

Mineração de Dados

Aula 1 – parte 2

Especialização em Ciência de Dados e suas Aplicações

- Um atributo é uma propriedade de um objeto
- Exemplos: cor de olho, temperatura, etc.
 - Atributo é também conhecido como **variável, campo, característica, feature**

Atributos

Objetos

<i>Tid</i>	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Medida de similaridade

- Medida numérica de quão parecidos dois objetos de dados são.
- É maior quando são mais parecidos
- Geralmente estão no intervalo $[0,1]$

- Medida de dissimilaridade

- Medida numérica de quão diferente dois objetos são
- Dissimilaridade mínima é geralmente 0

- Proximidade refere a similaridade ou dissimilaridade

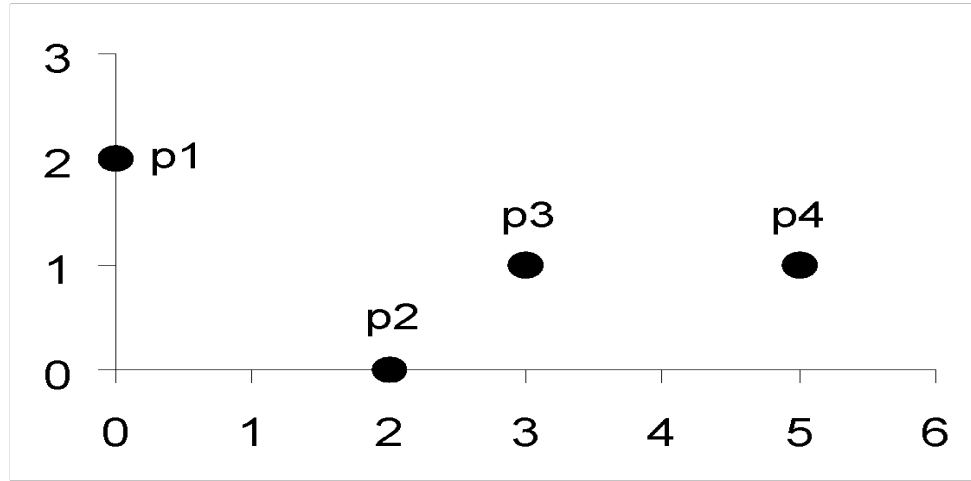
- Distância Euclideana

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$$

onde n é o número de dimensões (atributos) e x_k e y_k são os k^{th} atributos ou objetos de dados \mathbf{x} e \mathbf{y} .

- Padronização é necessária, se a escala difere

Distância Euclideana



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Matriz de distância

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}$$

$r = 1$. [Distância Manhattan](#) (City block, taxicab, L_1 norm).

n é o número de dimensões (atributos) e x_k e y_k são os $k^{\text{éssimos}}$ atributos ou objetos de dados \mathbf{x} e \mathbf{y} .

Manhattan

point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

L1	p1	p2	p3	p4
p1	0	4	4	6
p2	4	0	2	4
p3	4	2	0	2
p4	6	4	2	0

Euclideana

L2	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Matriz de
distância

- Situação comum: objetos p e q possuem apenas atributos binários
- Computar a similaridade assim:
 f_{01} = # de atributos onde p era 0 e q é 1
 f_{10} = # de atributos onde p era 1 e q é 0
 f_{00} = # de atributos onde p era 0 e q é 0
 f_{11} = # de atributos onde p era 1 e q é 1

Simple Matching (**SMC**) and Coeficiente de Jaccard (**J**)

$$\begin{aligned}\text{SMC} &= \text{número de correspondências "11" e "00" / número de atributos} \\ &= (f_{11} + f_{00}) / (f_{01} + f_{10} + f_{11} + f_{00})\end{aligned}$$

$$\begin{aligned}\text{J} &= \text{número de correspondências "11" / número de atributos não zero} \\ &= (f_{11}) / (f_{01} + f_{10} + f_{11})\end{aligned}$$

SMC vs Jaccard: Exemplo

$$\mathbf{x} = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0$$

$$\mathbf{y} = 0\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1$$

$$f_{01} = 2$$

$$f_{10} = 1$$

$$f_{00} = 7$$

$$f_{11} = 0$$

$$\begin{aligned}\text{SMC} &= (f_{11} + f_{00}) / (f_{01} + f_{10} + f_{11} + f_{00}) \\ &= (0+7) / (2+1+0+7) = \mathbf{0.7}\end{aligned}$$

$$\mathbf{J} = (f_{11}) / (f_{01} + f_{10} + f_{11}) = 0 / (2 + 1 + 0) = \mathbf{0}$$

- Se \mathbf{d}_1 e \mathbf{d}_2 vetores numéricos, então

$$\cos(\mathbf{d}_1, \mathbf{d}_2) = \langle \mathbf{d}_1, \mathbf{d}_2 \rangle / \|\mathbf{d}_1\| \|\mathbf{d}_2\| ,$$

onde $\langle \mathbf{d}_1, \mathbf{d}_2 \rangle$ indica o produto escalar dos vetores, \mathbf{d}_1 e \mathbf{d}_2 , e $\|\mathbf{d}\|$ é a magnitude do vetor \mathbf{d} .

- Ex:

$$\mathbf{d}_1 = 3 \ 2 \ 0 \ 5 \ 0 \ 0 \ 0 \ 2 \ 0 \ 0$$

$$\mathbf{d}_2 = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 2$$

$$\langle \mathbf{d}_1, \mathbf{d}_2 \rangle = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$$

$$\|\mathbf{d}_1\| = (3*3 + 2*2 + 0*0 + 5*5 + 0*0 + 0*0 + 0*0 + 2*2 + 0*0 + 0*0)^{0.5} = (42)^{0.5} = 6.481$$

$$\|\mathbf{d}_2\| = (1*1 + 0*0 + 0*0 + 0*0 + 0*0 + 0*0 + 0*0 + 1*1 + 0*0 + 2*2)^{0.5} = (6)^{0.5} = 2.449$$

