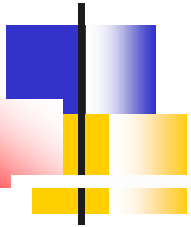


# Ciência de Dados

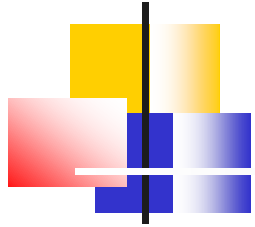
## Exploração de dados PARTE II



Profa. Roseli Ap. Francelin  
Romero – SCC

Material elaborado pelo Prof. Dr.  
André C. P. L. F. de Carvalho  
Dr. Isvani Frias-Blanco e  
complementado por Roseli Romero  
ICMC-USP

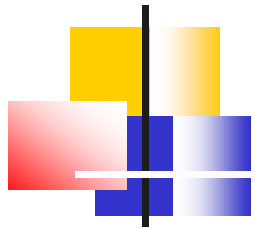




# Tópicos

---

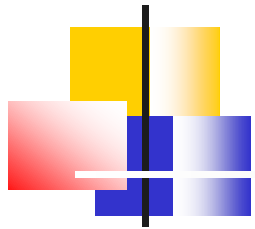
- Exploração de dados
  - Dados multivariados:
    - Matriz de Covariância
    - Coeficiente de Pearson
    - Matriz de Correlação



# Dados multivariados

---

- Possuem mais de um atributo
  - Cada atributo é uma variável
- Medidas de localização (tendência central)
  - Podem ser obtidas calculando medida de localização de **cada atributo** separadamente
  - Ex.: média, mediana, ...
    - Média dos objetos de um conjunto de dados com  $m$  atributos é dada por:  $\bar{x} = (\bar{x}_1, \dots, \bar{x}_m)$



# Dados multivariados

---

- Medidas de espalhamento (dispersão)
  - Podem ser calculadas para **cada atributo** independentemente dos demais
    - Usando qualquer medida de espalhamento
      - Intervalo, variância, desvio padrão
  - Para dados multivariados numéricos deve-se usar uma matriz de covariância
    - Cada elemento da matriz é a covariância entre dois atributos



# Dados multivariados

---

- Cálculo de cada elemento  $s_{ij}$  de uma matriz de covariância  $S$  para um conjunto de  $n$  objetos

$$s_{ij} = \text{covariância}(x_i, x_j) = \frac{1}{n-1} \sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j)$$

Onde:

$\bar{x}_i$  : Valor médio do  $i$ -ésimo atributo

$x_{ki}$  : Valor do  $i$ -ésimo atributo para o  $k$ -ésimo objeto

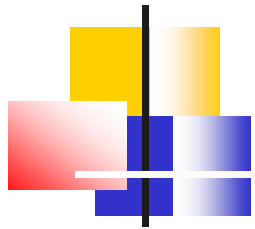
- Obs: covariância  $(x_i, x_i) = \text{variância}(x_i)$ 
  - Matriz de covariância tem em sua diagonal as variâncias dos atributos



# Dados multivariados

---

- Covariância de dois atributos
  - Mede o grau com que os atributos variam juntos
    - Valor próximo de 0:
      - Atributos não têm um relacionamento
    - Valor positivo:
      - Atributos diretamente relacionados
        - Quando o valor de um atributo aumenta, o do outro também aumenta
      - Valor negativo:
        - Atributos inversamente relacionados
    - Valor depende da magnitude dos atributos



# Dados multivariados

---

- Covariância de dois atributos
  - É difícil avaliar o relacionamento entre dois atributos olhando apenas a covariância
    - Sofre influência da faixa de valores dos atributos
    - Correlação entre dois atributos ilustra mais claramente a força da relação entre eles
      - Mais popular que covariância
      - Elimina influência da faixa de valores



# Dados Multivariados

---

- Matriz de Correlação
- Correlação Linear: Coeficiente de correlação de Pearson
- Correlação Não Linear:
  - Correlação  $\eta$  (eta)
  - Kendall
  - Spearman





# Matriz de Correlação

---

- Coeficiente de Pearson

$$\rho = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- Coeficiente de Spearman

$$\rho = 1 - \frac{6 \sum_{i=1}^n d_i^2}{(n^3 - n)}$$

onde  $d_i$  é a diferença entre cada posição de  $x$  e  $y$ .

- Coeficiente de Kendall

$$\tau = \frac{2(n_c - n_d)}{n(n-1)}$$



# Matriz de Correlação

---

- A matriz de correlação é simétrica.
- O valor 1 significa que as duas variáveis são exatamente correlacionadas. É o caso de uma relação linear entre as duas variáveis;
- O valor negativo, significa que elas são inversamente correlacionadas.
- Não necessariamente informa o tipo de relação que existe entre as variáveis.

