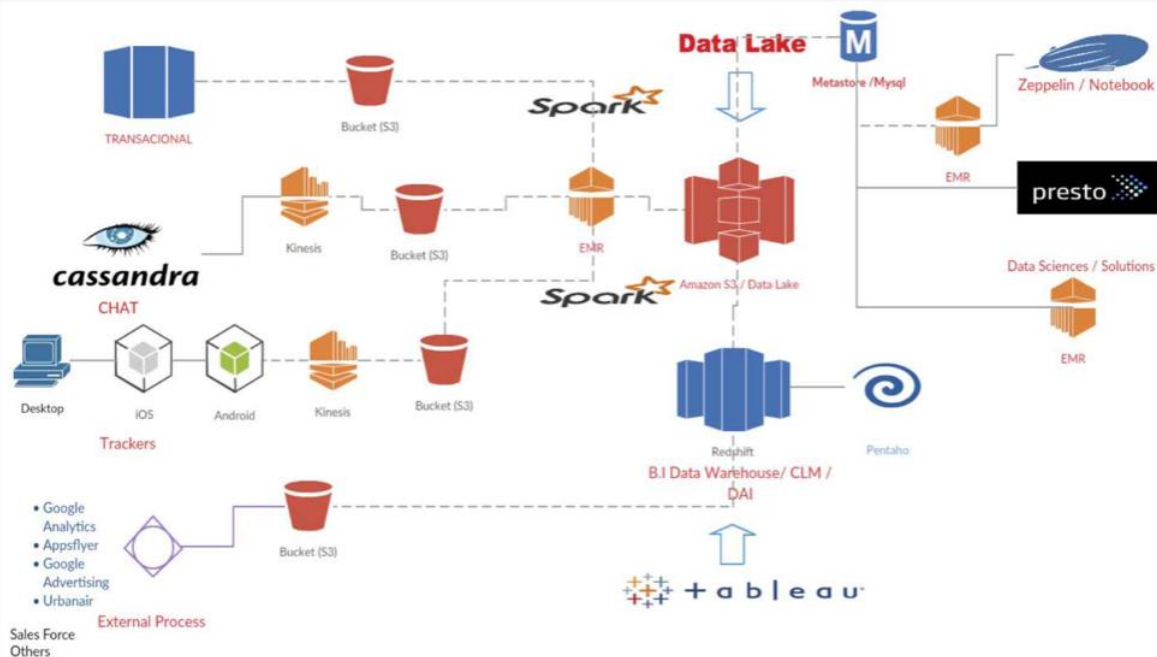


# Técnicas de Discriminação

## Exemplo de Árvore de Decisão



# Plataforma de Dados

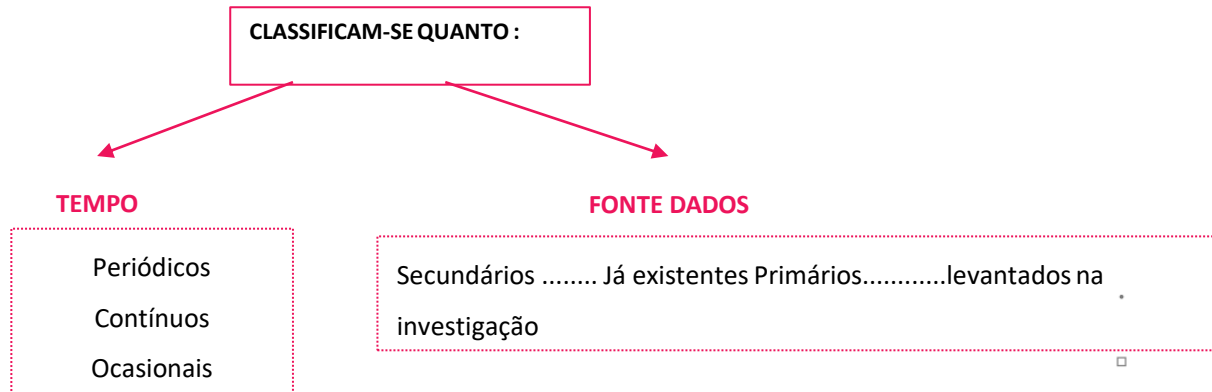




# DADOS

# LEVANTAMENTO DE DADOS

“É A OPERAÇÃO DE COLETA PARA DESCRIÇÃO E/OU ANÁLISE DAS CARACTERÍSTICAS DE UMA POPULAÇÃO”



## LEVANTAMENTO DE DADOS

Exemplos de dados secundários do IBGE :

- Pesquisa Mensal de Emprego.
- Pesquisa Industrial Mensal de Empregos e Salários.
- Pesquisa Mensal de Comércio.
- Pesquisa Nacional de Saúde.
- Censo Demográfico.
- Pesquisa de Orçamentos Familiares (POF).
- Pesquisa Nacional por Amostra de Domicílios (PNAD).
- Contagem Populacional.

