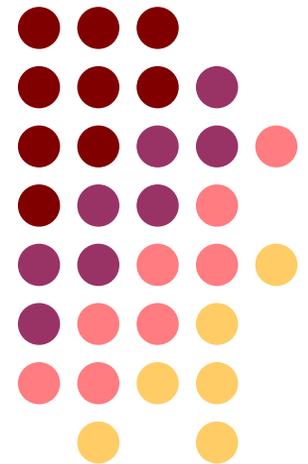


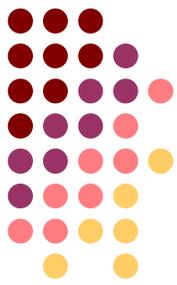
# Dados Simbólicos

Universidade Federal de  
Pernambuco

[CIn.ufpe.br](http://CIn.ufpe.br)



# Dados Simbólicos



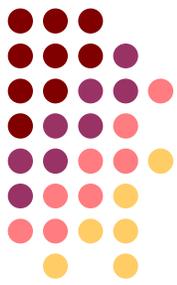
Dados simbólicos são informações complexas que são expressas por intervalos, conjuntos, frequências distribuições de probabilidade.

Exemplos de situações que necessitam desses tipos de dados

1 - Intervalo para indivíduos (ou objetos de primeira ordem)

$Y$  = tempo gasto com estudos por dia

$Y(k) = [0,6]$  (em horas)



2 - Conjunto para classes de indivíduos (objetos de segunda ordem, objetos agregados)

$Y$  - instituições bancárias existentes em uma cidade

$Y(k) = \{\text{Bandeirantes, Itau, Bradesco, BB}\}$

onde  $\omega$  é uma cidade (classe de indivíduos)

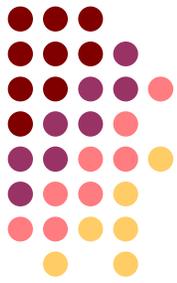
3 - Distribuição de frequências -  $Y$  = percentual de votos por partido político em um Estado

$Y(k) = ((\text{ABR};0,50),(\text{PDF};0,30),(\text{DEO};0,20))$

4 - Distribuição de probabilidade -  $Y$  = níveis de cinza de uma região de uma imagem

$Y(k) = \Gamma(20,30)$

# Tipos de Variáveis Simbólicas



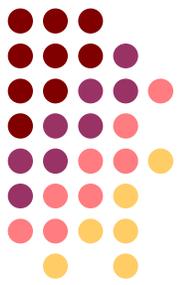
## Notação

Y - uma variável

E - conjunto de entidades

k - uma entidade (ex. uma classe de estudantes,  
uma empresa)

O - conjunto dos possíveis valores de Y



# 1 - Variáveis Multi-valorada

Uma variável  $Y$  é multi-valores é uma função

$$Y_i : E \rightarrow B = P(O_i)$$
$$y_i(k) \rightarrow U$$

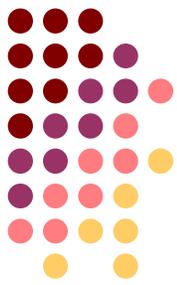
onde  $U \subseteq P(O_i)$  conjunto de todos os subconjuntos de  $O_i$ . No caso clássico  $|y(k)| = 1$ .

## Multi-valor categórico

a)  $U$  pode ser um subconjunto categórico.

Exemplo

$\text{cor}(k) = \{\text{azul}, \text{verde}, \text{amarelo}\}$



## Multi-valor quantitativo

b)  $U$  pode ser um subconjunto finito de números reais.

Exemplo

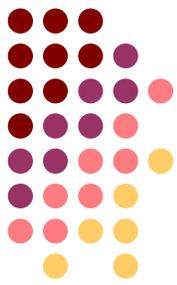
$Y$ =quantidade de partos realizados em 5 hospitais de uma cidade ( $k$ ) em uma semana.

$y(k) = \{100, 200, 150, 220, 180\}$

## Multi-valor ordinal

Se o domínio de  $Y$  é definido com uma ordenação conforme definido com os dados clássicos.

Na prática a variável tem uma ordenação.



## 2 - Variáveis Intervalo

Uma variável  $Y$  é denominada intervalo se para todo  $k \in E$  o subconjunto  $y(k) = [a, b]$  onde  $a \leq b$ .