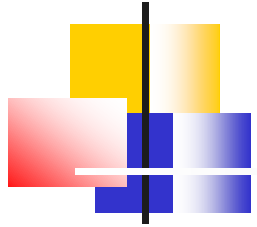
# Matriz de Correlação

Tabela 2. Matriz de correlação das variáveis de qualidade de água no médio Rio Pomba

Varial	OD	DBO	Temp	CE	Alcalin	DQO	PT	Norg	NH <sub>4</sub>	NTK	NT	ST	SST	SDT	SIS
OD	1,00														
DBO	<b>-0,58</b>	1,00													
Temp	<b>-0,57</b>	0,18	1,00												
CE	0,02	-0,24	<b>0,61</b>	1,00											
Alcalin	-0,23	-0,34	<b>0,63</b>	<b>0,77</b>	1,00										
DQO	<b>-0,70</b>	0,32	<b>0,64</b>	0,12	0,38	1,00									
PT	-0,09	0,31	0,03	0,21	-0,11	-0,10	1,00								
Norg	0,16	0,23	<b>-0,62</b>	-0,49	<b>-0,71</b>	-0,49	0,28	1,00							
NH <sub>4</sub>	<b>-0,76</b>	<b>0,53</b>	0,38	-0,01	0,12	0,31	0,48	0,08	1,00						
NTK	0,08	0,28	<b>-0,57</b>	-0,49	<b>-0,69</b>	-0,45	0,33	<b>0,99</b>	0,18	1,00					
NT	0,15	0,27	<b>-0,55</b>	-0,34	<b>-0,62</b>	<b>-0,50</b>	0,35	<b>0,97</b>	0,13	<b>0,97</b>	1,00				
ST	0,24	-0,06	-0,23	0,33	-0,02	<b>-0,61</b>	0,47	0,29	0,08	0,29	0,37	1,00			
SST	-0,01	0,35	-0,39	-0,14	-0,28	-0,43	<b>0,60</b>	<b>0,61</b>	0,37	<b>0,64</b>	<b>0,63</b>	<b>0,63</b>	1,00		
SDT	0,30	-0,38	0,04	<b>0,55</b>	0,20	-0,42	0,09	-0,15	-0,22	-0,17	-0,08	<b>0,75</b>	-0,04	1,00	
SIS	0,01	0,26	-0,47	-0,19	-0,37	-0,43	<b>0,54</b>	<b>0,61</b>	0,32	<b>0,63</b>	<b>0,63</b>	<b>0,75</b>	<b>0,84</b>	0,24	1,00

OD (mg L<sup>-1</sup>) – oxigênio dissolvido; DBO (mg L<sup>-1</sup>) – demanda bioquímica de oxigênio; Temp (°C) – temperatura; CE (μS cm<sup>-1</sup>) – condutividade elétrica; Alcalin (mg L<sup>-1</sup> CaCO<sub>3</sub>) – alcalinidade; DQO (mg L<sup>-1</sup>) – demanda química de oxigênio; PT (mg L<sup>-1</sup> – PO<sub>4</sub>) – fósforo total; Norg (mg L<sup>-1</sup> – N) – nitrogênio orgânico; NH<sub>4</sub> (mg L<sup>-1</sup> – N) – nitrogênio amoniacal; NTK (mg L<sup>-1</sup> – N) – nitrogênio Kjeldahl; NT (mg L<sup>-1</sup> – N) – nitrogênio total; ST (mg L<sup>-1</sup>) – sólidos totais; SST (mg L<sup>-1</sup>) – sólidos suspensos totais; SDT (mg L<sup>-1</sup>) – sólidos dissolvidos totais; SIS (mg L<sup>-1</sup>) – sólidos inorgânicos suspensos