Link: <https://www.ibge.gov.br/>

## LEVANTAMENTO DE DADOS

- Exemplos de dados secundários da Agência Nacional de Saúde suplementar (ANS)

<http://www.ans.gov.br/anstabnet/>

- Pesquisa Mensal de Comércio
- Susep

<http://www2.susep.gov.br/menuestatistica/Autoseg/menu1.aspx>

# ESTATÍSTICA DESCRITIVA







# ESTATÍSTICA

ESTATÍSTICA  
DESCRITIVA

ORGANIZAR  
DESCREVER  
APRESENTAR

DISTRIBUIÇÕES  
TABELAS  
GRÁFICOS  
MEDIDAS



# Estatística Descritiva

Tem por objetivo organizar, descrever e apresentar os dados, de uma determinada população, em tabelas, gráficos e medidas de resumo.

# População



População

Elementos (N=8)

Variáveis (atividade física, sexo, idade, filhos ...)



Quais as ocorrências possíveis para atividade física?

Como você representaria essas ocorrências?

# Apresentação dos dados

## Arquivo

estrutura matricial : linhas e colunas

ordem	Sexo	Atividade física	Estado civil	Grupo
1	F	Sim	Solteira	1
2	M	Sim	solteiro	1
3	F	Não	Casada	2
4	M	Não	Casado	2
5	F	Não	Casada	3
6	M	Não	Casado	3
7	F	Não	Solteira	3
8	M	Não	Solteiro	3

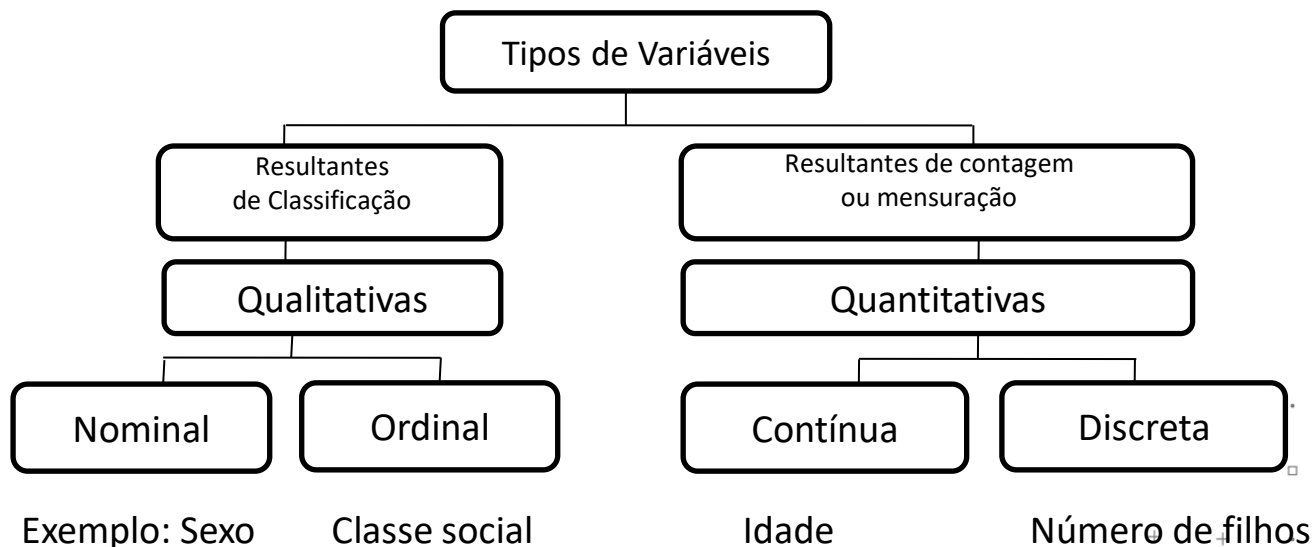
# Apresentação dos dados

## Arquivo

estrutura matricial : linhas e colunas

Grupo	Masculino	Feminino	Atividade física_Sim	Atividade física_Nao	Solteira	Casada
1	1	1	2	0	2	0
2	1	1	0	2	0	2
3	2	2	0	4	2	2