Exemplo

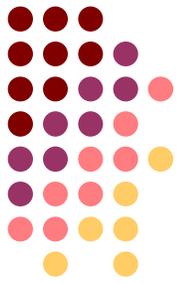
$E$  = conjunto de empresas

$Y$  = valor gasto com impostos

$k$  = uma empresa

$y(k) = [2000, 3000]$

# Dois níveis de paradigmas obtidos por agregação



Exemplo

$E = \{C_1, C_2, \dots, C_m\}$  conjunto de classes de estudantes

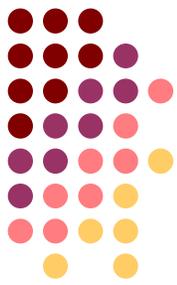
$Y$  = descreve a altura na classe  $C_i$

a)  $Y$  - variável multi-valor

$$y(C_i) = \{1.50, 1.56, 1.58, 1.70\}$$

b)  $Y$  - variável intervalo

$$y(C_i) = [1.50, 1.70]$$



### 3 - Variáveis Modais

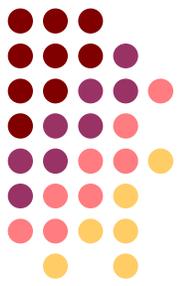
Uma variável modal  $Y$  com domínio é uma função

$$Y_i : E \rightarrow B=M(O_i)$$
$$y(k) \rightarrow (U(k), \pi_k)$$

onde

- $\pi_k$  é uma medida ou uma distribuição (probabilidade, frequência ou peso)
- $U(k) \subset O_i$
- $M(O_i)$  é uma família de medidas  $\pi$  em  $O_i$

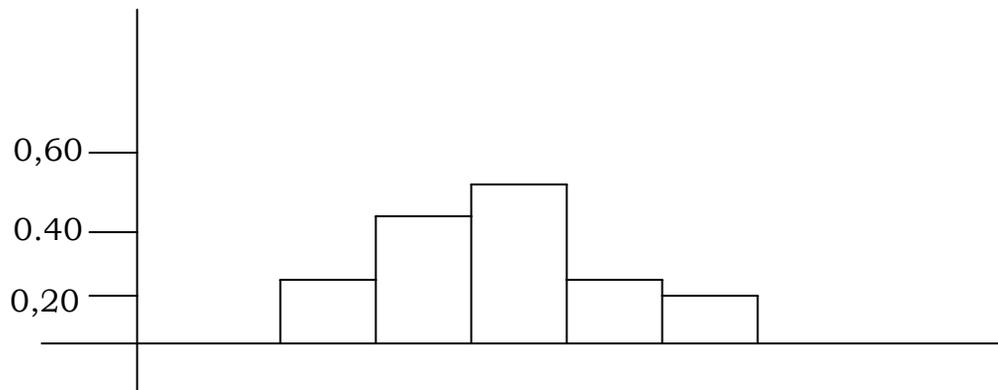
# Exemplo



$C = \{C_1, \dots, C_{10}\}$  um conjunto de 10 empresas com 30 funcionários

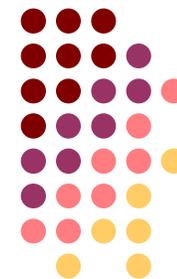
Y - descreve altura na empresa  $C_i$

1 - Um histograma com os intervalos:



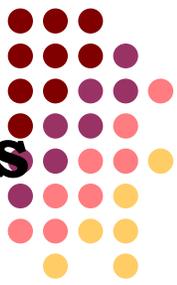
2 - Uma distribuição normal  $N(168, 48.4)$  com média 168 e variância 48.4

# Tabela de dados simbólicos



classe	Marca	Dist. de freq. do tempo de uso	Preço
$e_1$	{CITR, PEUG}	( A 0.4 B0.3 C0.4)	[12000,30000]
$e_2$	{BMW,FIAT}	( A 0.2 B0.5 C0.3)	[15000,35000]
$e_3$	{BMW,FIAT, ALFA}	( A 0.6 B0.2 C0.2)	[8000,45000]
:	:	:	:

# Dependências entre Variáveis Simbólicas



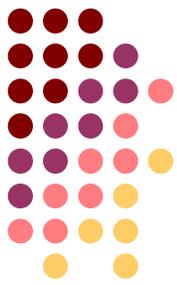
Entre variáveis simbólicas são adotados dois tipos de dependências que são expressas através de regras

## 1 - Dependência Hierárquica

Uma variável  $Y_i$  pode tornar-se inaplicável se outra variável  $Y_j$  assume valores em um subconjunto  $S_j \subset O_j$ .

A dependência é expressa pelas regras

$$r: [[y_j \in S_j]] \xrightarrow{DH} [y_i = NA] \quad (\text{equivalência lógica})$$



ou equivalentemente pelas regras

$$r_1: [[y_j \in S_j]] \rightarrow [y_i = \text{NA}]$$

$$r_2: [y_i = \text{NA}] \rightarrow [[y_j \in S_j]]$$

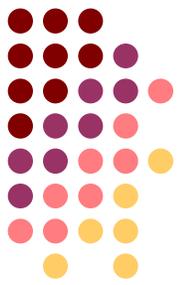
## Exemplo

$Y_1 = \text{Sexo} = \{M, F\}$  e  $Y_2 = \text{Parto} = \{\text{sim}, \text{n\~{a}o}\}$

A depend\~{e}ncia entre  $Y_1$  e  $Y_2$  \u00e9 expressa pelas regras

$r_1$ : se  $[y_1 = M]$  **ent\~{a}o**  $[y_2 = \text{NA}]$

$r_2$ : se  $[y_2 = \text{NA}]$  **ent\~{a}o**  $[y_1 = M]$



## 2 - Dependência lógica

Um subconjunto  $S_i \subset O_i$  da variável  $Y_i$  pode estar em correspondência com o subconjunto  $S_j \subset O_j$  da variável  $Y_j$ .