Fonte: Scielo



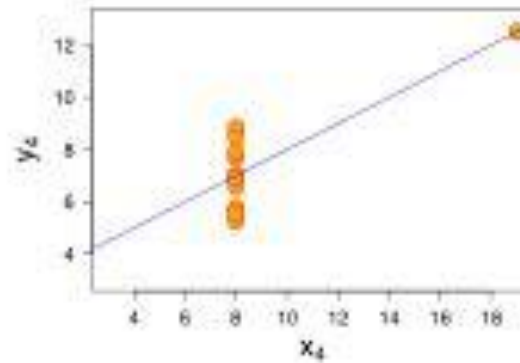
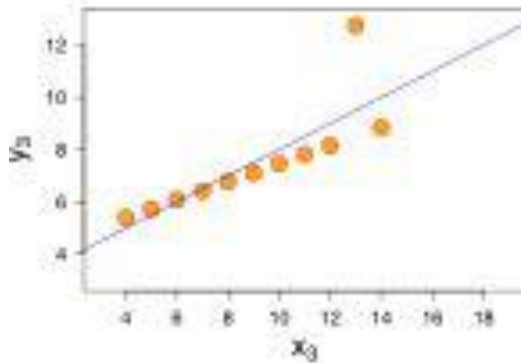
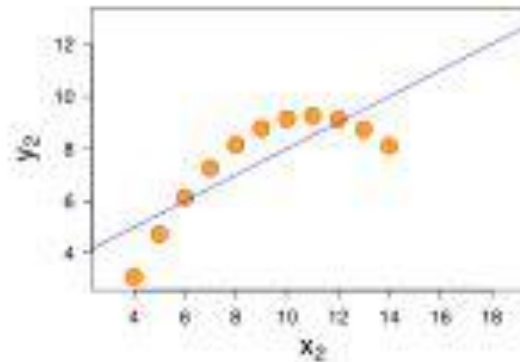
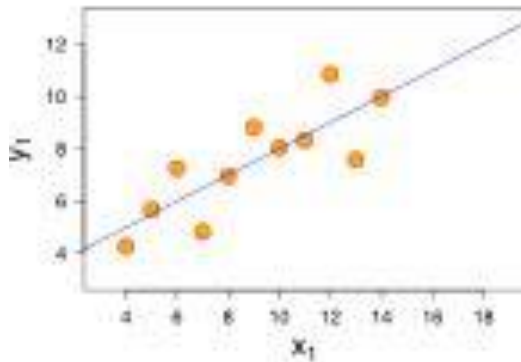
# Exercício

---

- Calcular a matriz de covariância para o conjunto de dados:

Peso	Altura	Temperatura
73	170	37
67	165	38
90	190	34
49	152	31

# Dados Multivariados



Todos  
correlação =  
0,816

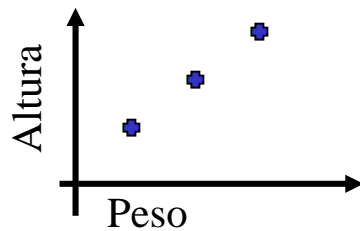
Fonte: Wikipedia



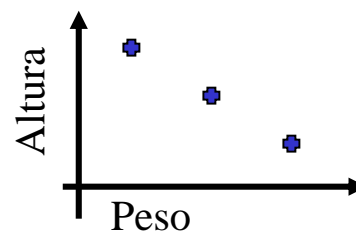
# Exemplo

---

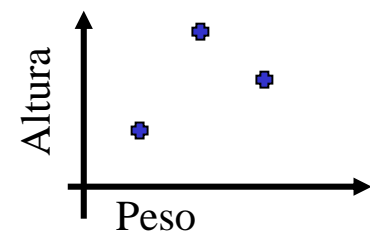
Peso	Altura
60	170
70	180
80	190

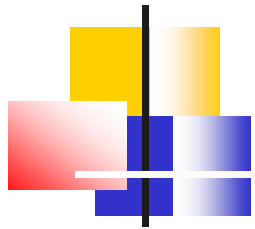


Peso	Altura
60	190
70	180
80	170



Peso	Altura
60	170
70	190
80	180

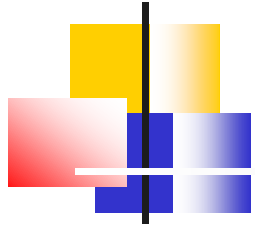




# Perguntas

---





# Outras formas de sumarizar dados

---

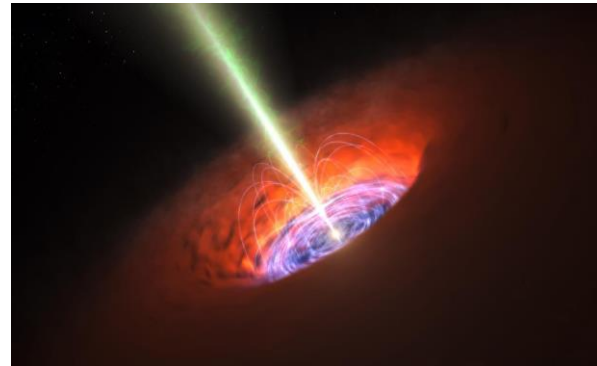
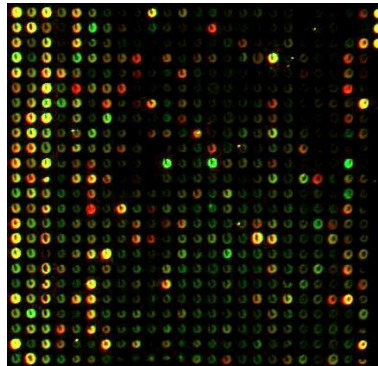
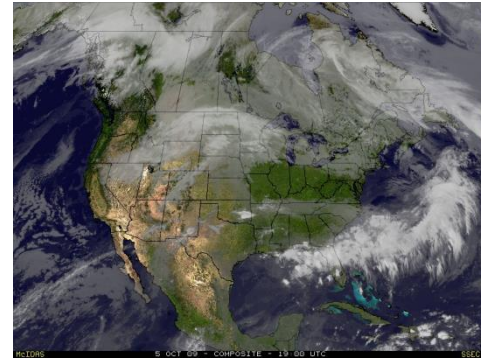
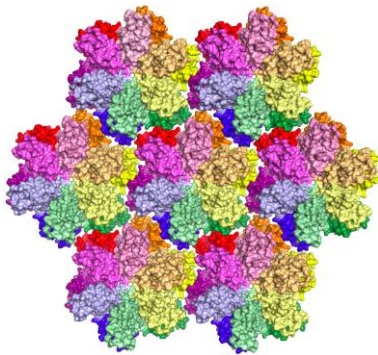
- Visualização gráfica
  - Em vários casos, facilita compreensão de padrões mais complexos nos dados
- Exemplos simples
  - Histograma
  - Diagrama de torta
  - *Scatter plot*
  - Faces de Chernoff

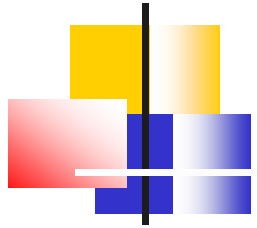




# Exemplos

---





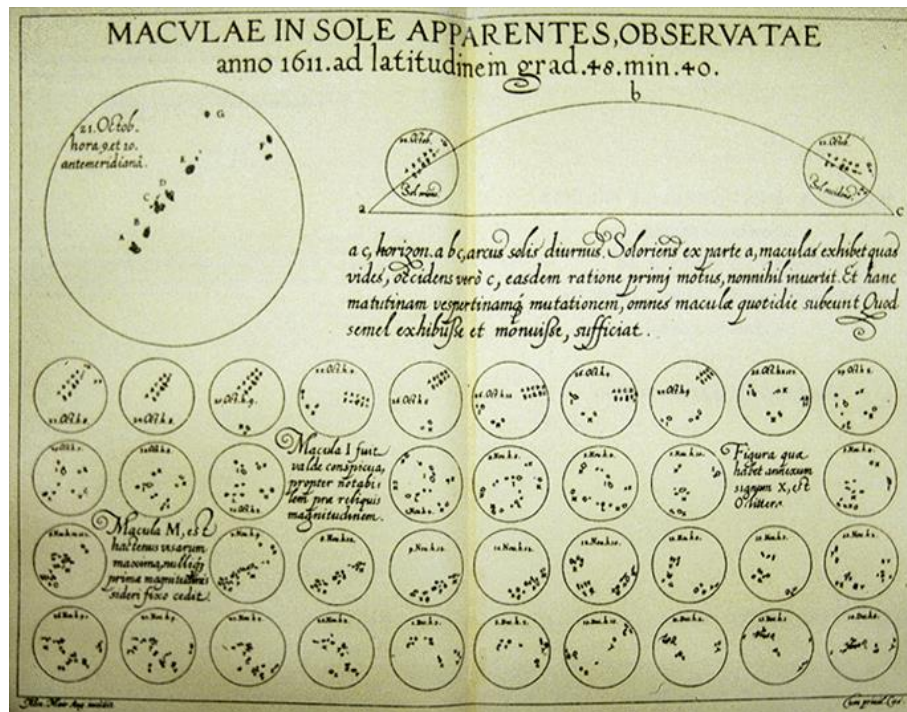
# Visualização

---

- Visualização tem um papel importante em análise de dados
  - Uma das técnicas mais poderosas para exploração dos dados
    - Facilita visualização de dados e resultados
    - *Visual data mining*
      - Usa técnicas de visualização em mineração de dados
      - Importante área de CD