$$\cos(\mathbf{d}_1, \mathbf{d}_2) = 0.3150$$

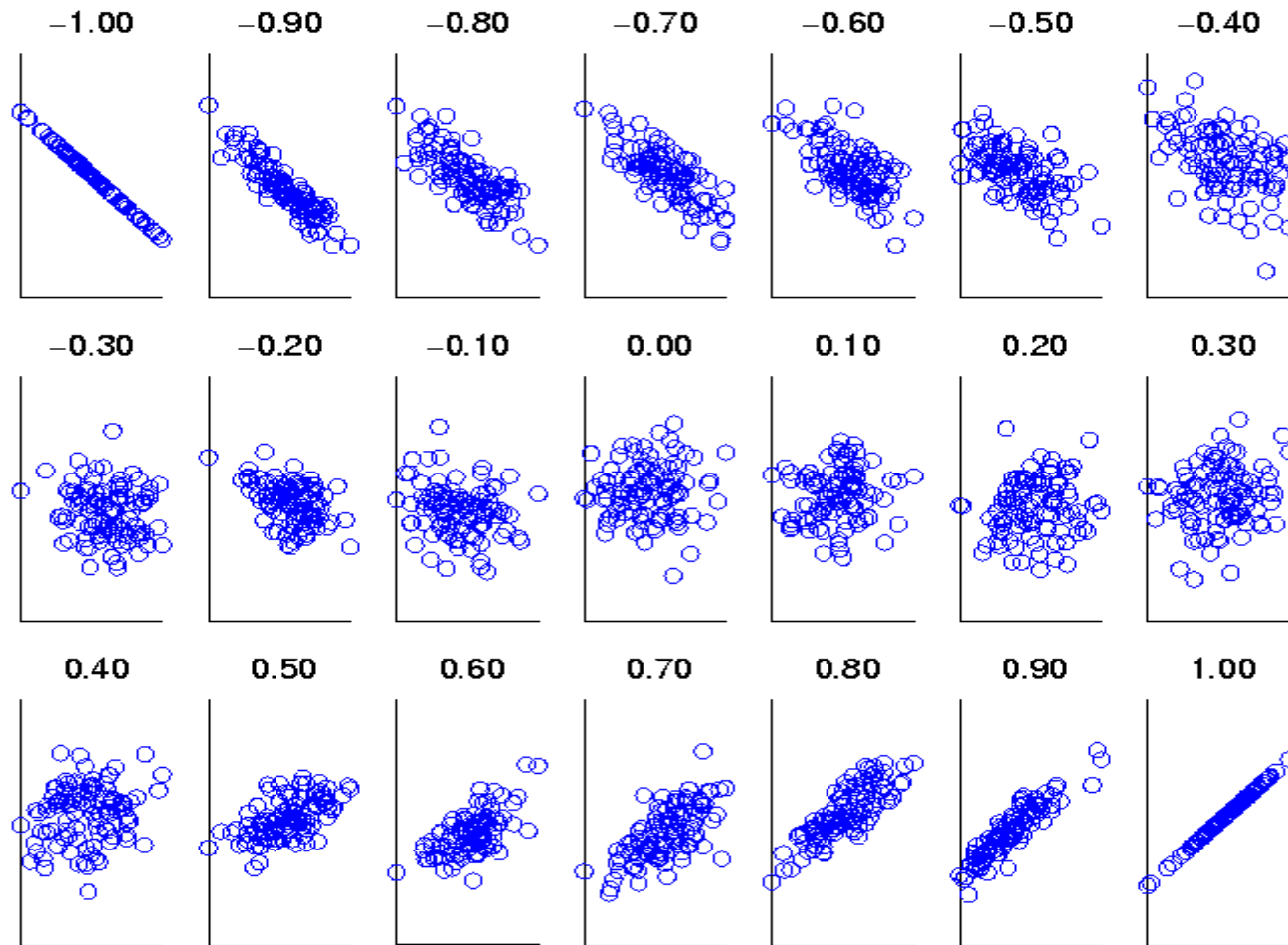
$$\text{corr}(\mathbf{x}, \mathbf{y}) = \frac{\text{covariance}(\mathbf{x}, \mathbf{y})}{\text{standard_deviation}(\mathbf{x}) * \text{standard_deviation}(\mathbf{y})}$$

p=1 Significa uma correlação perfeita positiva entre as duas variáveis.

p=-1 Significa uma correlação negativa perfeita entre as duas variáveis - Isto é, se uma aumenta, a outra sempre diminui.

p=0 Significa que as duas variáveis não dependem linearmente uma da outra. No entanto, pode existir uma dependência não linear.

Avaliando correlação visualmente



**Scatter plots
mostrando a
similaridade de –
1 a 1.**

- $\mathbf{x} = (-3, -2, -1, 0, 1, 2, 3)$
- $\mathbf{y} = (9, 4, 1, 0, 1, 4, 9)$

$$y_i = x_i^2$$

- $$\text{corr} = \frac{(-3)(5) + (-2)(0) + (-1)(-3) + (0)(-4) + (1)(-3) + (2)(0) + 3(5)}{6 * 2.16 * 3.74} = 0$$

Registro

- Matrizes
- Documentos
- Transações

Grafos

- WWW

Ordenados

- Séries temporais
- Dados sequenciais
- ...

Em uma coleção de registros todos os objetos possuem um conjunto fixo de atributos

<i>Tid</i>	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Se os objetos de dados possuem o mesmo conjunto fixo de **atributos numéricos**, então eles podem ser considerados **pontos em um espaço multidimensional**, onde cada dimensão representa um atributo distinto
- Podem ser representados por uma matriz $m \times n$, onde m são as linhas, uma para cada objeto, e n são colunas, uma para cada atributo

Projection of x Load	Projection of y load	Distance	Load	Thickness
10.23	5.27	15.22	2.7	1.2
12.65	6.25	16.22	2.2	1.1

- Caso especial de Matriz. Somente não zeros é importante.
- E.g., dados de documento: Cada documento vira um “vetor de termos”.
- Ex: número de vezes que um termo aparece no documento

	team	coach	play	ball	score	game	win	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0