A dependência é expressa pela regra

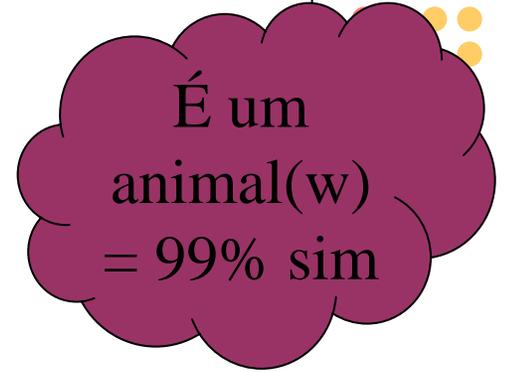
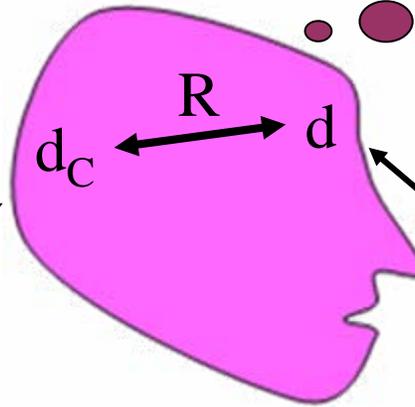
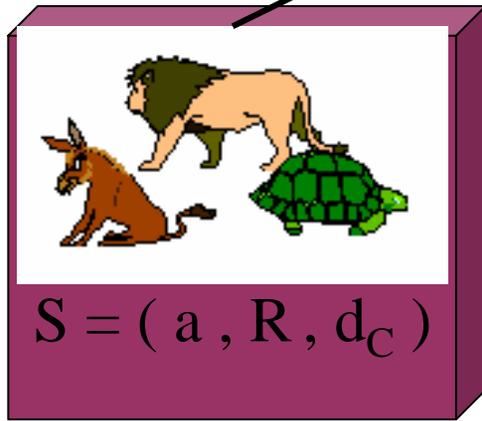
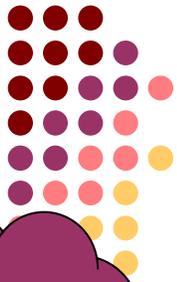
$$r: [[y_j \in S_j]] \stackrel{DL}{\rightarrow} [y_j \in S_j]$$

Exemplo

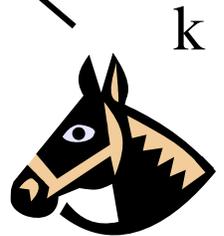
$Y_1 = \text{idade}[0,90]$  e  $Y_2 = \text{altura (em cm)}[0,4;2,00]$

A dependência entre  $Y_1$  e  $Y_2$  é expressa pela regra  
r: se  $[y_1 = [0,10]]$  **então**  $[y_2 = [0,40;1,60]]$

# Objetos Simbólicos

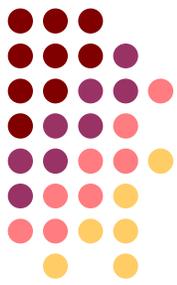


y



$$a(w) = [y(k) R d_C]$$

# Objetos Simbólicos Booleanos

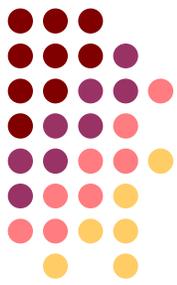


## Definição

Um objeto simbólico booleano é uma tripla  $s = (a, R, d)$  onde  $R$  é uma relação entre descrições,  $d \in D$  uma descrição e  $a$  uma função binária  $a_s : E \rightarrow \{0, 1\}$  com

$$a_s(k) = \bigwedge_{j=1}^p [y_j(k) R d] = 1 \quad \text{se e somente se} \quad [y_j(k) R d] = 1$$

$$\forall j \in \{1, \dots, p\}$$



A **extensão** de  $s$  é definida como

$$\text{Ext}(s) = \{k \in E / a_s(k) = 1\}$$

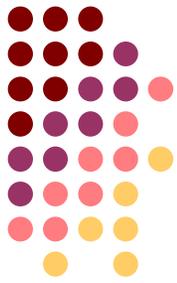
## Exemplo

Considere uma tabela de dados clássicos e  $s=(a,R,d)$  um Objeto simbólico onde  $y(k) = [\text{altura}(k), \text{peso}(k)]$ ,  $d = [[140, 160], [50, 60]]$  e  $R = \in$  então