# Um dos primeiros

## ■ Mapa solar de Galileu



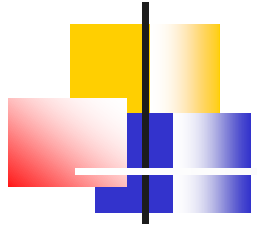
# Outro dos primeiros usos



Napoleão

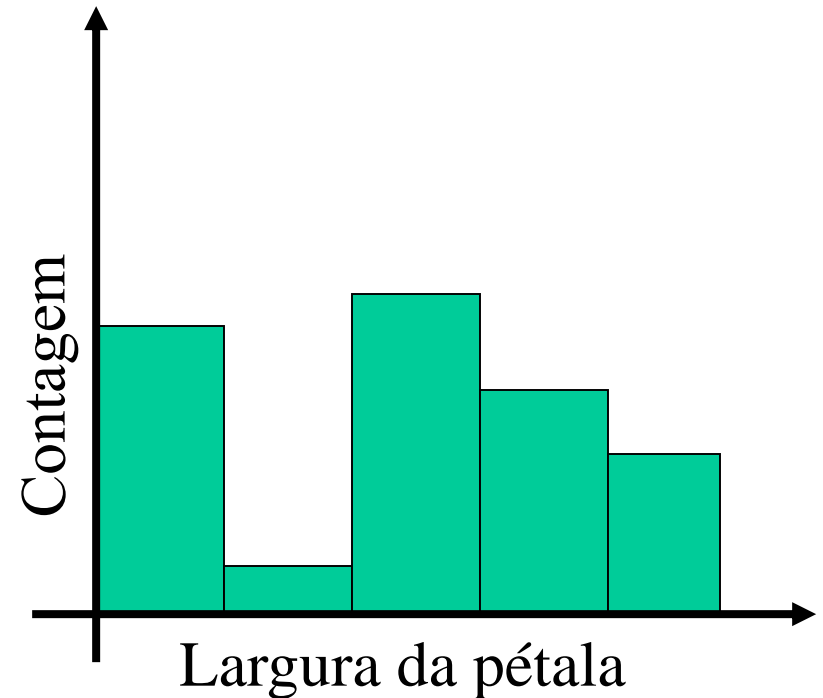
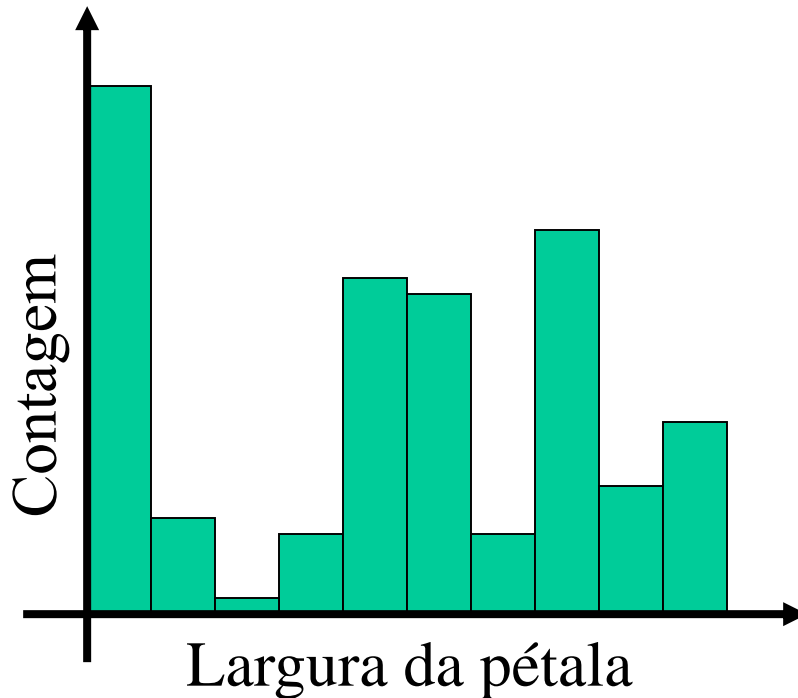
Exército francês, comandado por Napoleão  
Invade a Rússia, em 1812