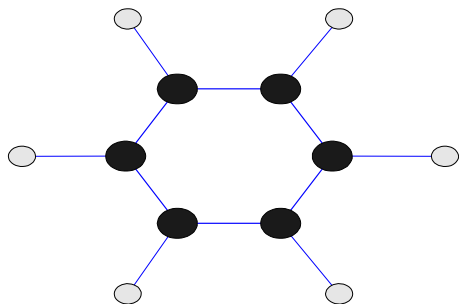
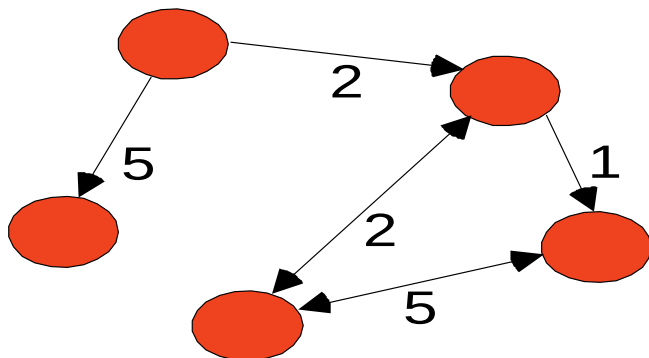
Representação conhecida como

Document-term matrix

- Um tipo especial de registro, onde:
 - Cada registro (transação) envolve um conjunto de itens
 - Ex: supermercado

<i>TID</i>	<i>Items</i>
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

- Ex: Interações humanas, uma molécula e páginas Web



Benzene Molecule: C₆H₆

Useful Links:

- [Bibliography](#)
- Other Useful Web sites
 - [ACM SIGKDD](#)
 - [KDnuggets](#)
 - [The Data Mine](#)

Knowledge Discovery and Data Mining Bibliography

(Gets updated frequently, so visit often!)

- [Books](#)
- [General Data Mining](#)

Book References in Data Mining and Knowledge Discovery

Usama Fayyad, Gregory Piatetsky-Shapiro, Padhraic Smyth, and Ramasamy uthurasamy, "Advances in Knowledge Discovery and Data Mining", AAAI Press/the MIT Press, 1996.

J. Ross Quinlan, "C4.5: Programs for Machine Learning", Morgan Kaufmann Publishers, 1993.
Michael Berry and Gordon Linoff, "Data Mining Techniques (For Marketing, Sales, and Customer Support)", John Wiley & Sons, 1997.

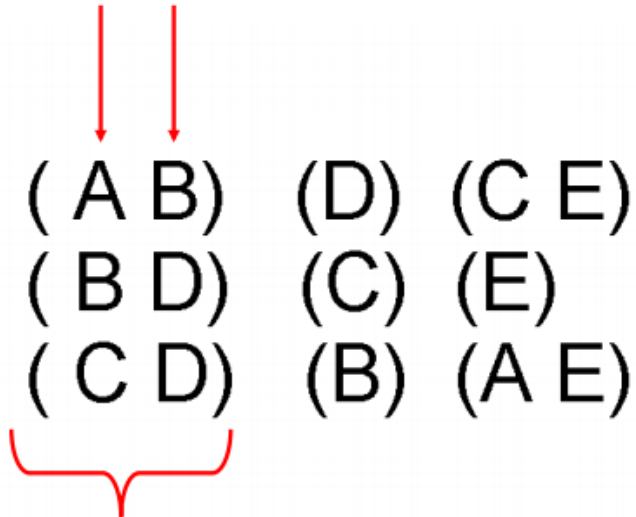
General Data Mining

Usama Fayyad, "Mining Databases: Towards Algorithms for Knowledge Discovery", Bulletin of the IEEE Computer Society Technical Committee on data Engineering, vol. 21, no. 1, March 1998.

Christopher Matheus, Philip Chan, and Gregory Piatetsky-Shapiro, "Systems for knowledge Discovery in databases", IEEE Transactions on Knowledge and Data Engineering, 5(6):903-913, December 1993.

- Sequências de transações, por exemplo.

Items/Events

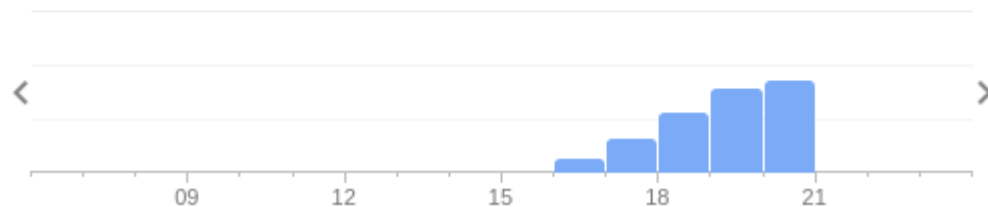


An element of
the sequence

- Série temporal
- Bar do Didi - Curitiba

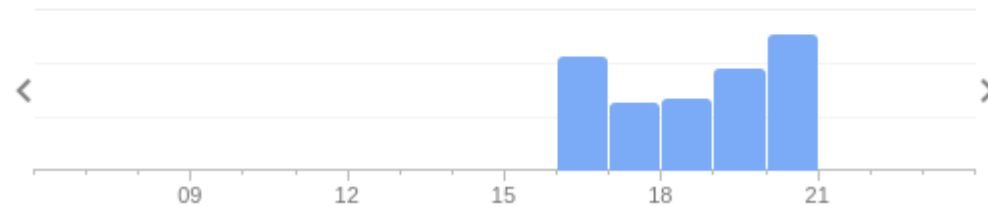
Horários de pico ?

Quintas-feiras



Horários de pico ?

Sextas-feiras



- Dados com baixa qualidade afetam negativamente muitos esforços de processamento de dados
- **Ex:** um modelo de classificação para detectar quem é um risco para conceder empréstimos com dados de baixa qualidade
 - Alguns bons candidatos possuem crédito negado
 - Mais crédito é concedido para indivíduos padrão

Exemplos de problema na qualidade dos dados:

Ruído e *outliers*

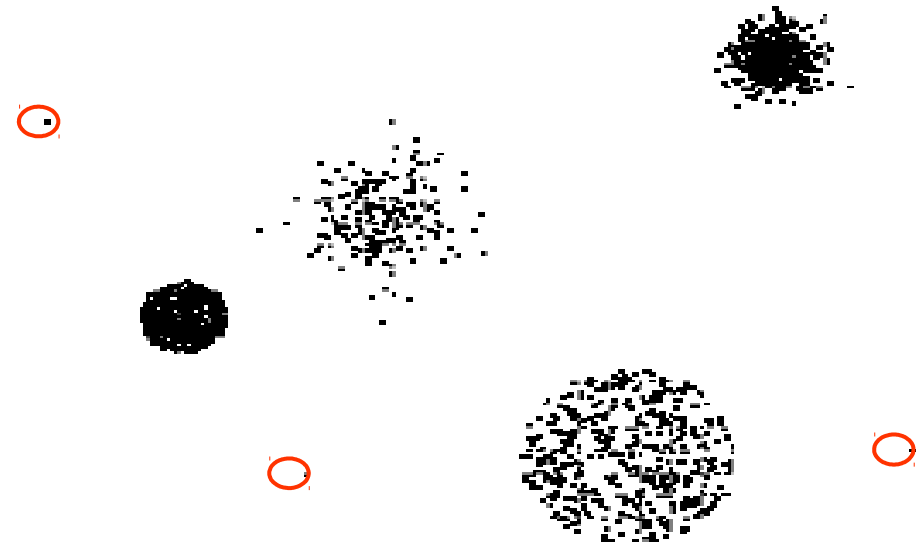
Dados faltantes

Dados duplicados

Dados errados

Outliers são objetos de dados com características que são consideravelmente diferentes da maioria dos outros objetos

- **Caso 1:** Outliers são ruído que interferem nas análises
- **Caso 2:** Outliers são o objetivo das análises
 - ◆ Fraude de cartão de crédito
 - ◆ Detecção de intrusos
- Causas?



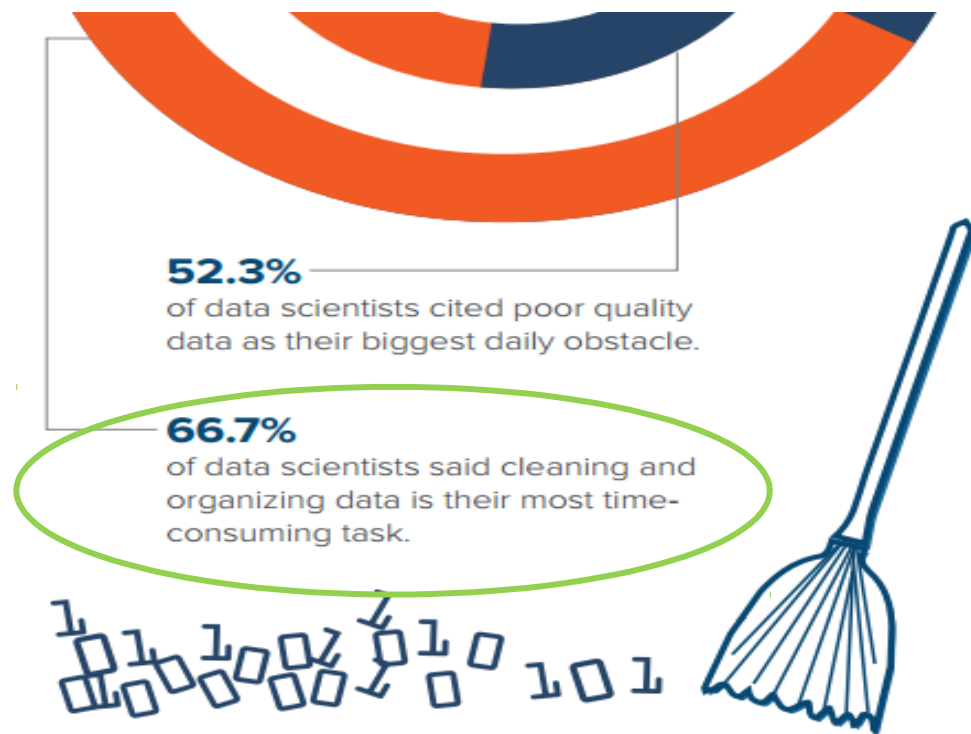
- Razões para dados faltantes:

- Informação não é coletada: e.g., pessoa se recusa a fornecer a idade
- Atributos podem não ser aplicados a todos os casos: e.g., renda anual não é aplicável para crianças

- Trabalhando com dados faltantes

- Eliminar objetos de dados ou atributos
- Estimar dados faltantes
 - ◆ Ex: séries temporais de temperatura
 - ◆ Ex: resultados do censo
- Ignorar dados faltantes durante as análises

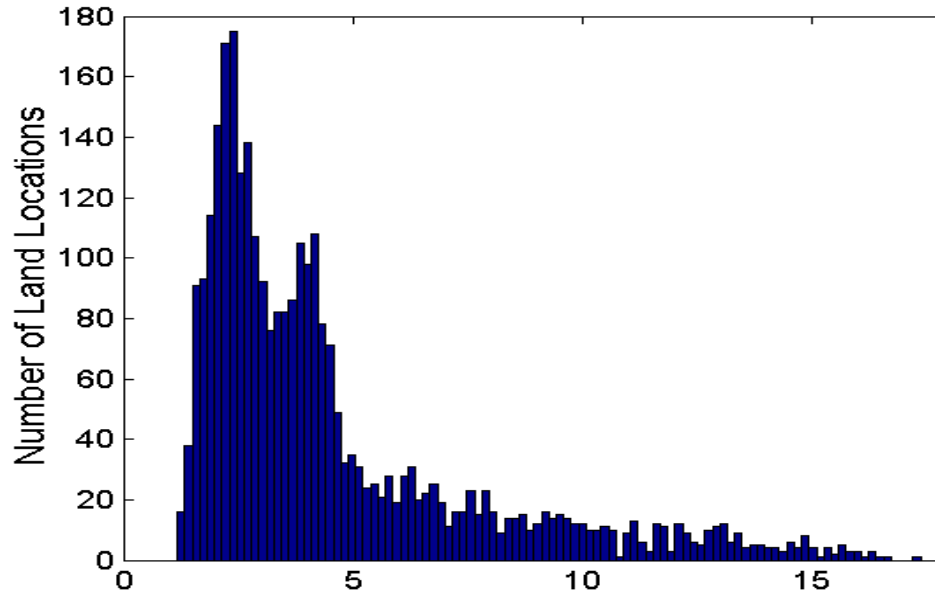
- Datasets podem incluir objetos duplicados (ou parcialmente duplicados)
 - Principalmente ao fundir dados de fontes diferentes
- Ex:
 - A mesma pessoa com o mesmo ID em diferentes redes.