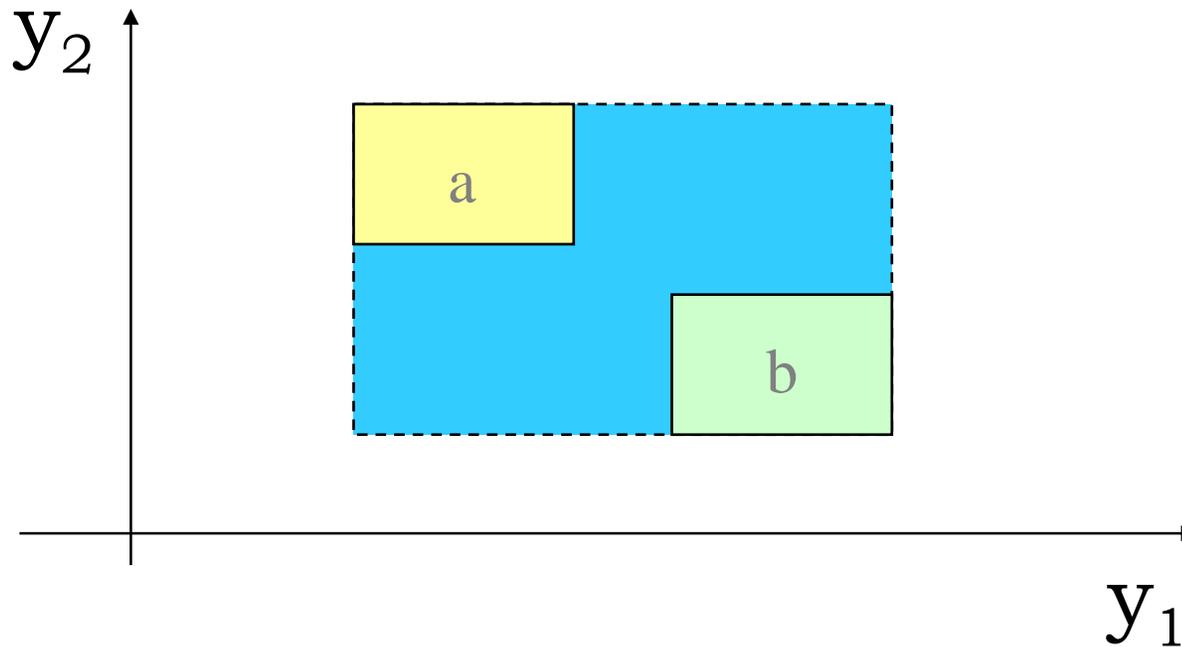
$$a(k) = [\text{altura}(k) \in [165, 175]] \wedge [\text{peso}(k) \in [50, 60]]$$

$$\text{Ext}(s) = \{\omega_1, \omega_3\}$$

# Operadores Simbólicos

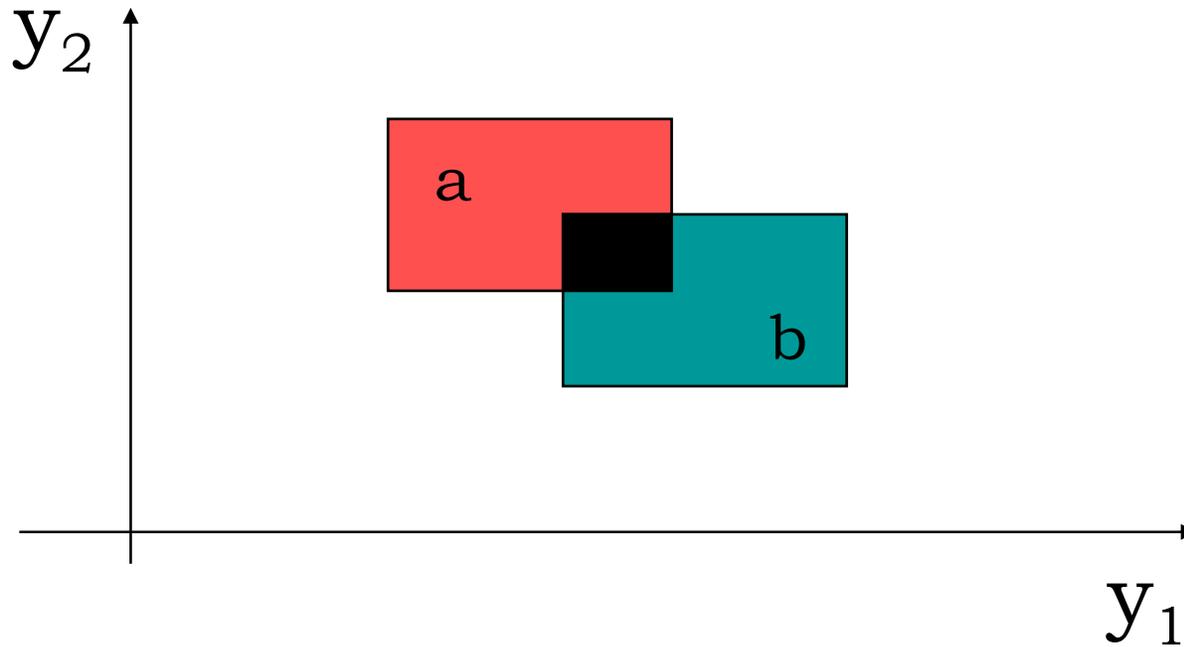
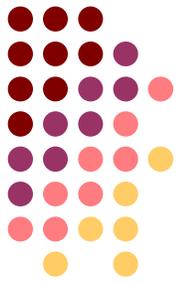