# Escala de Mensuração

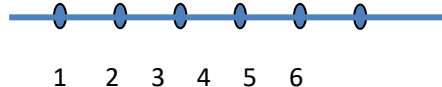


# Escala de Mensuração

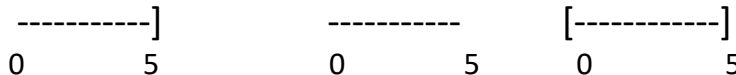
Variável qualitativa nominal: não existe nenhuma ordenação nos possíveis resultados. CATEGORIAS

Variável qualitativa ordinal: os possíveis resultados são ordenados. POSTOS

Variável quantitativa discreta: resultam de operação de contagem



Variável quantitativa contínua: possíveis resultados (valores) formam um intervalo de números reais



# Aplicando conhecimento

Classifique cada variável de acordo com seu tipo:

Variável	Ocorrência	Tipo (escala de mensuração)
Estado civil	Solteiro	Qualitativa Nominal
	Casado	
	Viúvo	
	Divorciado	
Faz atividade física	0=Não ; 1=Sim	Qualitativa Nominal
Idade (anos)	[0 – 110]	Quantitativa contínua
Anos de estudo	[0 – 99]	Quantitativa contínua



Exercitando!!!!



Base  
Cadastro

# Escala de Mensuração

Exemplo: Escala de questionário:

➡ péssimo      regular      bom      ótimo      excelente  
 ( )                      ( )                      ( )                      ( )                      ( )

Variável  
Qualitativa  
ordinal

➡ 1                      2                      3                      4                      5

Certamente  
Não compraria

Discordo  
Totalmente

Certamente  
Compraria

Concordo  
Totalmente

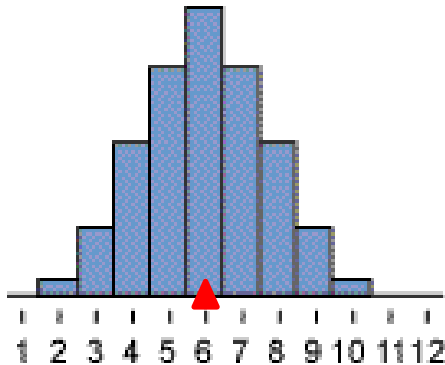
Variável  
discreta

# Medidas Resumo

São estatísticas que resumem, em um único valor, a tendência central (média, mediana, moda), a variabilidade (variância, desvio padrão) e a forma da distribuição (simétrica ou assimétrica) da variável.

# Medidas Resumo

## Distribuição simétrica



Distribuição do tempo de uso de internet (horas)

### Medidas de tendência central:

- Média
- Mediana
- Moda

Indicam o centro da distribuição de frequências ou a região de maior concentração de frequência na distribuição.

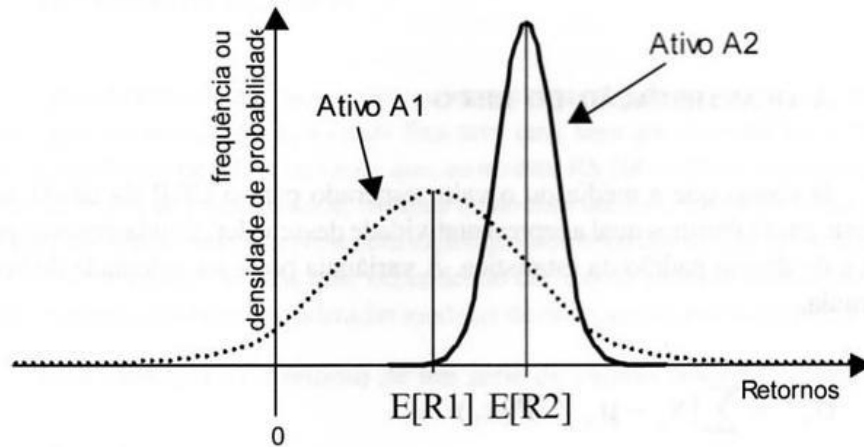
### Medidas de dispersão:

- Variância
- Desvio padrão

Indicam o grau de homogeneidade dos valores, até que ponto eles se encontram concentrados ou dispersos da média.

# Medidas Resumo

Decisão pela média



Qual ativo você escolheria para investir? Justifique sua escolha.

# Medidas Resumo

## Exemplo 2

Durante uma verificação de qualidade no conteúdo de seis recipientes de café instantâneo, foram obtidas as seguintes notas:

6,03 5,59 6,40 6,00 5,99 6,02

Qual a média e a mediana encontrada?