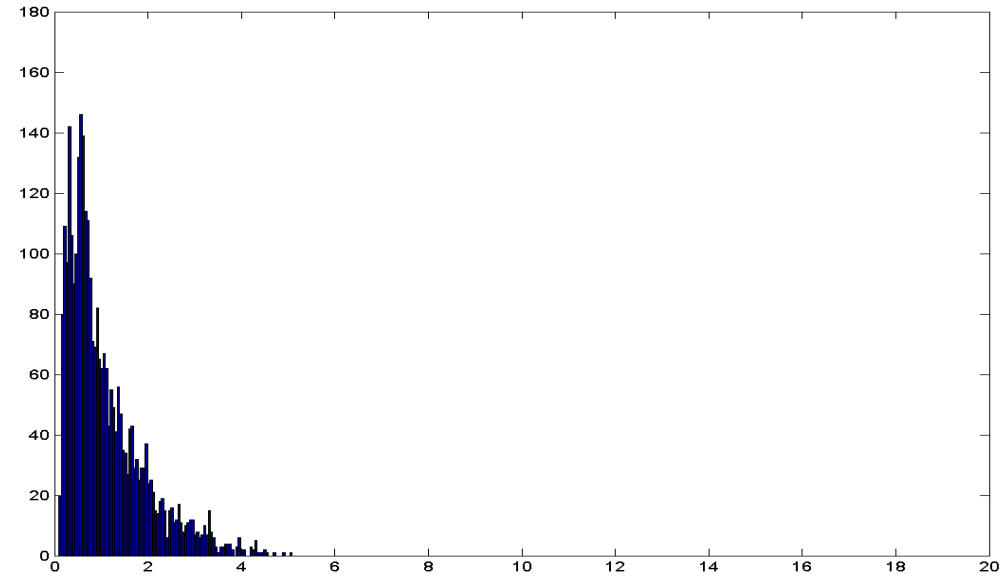
Estas tarefas tendem a ser demoradas e tediosas

- Combina 2 ou mais atributos em um único atributo
- **Objetivos**
 - Redução de dados
 - Reduzir o número de atributos
 - Mudar a escala
 - ◆ Cidades agregadas em regiões ou estados, etc.
 - ◆ Dias em semanas, meses ou anos
 - Dados mais “estáveis”
 - ◆ Dados agregados tendem a ter menos variabilidade

Variação da precipitação na Austrália (grid cells)



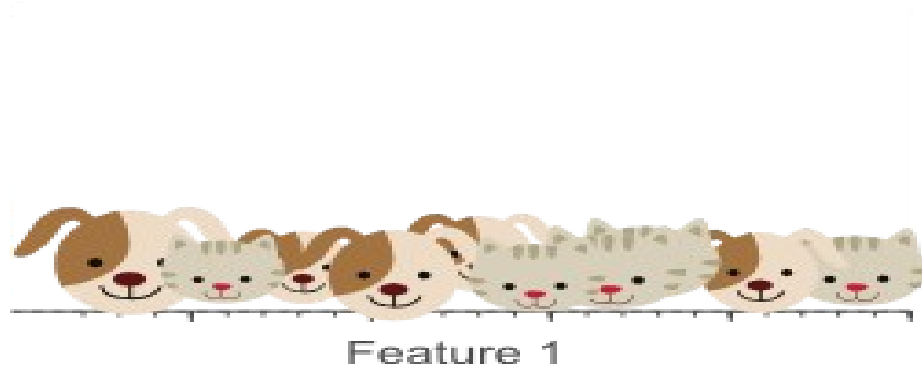
Desvio padrão médio da precipitação
por mês



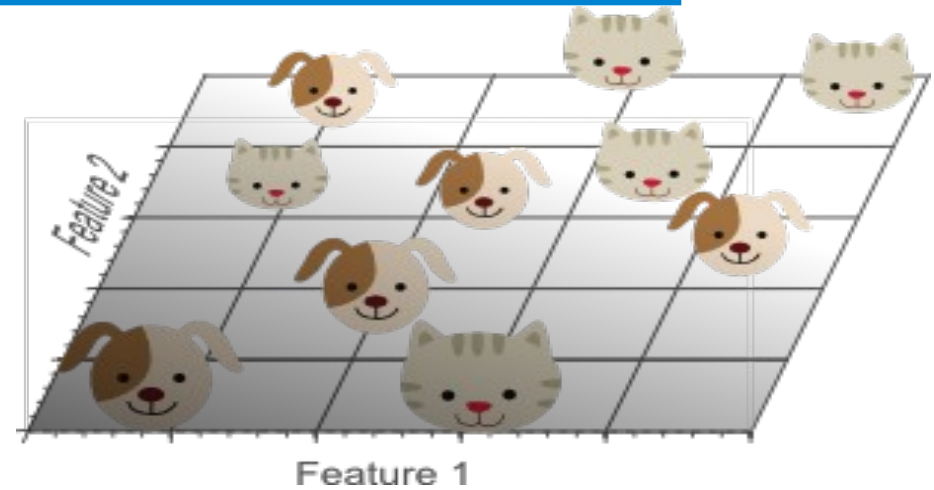
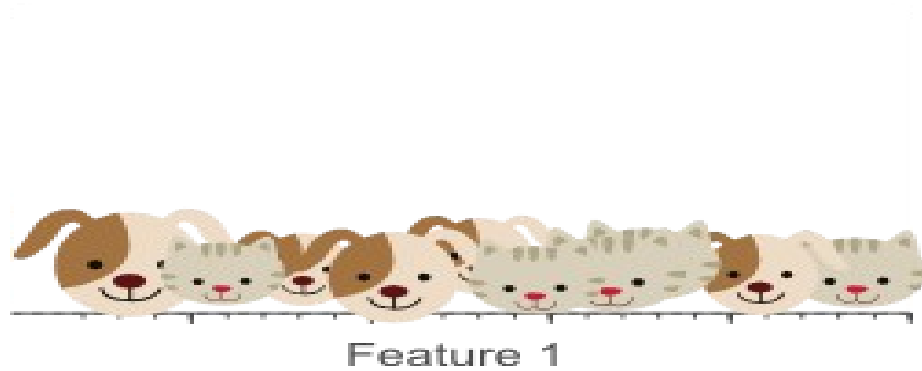
Desvio padrão médio da
precipitação por ano