## Junção



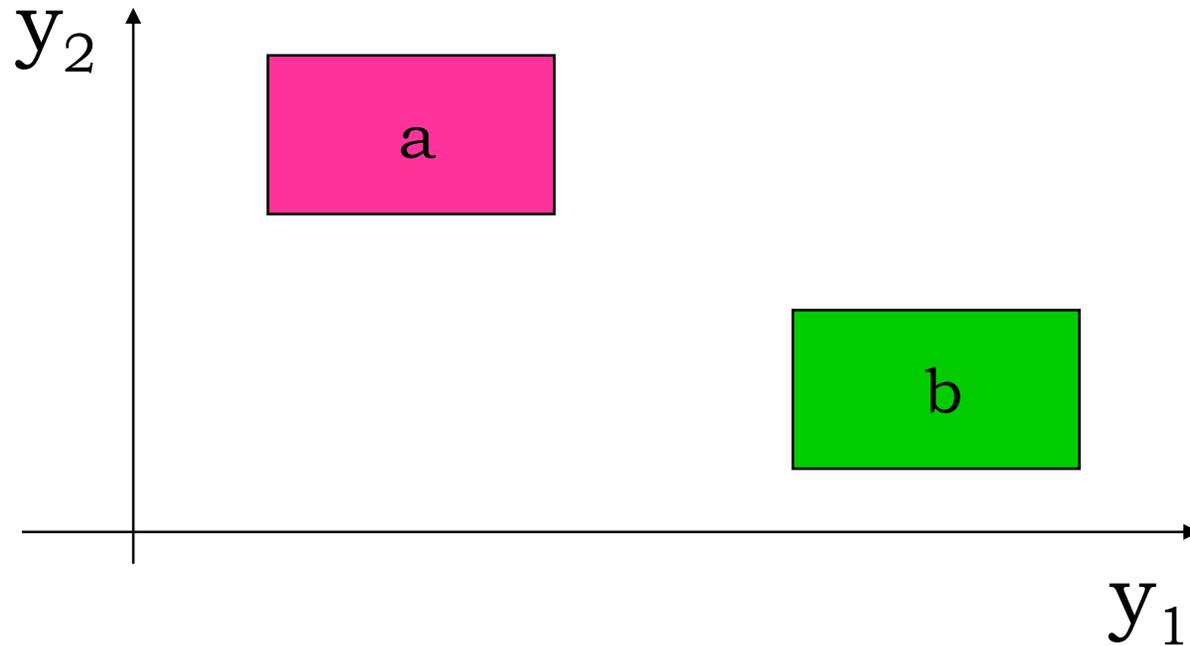
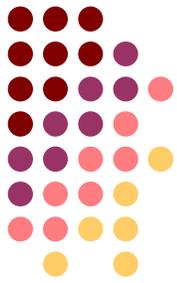
$$a \oplus b = [y_1 \in C_1] \wedge [y_2 \in C_2]$$