Média aritmética:  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \Rightarrow \bar{x} = \frac{6,03 + 5,59 + 6,40 + 6,00 + 5,99 + 6,02}{6} \Rightarrow \bar{x} = 6,00$

Mediana: 5,59 5,99 6,00 6,02 6,03 6,40

$$mediana = \frac{6,00 + 6,02}{2} = 6,01$$

# Medidas Resumo

## Exemplo 1

Durante uma verificação de qualidade no conteúdo de seis recipientes de café instantâneo, foram obtidas as seguintes notas:

6,03 5,59 6,40 6,00 5,99 6,02

Qual a média e a mediana encontrada?

$$\bar{x} = 6,00 \quad \text{mediana} = 6,01$$

Suponha que o terceiro valor tenha sido incorretamente medido e que na verdade seja de 6,04. Determine novamente a nota média e mediana.

Média aritmética:

$$\bar{x} = \frac{6,03 + 5,59 + 6,04 + 6,00 + 5,99 + 6,02}{6} = 5,95$$

Mediana:

5,59 5,99 6,00 6,02 6,03 6,04

$$\text{mediana} = \frac{6,00 + 6,02}{2} = 6,01$$

# Medidas Resumo

Comparação entre Média, Mediana e Moda

	VANTAGENS	LIMITAÇÕES	TIPO DE VARIÁVEIS
MÉDIA	Reflete todos os valores da amostra	É influenciada por valores extremos	Contínua e discreta
MEDIANA	Menos sensível a valores extremos que a média	Mais difícil de ser determinada para grande quantidade de dados	Contínua e discreta
MODA	Representa um valor típico	Não tem função em certos conjuntos de dados	Contínua, discreta, nominal e ordinal



# Medidas Resumo

## MEDIDAS DE POSIÇÃO - MÉDIA

- Média Aritmética Simples:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

- Média Aritmética Ponderada:

$$\bar{x} = \frac{\sum_{i=1}^n x_i \cdot F_i}{n}$$

- Média Geométrica (evolução):

$$Mg = \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n}$$

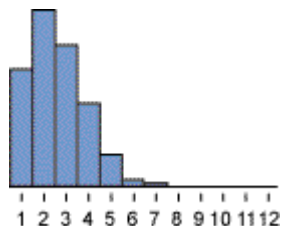
- Média Quadrática:

$$\bar{x}^2 = \frac{\sum_{i=1}^n x_i^2}{n}$$

# Medidas Resumo

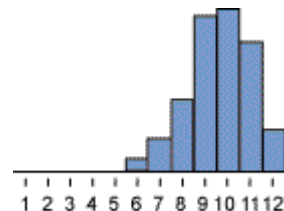
Decisão pela média ?????

Assimétrico à direita



Média > Mediana

Assimétrico à esquerda

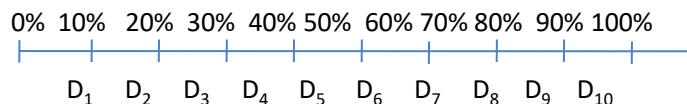


Média < Mediana

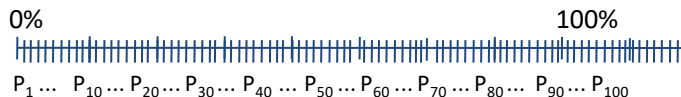
# Medidas Resumo

## • Outras Medidas de Posição

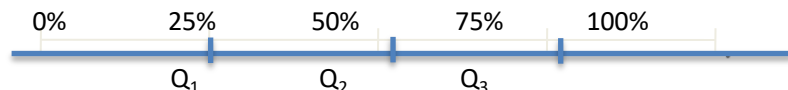
Decis: dividem um conjunto de dados em dez partes iguais.



Percentis (P): dividem a série em cem partes, de modo que p% ficam abaixo dele (P).



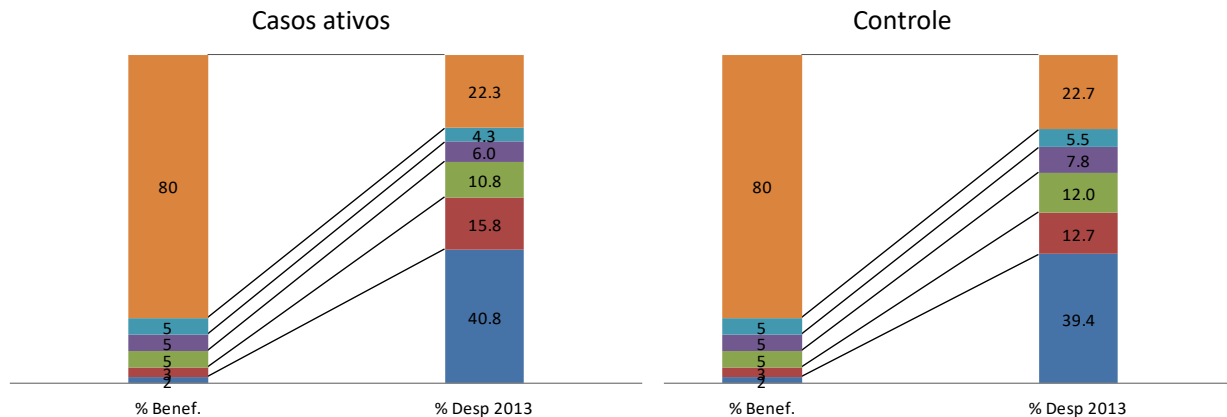
Quartis: dividem a série em quatro partes iguais.