A precipitação média por ano varia menos

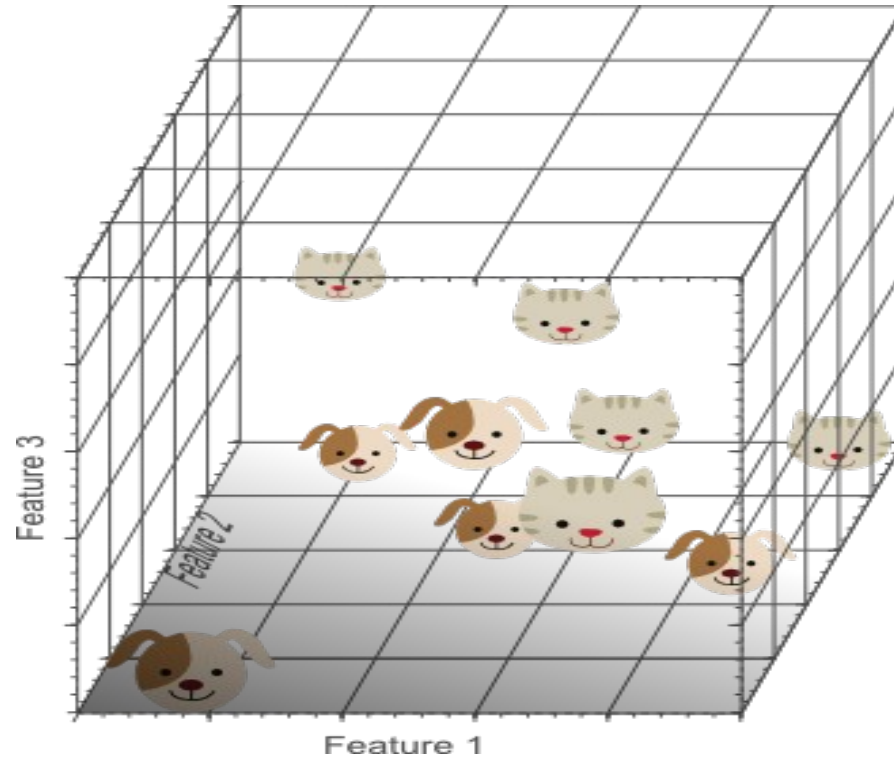
Importância das dimensões



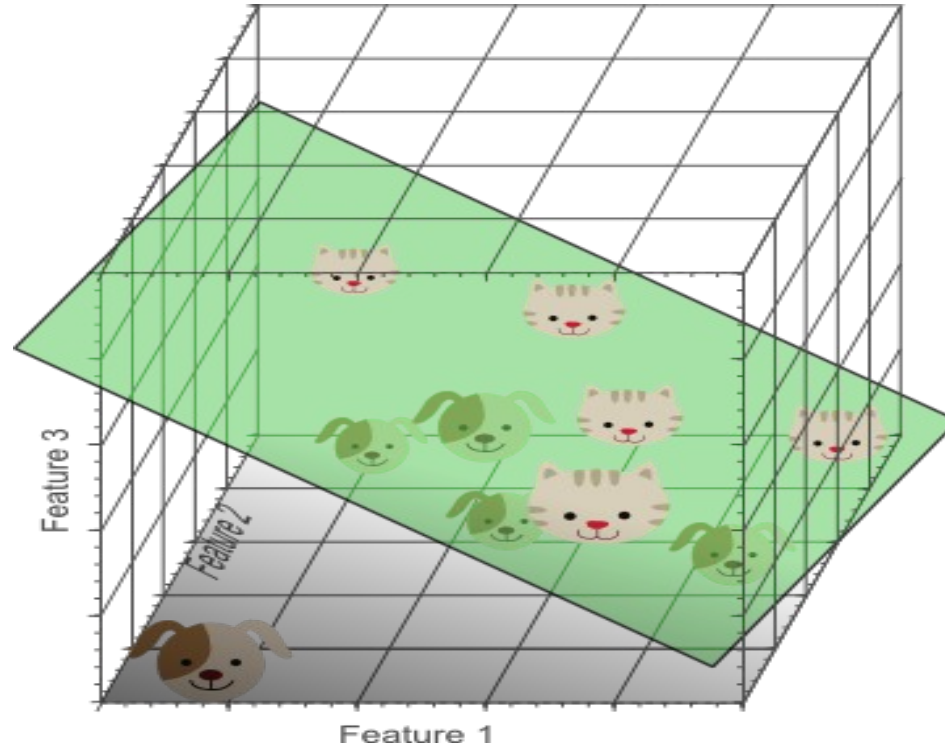
Importância das dimensões



Importância das dimensões



Importância das dimensões



Objetivo: capturar informações importantes no dataset mais eficientemente do que com os atributos originais