# Conjunção



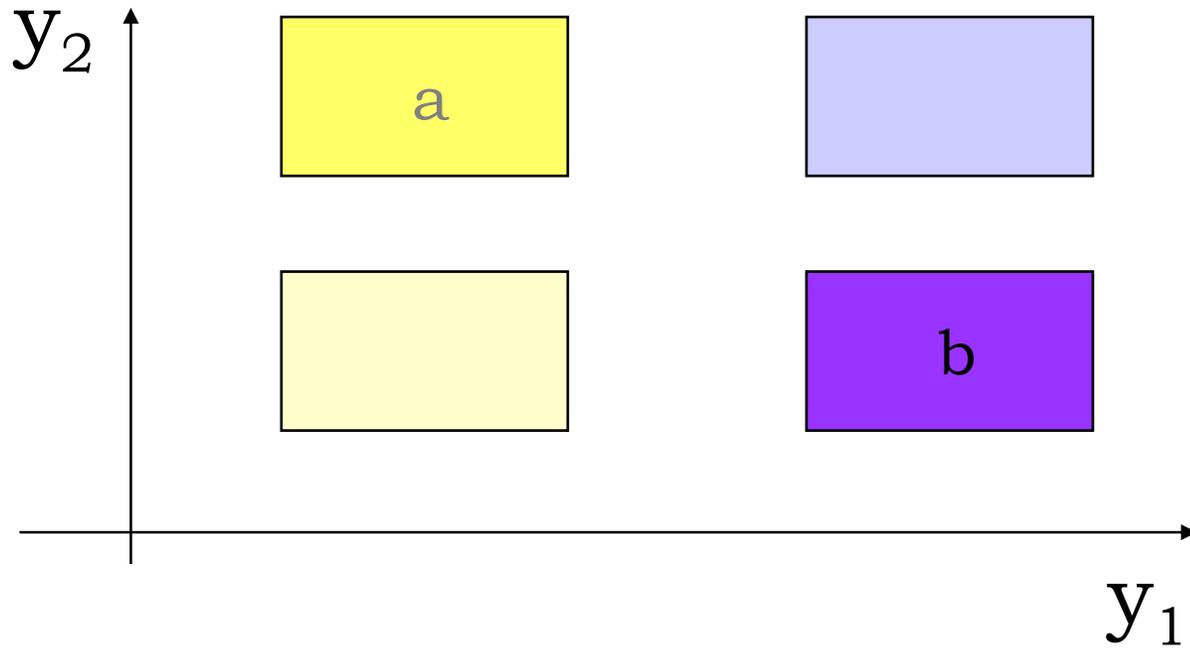
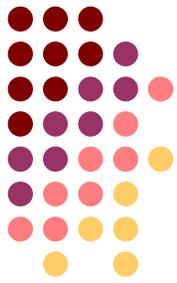
$$a \wedge b = [y_1 \in (A_1 \cap B_1)] \wedge [y_2 \in (A_2 \cap B_2)]$$

# Disjunção

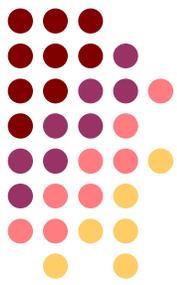


$$a \vee b = \{ [y_1 \in A_1] \wedge [y_2 \in A_2] \} \vee \{ [y_1 \in B_1] \wedge [y_2 \in B_2] \}$$

# União



$$a \cup b = [y_1 \in A_1 \cup B_1] \wedge [y_2 \in A_2 \cup B_2]$$



## Potencial de descrição

O potencial de descrição é uma medida positiva definida no objeto que representa o volume (De Carvalho, 1995) .

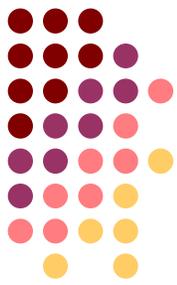
$$\pi (s ) = \prod_{j=1}^p \mu (d_j )$$

onde  $\mu(d_i)$  = cardinal de  $d_i$ , se  $d_i$  é um conjunto  
amplitude de  $d_i$ , se  $d_i$  é um intervalo

### Exemplo

Seja o objeto  $s$  apresentado anteriormente, então o potencial de  $s$  é  $\pi(s) = (175-165) \times (60-50) = 100$

# Objetos Simbólicos Modais

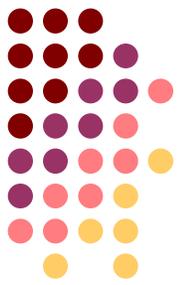


## Definição

Um objeto simbólico modal é uma tripla  $s = (a, \Phi, d)$  onde  $\Phi$  é uma relação de fuzzy e  $a$  uma função  $a_s : E \rightarrow [0, 1]$ .

A **extensão** de  $s$  é definida como

$$\text{Ext}(s) = \{k \in E / a_s(k) > \alpha\}$$



## Exemplo

Seja  $E$  um conjunto de espécies de animais,  $k$  uma classe de animais  $Y_1 = \text{cor}$  e  $Y_2 = \text{mês de nascimento}$ . Então  $O_1 = \{\text{branco, cinza, marrom, preto, azul, ...}\}$  e  $O_2 = \{\text{jan, ..., dez}\}$

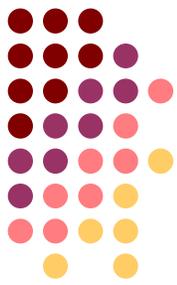
Um objeto  $s = (a, \Phi, d)$  onde  $d = [\{\text{preto, marrom}\}, \{\text{maio, junho}\}]$  e  $\Phi$  é uma relação de fuzzy definida como

$$\Phi(y(k), d) = \frac{|y(k) \cap d|}{|y(k) \cup d|}$$

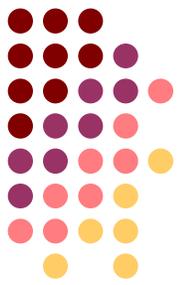
Então para  $y(k) = [\{\text{branco, marrom}\}, \{\text{jun, jul}\}]$   $\Phi(y(k), d) = 1/7$ .