# Medidas Resumo

## Exemplo: Despesas

Gráfico de Pareto de despesas



# Medidas **Resumo**

## Economia nacional

### São Paulo, Rio e Brasília respondem por 21% do PIB brasileiro

*Andrea Bruxellas*

*Direto do Rio de Janeiro*

*Especial para o Terra*

Os municípios de São Paulo, Rio de Janeiro e Brasília respondiam por 21% do Produto Interno Bruto brasileiro em 2007. Segundo dados divulgados pelo Instituto Brasileiro de Geografia e Estatística (IBGE) nesta quarta-feira, a capital paulista responde pela maior fatia do PIB brasileiro, gerando 12% de toda riqueza produzida no País, seguida do Rio de Janeiro (5,2%), Brasília (3,8%), Belo Horizonte (1,4%) e Curitiba (1,4%).

"Com os dados de 2007 a gente pode notar uma estabilidade na série. Ou seja, na série inteira a gente vê que a renda ainda está muito concentrada em alguns municípios e isso é bastante estável. Nas cinco principais cidade a gente tem um quarto do PIB. Tirando essas cidades, a economia esta concentrada em 50 cidades que geram 50% da riqueza do País", disse a coordenadora do IBGE Sheila Cristina Zani.

Já os menores PIB do Brasil foram verificados em Santo Antônio dos Milagres (PI), São Miguel da Baixa Grande (PI), Areia de Barúnas (PB), São Luís do Piauí (PI) e Olho D'Água do Piauí (PI). Segundo o IBGE, a soma dos PIB destes cinco municípios representava 0,001% da riqueza produzida em todo País em 2007.

# Medidas Resumo

- Medidas de Dispersão

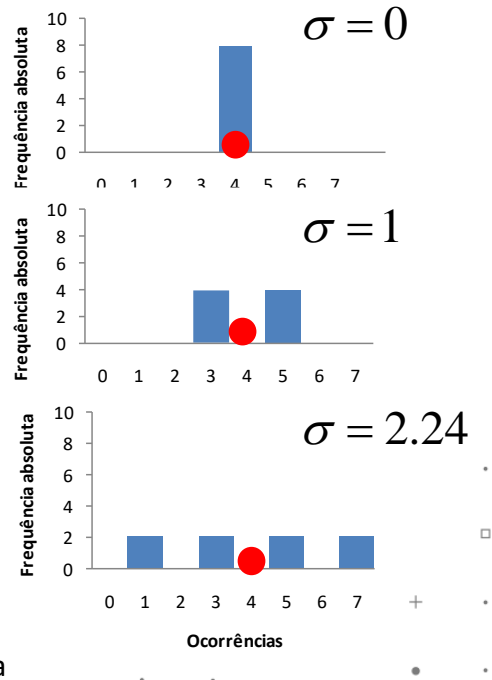
Exemplo 8:

A: 4, 4, 4, 4, 4, 4, 4, 4, 4

B: 3, 3, 3, 3, 5, 5, 5, 5

C: 1, 1, 3, 3, 5, 5, 7, 7

Qual o desvio padrão?



# Medidas de Dispersão

Medidas de Dispersão: variância e desvio padrão

Exemplo C

X	Média	(X-Média)	(X-Média) <sup>2</sup>
1	4	-3	9
1	4	-3	9
3	4	-1	1
3	4	-1	1
5	4	1	1
5	4	1	1
7	4	3	9
7	4	3	9
Soma	-	0	40

Variância:

$$\sigma^2 = \frac{40}{8} = 5$$

Desvio padrão:

$$\sigma = \sqrt{\sigma^2} = \sqrt{5} = 2.24$$

# Medidas de Dispersão

O quanto os pontos (dados) estão distantes da média (ponto central)

➤ **variância da população**

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

➤ **variância da amostra**

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



# Medidas Resumo

EXEMPLO

## Controle estatístico do processo

O **Controle Estatístico de Processos** (CEP) é uma ferramenta da qualidade utilizada nos processos produtivos (e de serviços) com objetivo de fornecer informações para um diagnóstico mais eficaz na prevenção e detecção de defeitos/problemas nos processos avaliados e, conseqüentemente, auxilia no aumento da produtividade/resultados da empresa, evitando desperdícios de matéria-prima, insumos, produtos etc.