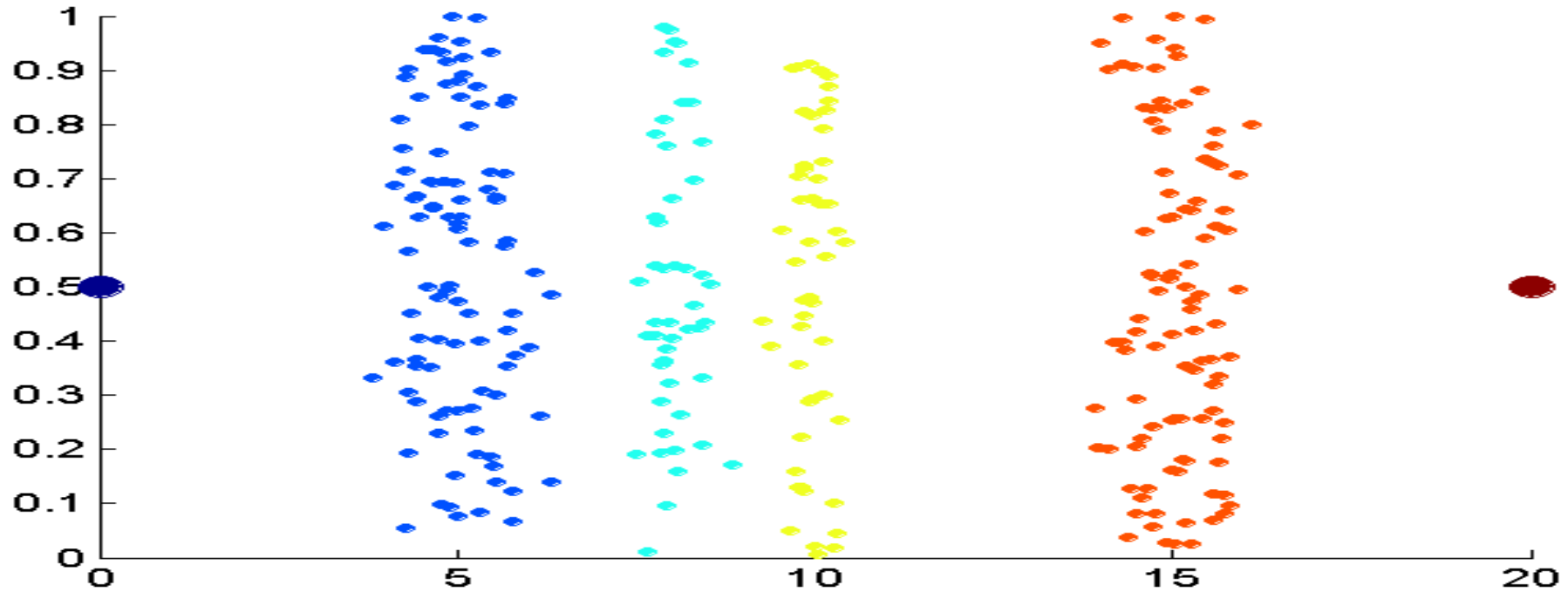
- Três metodologias gerais:

- Extração de features
 - Ex: bordas de imagens
- Construção de features
 - Ex: dividir a massa pelo volume para obter a densidade
- Mapear dados em um novo espaço
 - ◆ Ex: Análises de Fourier e wavelet

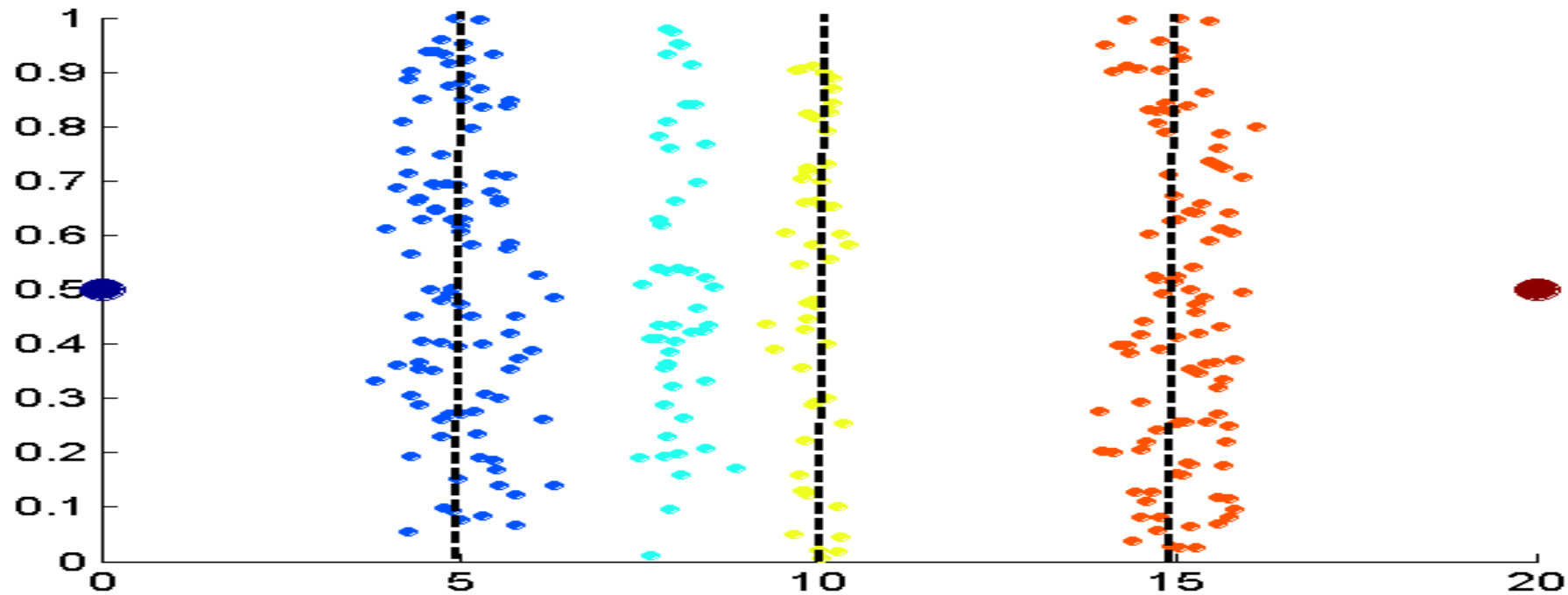
Processo de converter um atributo contínuo em um atributo categórico (ordinal ou nominal)

Potencialmente, um número infinito de valores numéricos são mapeados em um número menor de categorias