Adotando  $\alpha = 0.5$  logo  $a(k) = 0.13 < 0.5$  logo  $k \notin \text{Ext}(s)$

# Construindo descrições: O processo de generalização



- Sendo  $D_c$  = conjunto das descrições dos indivíduos de  $C$
- Aplicação dos operadores
  - O intervalo  $G_y(C) = [\inf(D_c), \sup(D_c)]$  constitui uma boa generalização de  $C$
- Seja
  - $C = \{w_1, w_2, w_3\}$  e
  - $D_c = \{y(w_1), y(w_2), y(w_3)\} = y(C)$
  - A generalização de  $C$  para a variável  $y$  é  $G_y(C)$
- Exemplos
  - Seja  $y$  uma variável numérica tal que:
    - $y(w_1) = 2.5, y(w_2) = 3.6, y(w_3) = 7.1$
  - Seja  $D$  o conjunto de valores incluídos em  $[1, 100]$ 
    - $\rightarrow G_y(C) = [2.5, 7.1]$  é a generalização de  $D_c$  para  $y$



# O processo de generalização

## ■ Exemplos

□ Seja  $y$  uma variável modal:

■  $y = \{\text{pequeno, grande}\}$  tal que:

■  $y(w1) = (1(1/3), 2(2/3))$

□ em que  $2(2/3)$  significa que a frequência da categoria 2 é  $2/3$

■  $y(w2) = (1(1/2), 2(1/2))$

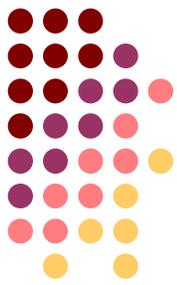
■  $y(w3) = (1(1/4), 2(3/4))$

□ Então  $Gy(C) = [ [ [1(1/4)], [1(1/2)] ], [ [2(1/2)], [2(3/4)] ] ]$   
é a generalização de  $D_c$  para  $y$

■ Isto significa que da categoria 1 o intervalo é  $[1/4, 1/2]$

■ Isto significa que da categoria 2 o intervalo é  $[1/2, 3/4]$

# O processo de generalização

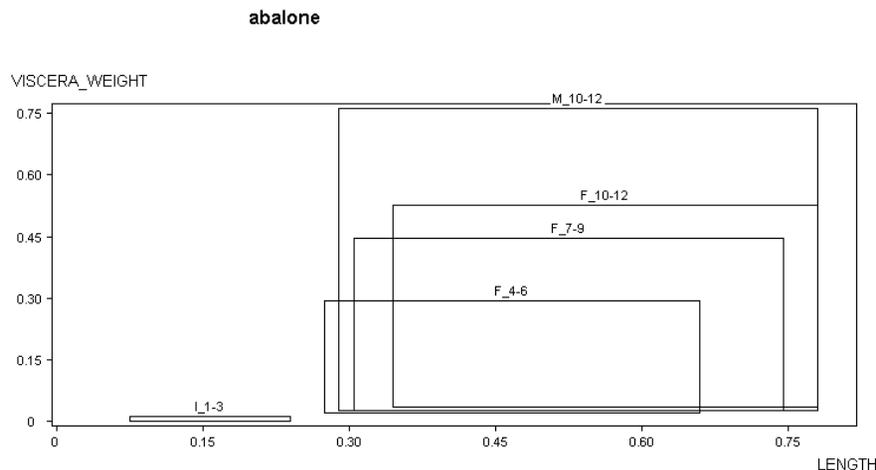


## ■ Exemplos

- Seja  $y$  uma variável intervalar tal que:
  - $y(w1) = [1.5, 3.2]$
  - $y(w2) = [3.6, 4]$
  - $y(w3) = [7.1, 8.4]$
  
- Então  $Gy(C) = [1.5, 8.4]$  é a generalização de  $D_c$  para a variável  $y$
- Seja  $C = \{w1, w2, w3\}$  e  $y$  é uma variável categórica não ordenada
- $y(w1) = 2, y(w2) = 2, y(w3) = 1$
- $D$  é o conjunto de ocorrência dos valores 1 e 2
- $G'(y(C)) = [1(1/3), 2(2/3)]$ 
  - $\rightarrow D_c = \{y(w1), y(w2), y(w3)\} = y(C)$  para a variável  $y$ .

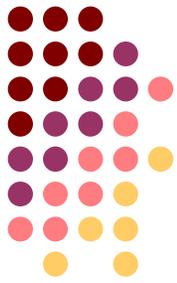
# O processo de generalização

- Uso dos operadores
  - Simples
  - Pode gerar outliers
- Estratégia
  - Reduzir os limites para reduzir ocorrência de *outliers*
    - *Teste komogorov Sminorv (distribuição uniforme)*
    - *Algoritmo de agrupamento*
  - Utilizar freqüências para o caso da variáveis categóricas