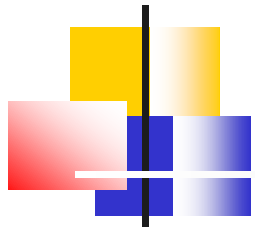




# Histogramas

- Conjunto de dados Iris
  - Largura das pétalas usando 10 e 5 cestas





# Diagrama de torta

---

- Frequências relativas podem ser vistas no diagrama circular

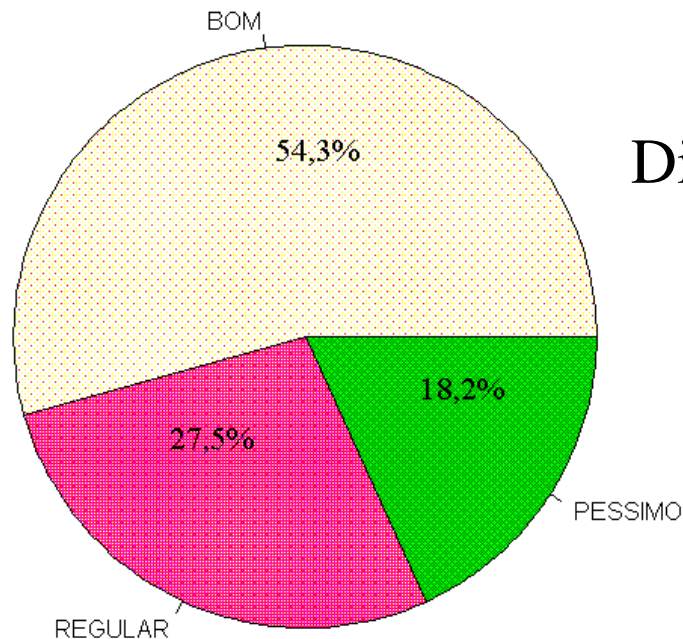
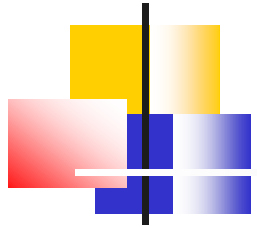


Diagrama de torta (pizza)





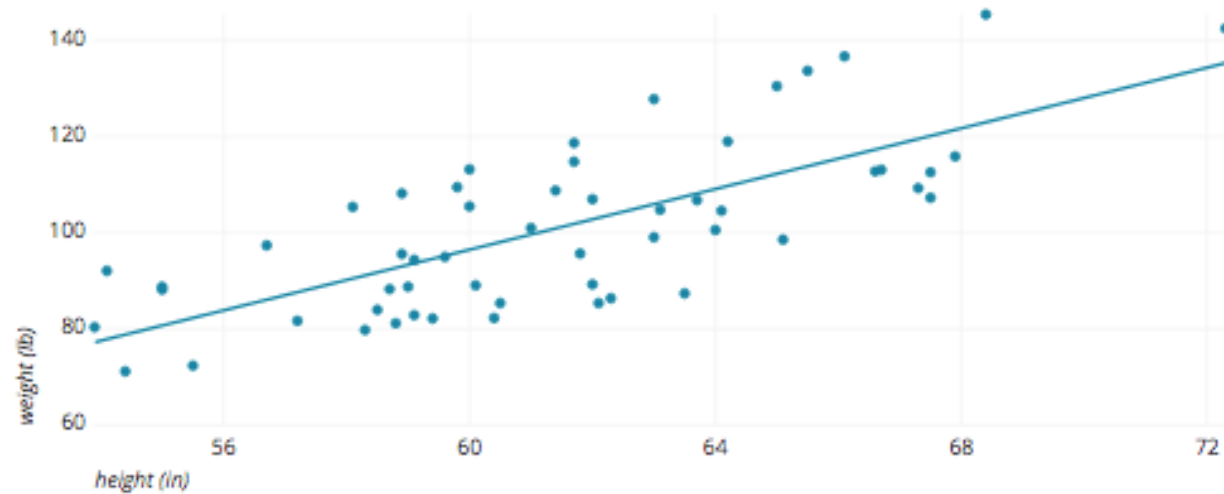
# Scatter Plot

---

- Usado para ilustrar graficamente **Correlação Linear** entre dois atributos
- Cada objeto é associado a uma posição em um gráfico
  - Valores dos atributos definem sua posição
  - Valores podem ser inteiros ou reais
- Matrizes de scatter plot resumem relação para vários pares de atributos