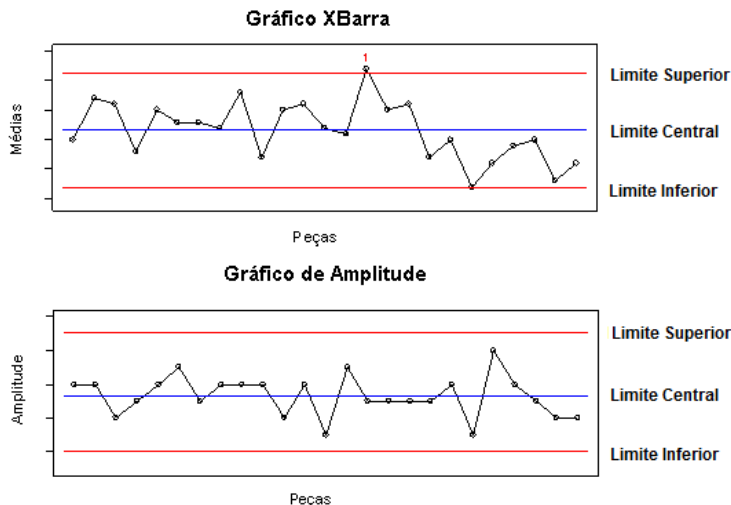
(Fonte: [https://pt.wikipedia.org/wiki/Controle\\_estat%C3%ADstico\\_de\\_processos](https://pt.wikipedia.org/wiki/Controle_estat%C3%ADstico_de_processos))

Exemplo: Fábrica de Café em Pó

# Medidas Resumo

## Controle estatístico do processo

### Gráfico de controle



“Mostrar evidências de que um processo esteja operando em estado de controle estatístico e dar sinais de presença de causas especiais de variação para que medidas corretivas apropriadas sejam aplicadas”.

“Manter o estado de controle estatístico estendendo a função dos limites de controle como base de decisões”.

“Apresentar informações para que sejam tomadas ações gerenciais de melhoria dos processos”.

O gráfico é construído a partir das medidas estatística como:

[Média aritmética.](#)  
[Desvio padrão.](#)  
[Média das médias.](#)  
[Somatórios](#) etc.

Fonte: <http://www.portalação.com.br/controle-estatistico-do-processo/graficos-ou-cartas-de-controle>



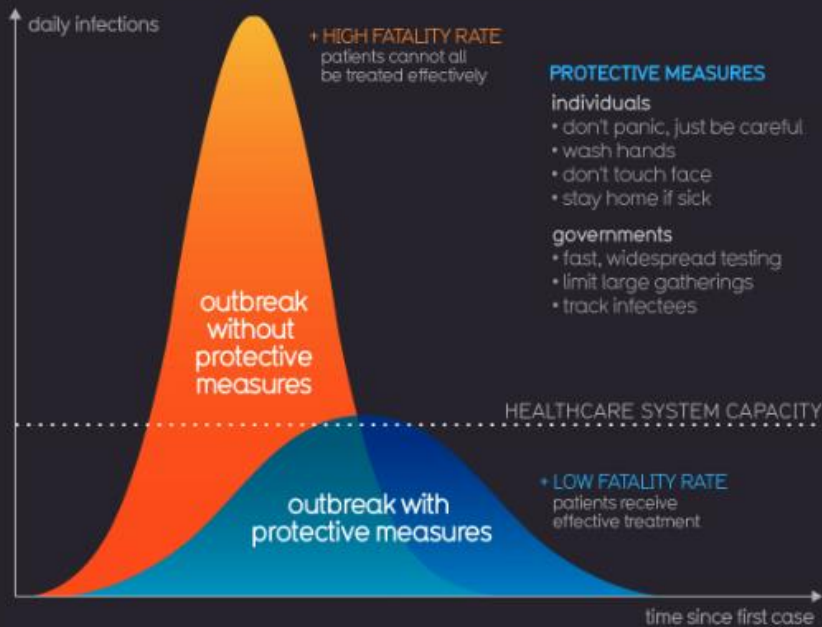
# MEDIDAS DE ASSIMETRIA



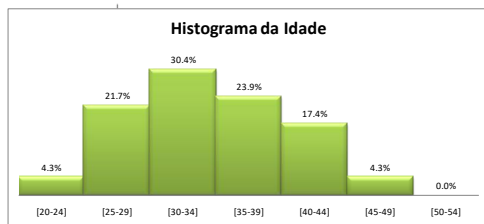
# Medidas Resumo

## Flattening the Curve

Fast, intelligent action slows pandemic effects, stops the overwhelm of healthcare systems