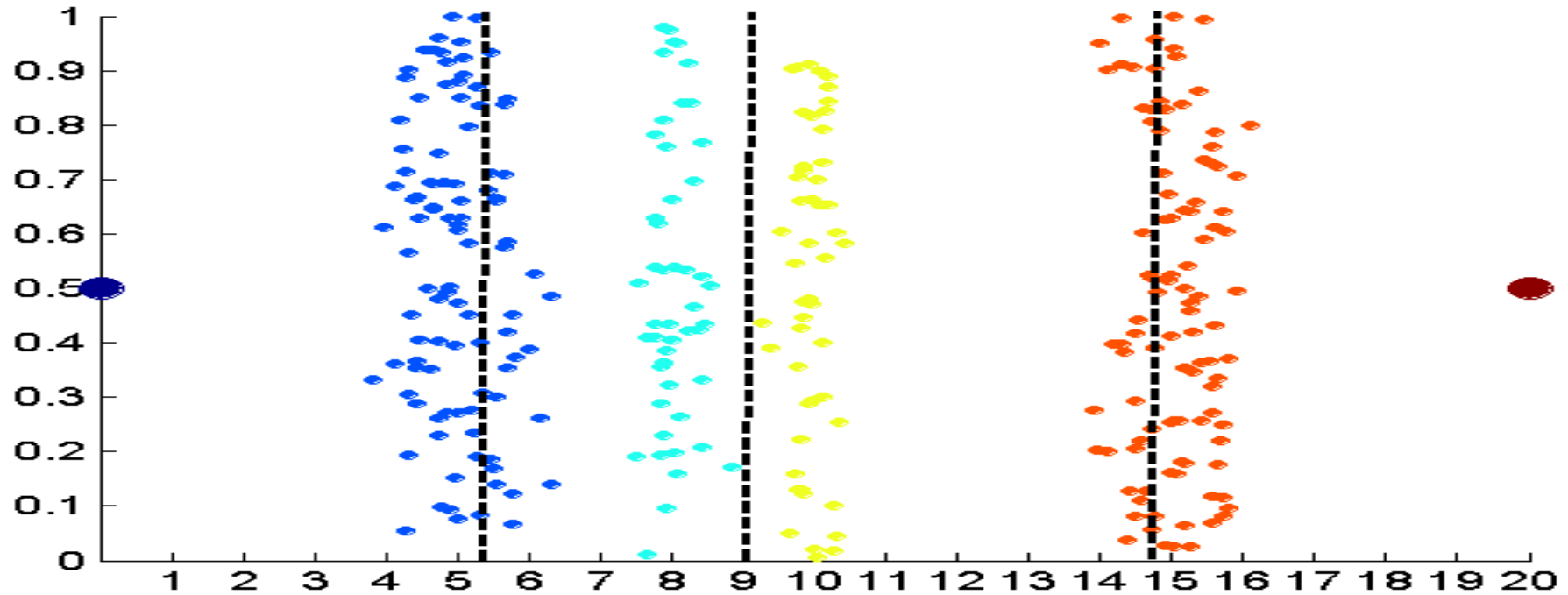
- Vários algoritmos de classificação funcionam melhor se variáveis possuem menos valores



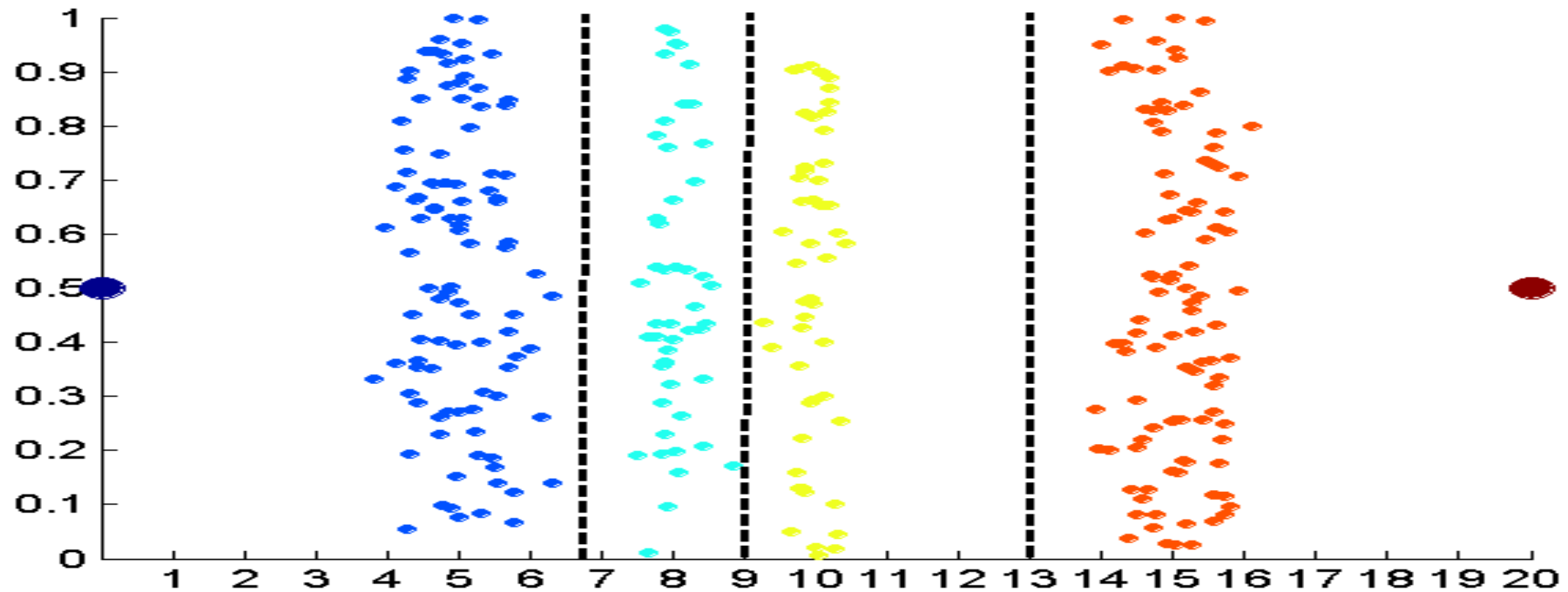
Dados consistem de quatro grupos de pontos e dois outliers



Intervalos de largura iguais



Igual frequência



Agrupamento

- Mapeia um atributo contínuo ou categórico em uma ou mais variáveis binárias
- Usado bastante para análises de associação
- Geralmente converte um atributo contínuo em categórico e depois em um conjunto binário

Table 2.6. Conversion of a categorical attribute to five asymmetric binary attributes.

Categorical Value	Integer Value	x_1	x_2	x_3	x_4	x_5
<i>awful</i>	0	1	0	0	0	0
<i>poor</i>	1	0	1	0	0	0
<i>OK</i>	2	0	0	1	0	0
<i>good</i>	3	0	0	0	1	0
<i>great</i>	4	0	0	0	0	1

Parte deste material é derivado do livro:
Introduction to Data Mining - Tan, Steinbach, Kumar