# Relação Linear

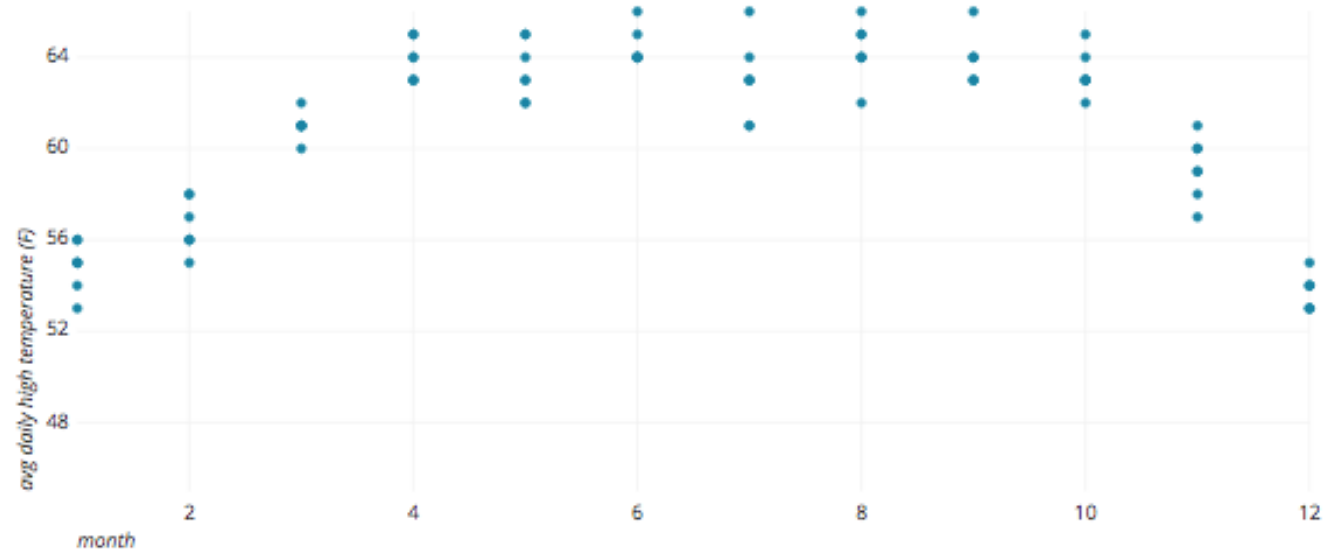
Weight and Height of Children





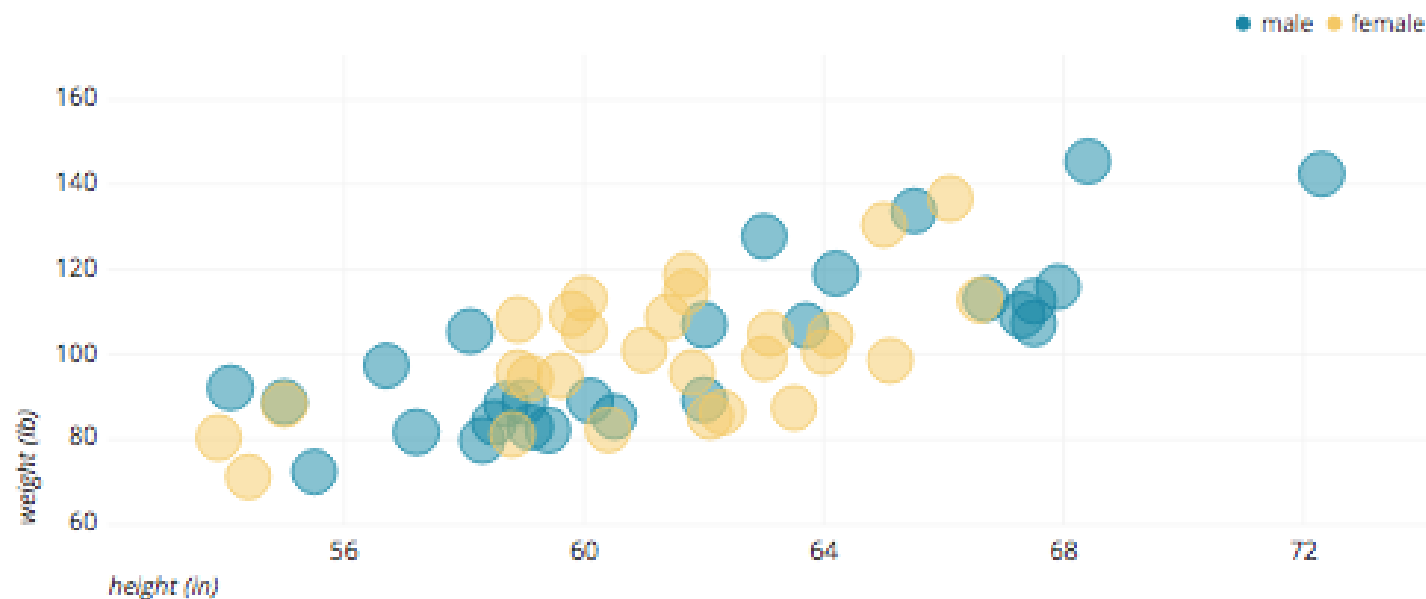
# Relação Não-Linear

Average Daily High Temperature by Month



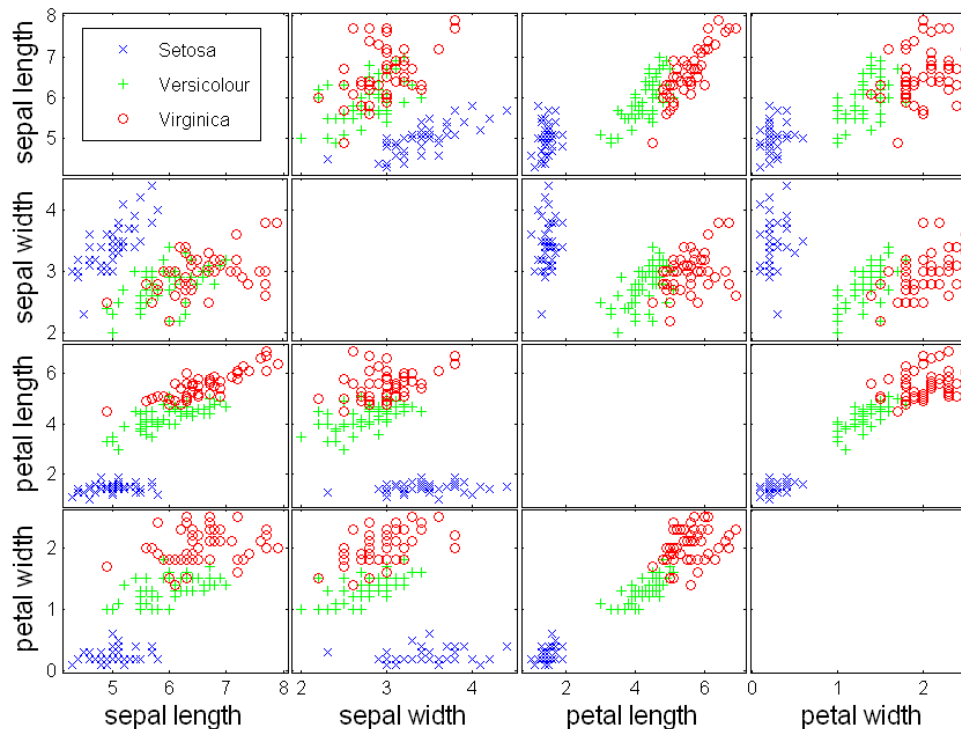
# Pode adicionar mais info

Weight and Height of Children by Gender

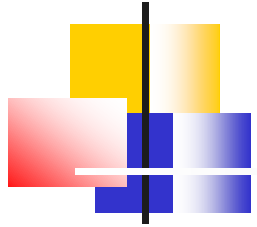


# Scatter Plot

- Matriz para atributos do conjunto iris



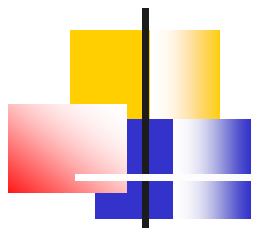
Diferentes classes  
são indicadas por  
cores diferentes



# Faces de Chernoff

---

- Criado por Herman Chernoff
- Mapeia os valores dos atributos para imagens familiares para seres humanos: faces
  - Cada objeto é representado por uma face
  - Cada atributo é associado a uma característica específica da face