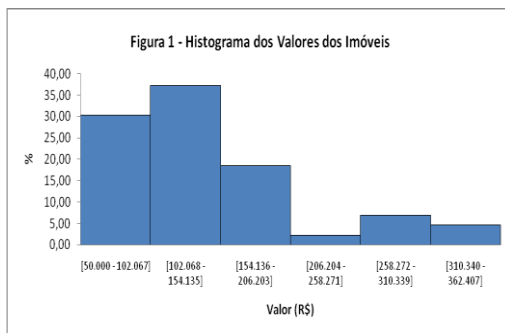


## Exemplo de estatística descritiva



<i>idade</i>	
Média	34.6
Erro padrão	1.1
Mediana	34.5
Modo	26
Desvio padrão	6.74
Variância da amostra	45.39
Curtose	-0.54
Assimetria	-0.07
Intervalo	28
Mínimo	20
Máximo	48
Soma	1245
Contagem	36

## Exemplo de estatística descritiva



Fonte: Estudo de Caso no Centro de Florianópolis

Valor (R\$)	
Média	144618.3
Erro padrão	10992.8
Mediana	120000.0
Modo	110000.0
Desvio padrão	72084.7
Variância da amostra	5196201097.5
Curtose	1.4
Assimetria	1.4
Intervalo	312400.0
Mínimo	50000.0
Máximo	362400.0
Soma	6218585.0
Contagem	43

# Medidas de Assimetria

As medidas de assimetria referem-se à forma da curva que representa a distribuição de frequência. A assimetria é o afastamento da simetria de uma frequência.

- Curvas de frequência simétrica ou em forma de sino: caracterizam-se pelo fato das observações equidistantes do ponto central terem a mesma frequência (curva normal)
- Curvas de frequência moderadamente assimétricas ou desviadas: a cauda de um lado da ordenada máxima é mais longa do que do outro. Se o ramo mais alongado fica à direita, a curva é dita de assimetria positiva, enquanto que, se ocorre o inverso, diz-se que a curva é de assimetria negativa.

# Medidas de Assimetria

## Coeficientes de Assimetria (Skewness)

$$\rightarrow As = \frac{m^3}{\sigma^3} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2]^{\frac{3}{2}}}$$

$As=0 \rightarrow$  simétrica

$As>0 \rightarrow$  assimétrica positiva

$As<0 \rightarrow$  assimétrica negativa

## Índice de Assimetria (Pearson)

$$\rightarrow A = \frac{\text{média} - \text{moda}}{\text{desvio padrão}}$$

$|A| < 0,15 \rightarrow$  simétrica

$0,15 < |A| < 1 \rightarrow$  assimetria moderada

$|A| > 1 \rightarrow$  assimetria forte



# Medidas de Assimetria

## Curtose (Kurtosis)

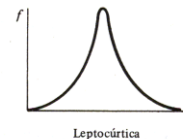
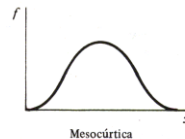
$$As = \frac{m^4}{\sigma^4} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{\left[ \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^2}$$

➤ Curtose: grau de achatamento em relação a uma curva Normal

➤ Leptocúrtica (afilado) →  $K > 3$

➤ Mesocúrtica →  $K = 3$

➤ Platicúrtica (achatado) →  $K < 3$



	MEAN		
1	$\frac{\sum x}{n}$		
		VARIANCE	
2	$\frac{\sum x^2}{n}$	$\frac{\sum (x - \mu)^2}{n}$	
			SKEWNESS
3	$\frac{\sum x^3}{n}$	$\frac{\sum (x - \mu)^3}{n}$	$\frac{1}{n} \frac{\sum (x - \mu)^3}{\sigma^3}$
			KURTOSIS
4	$\frac{\sum x^4}{n}$	$\frac{\sum (x - \mu)^4}{n}$	$\frac{1}{n} \frac{\sum (x - \mu)^4}{\sigma^4}$

# Medidas Resumo

## Outras Medidas de Dispersão

- Coeficiente de Variação
- Amplitude
- Amplitude Inter-Quartílica

# Medidas Resumo

## Outras Medidas de Dispersão

### Coeficiente de variação (CV)

É o quociente entre o desvio padrão e a média.

$$CV = \frac{\sigma}{\bar{X}}$$

Vantagem: caracterizar a dispersão dos dados em termos relativos a seu valor médio.