- Baseia-se na habilidade humana de distinguir faces



# Faces de Chernoff

---



Setosa

1

2

3

4

5



Versicolour

51

52

53

54

55



Virginica

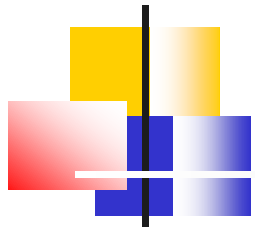
101

102

103

104

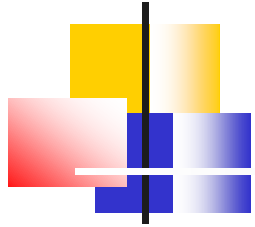
105



# Exercício

- Representar os dados a seguir usando faces de Chernoff

Febre	Idade	Batimento	Dor	Diagnóstico
sim	23	elevado	sim	doente
não	9	baixo	não	saudável
sim	61	elevado	não	saudável
sim	32	baixo	sim	doente
sim	21	elevado	sim	saudável
não	48	elevado	sim	doente



# Considerações Finais

---

- Caracterização de dados
  - Objetos e atributos
  - Tipos de dados
- Exploração de dados
  - Dados univariados
  - Medidas de localidade, espalhamento e distribuição
  - Dados multivariados
  - Visualização de dados e de resultados