# Medidas Resumo

Qual o coeficiente de variação?

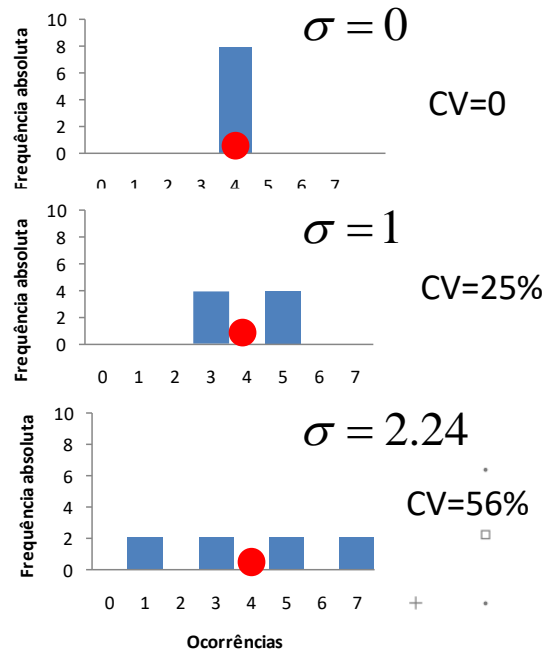
## • Medidas de Dispersão

Exemplo 8:

A: 4, 4, 4, 4, 4, 4, 4, 4, 4

B: 3, 3, 3, 3, 5, 5, 5, 5

C: 1, 1, 3, 3, 5, 5, 7, 7



● Média

# Medidas Resumo

## Outras Medidas de Dispersão

### Amplitude

É definida como a diferença entre o maior e o menor valor de um conjunto de dados.

Fortemente relacionado com a dispersão dos dados.

A amplitude pode levar a erros de avaliação, pois não representa o conjunto dos dados. Muitas vezes reflete muito mal a dispersão dos mesmos.

# Medidas Resumo

- Outras Medidas de Dispersão

## Amplitude Inter-quartílica

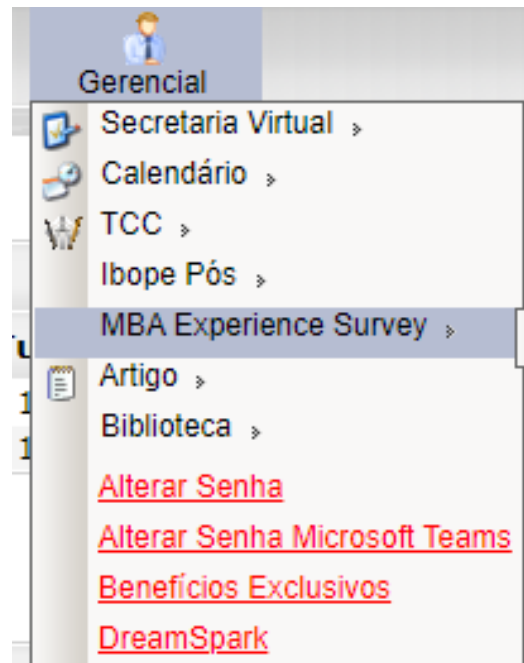
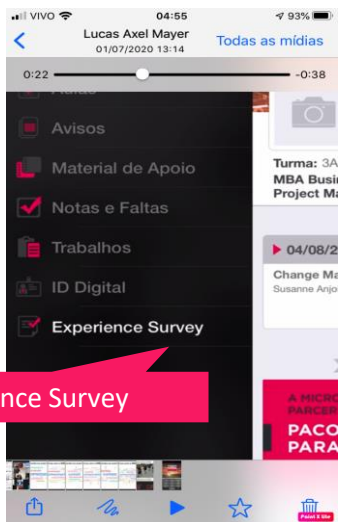
É a diferença entre o terceiro e o primeiro quartil ( $Q3 - Q1$ ).

Usada em análise exploratória de dados – gráficos Box Plot.

# O que você achou da aula de hoje?

Pelo aplicativo da FIAP

(Entrar no FIAPP, e no menu clicar em Experience Survey)







# OBRIGADA



/ Regina T. I. Bernal

FIAP

Copyright © 2023 | Professora Dra. Regina Tomie Ivata Bernal  
Todos os direitos reservados. Reprodução ou divulgação total ou parcial deste documento, é expressamente proibido sem consentimento formal, por escrito, do professor/autor.

FIAP