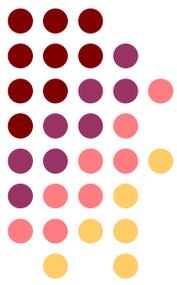


# O processo de generalização

- O processo se preocupa com pelo menos 2 questões
  - Overgeneralization
    - Ex.  $\text{age}(\text{alunos\_renata}) = [18,36]$
  - Perda de informação estatística, como média, variância, correlações entre variáveis
    - Solução  $\rightarrow$  criação de novas variáveis simbólicas contendo os respectivos valores



# Um critério de avaliação da descrição generalizada



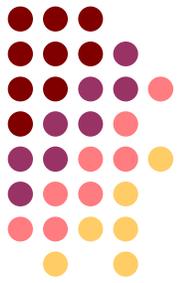
Pode-se avaliar a qualidade de uma generalização  $g_i$  associada a uma descrição  $d_i$  da classe  $C_i$  através da medida

$$\text{dens}(s_i) = \frac{|\text{Ext}(s_i / C_i)|}{\pi(d_i)}$$

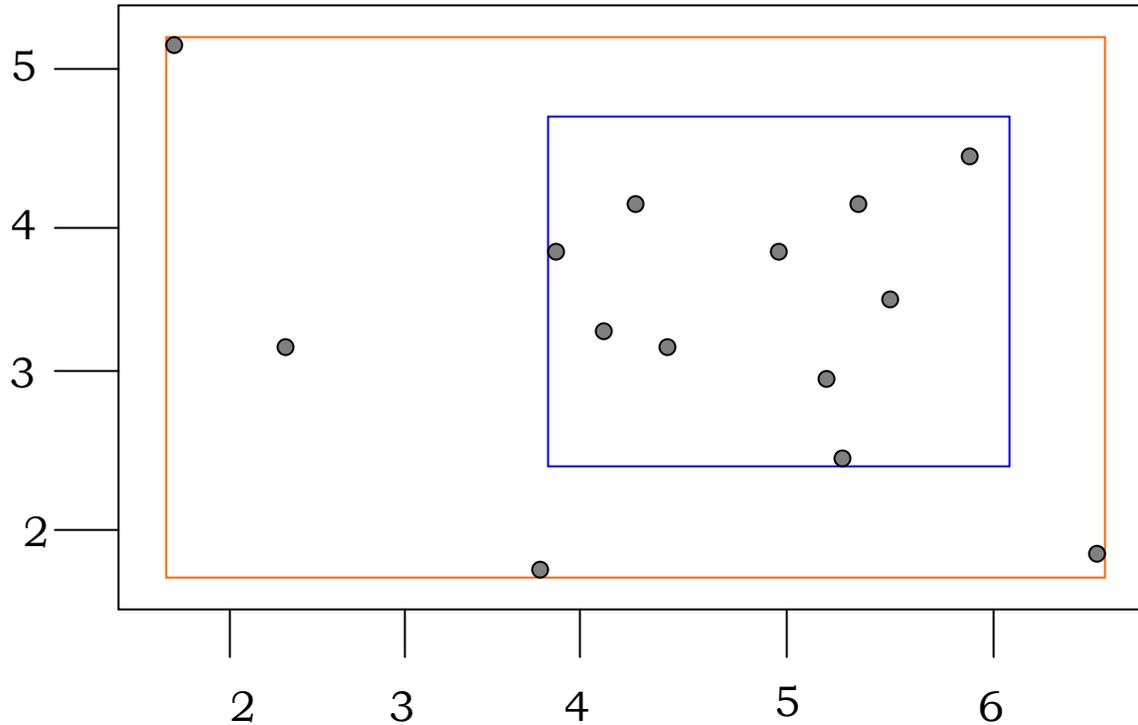
denominada de **densidade do objeto simbólico  $s_i$**

$\text{Ext}(s_i / C_i)$  define o conjunto de indivíduos que satisfazem a descrição do objeto  $s_i$

$\pi(d_i)$  é o volume que produz o índice de generalização de  $d_i$



# Problema de sobre-generalização



$$a_1 = [y_1 \in [1.5, 6.5]] \wedge [y_2 \in [1.5, 5.5]]$$

**descrição inicial**

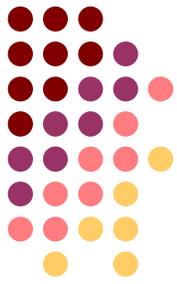
$$a_2 = [y_1 \in [3.8, 6.0]] \wedge [y_2 \in [2.5, 4.5]]$$

**descrição final**

**especialização**



# Etapa de Especialização



A idéia é encontrar um limiar  $\alpha$ -generalização de  $C_i$ , que corresponda o melhor trade-off entre a redução do volume do objeto  $s_i$  e o conjunto que satisfaz a descrição desse objeto.

Um  $\alpha$ -generalização de  $C_i$  é um objeto simbólico

$$s_i^\alpha = (a_i^\alpha, \epsilon, d_i^\alpha)$$

tal que

- $|\text{Ext}(s_i^\alpha / C_i)| \geq \alpha \times |C_i|$
- $d_i^\alpha$  tal que

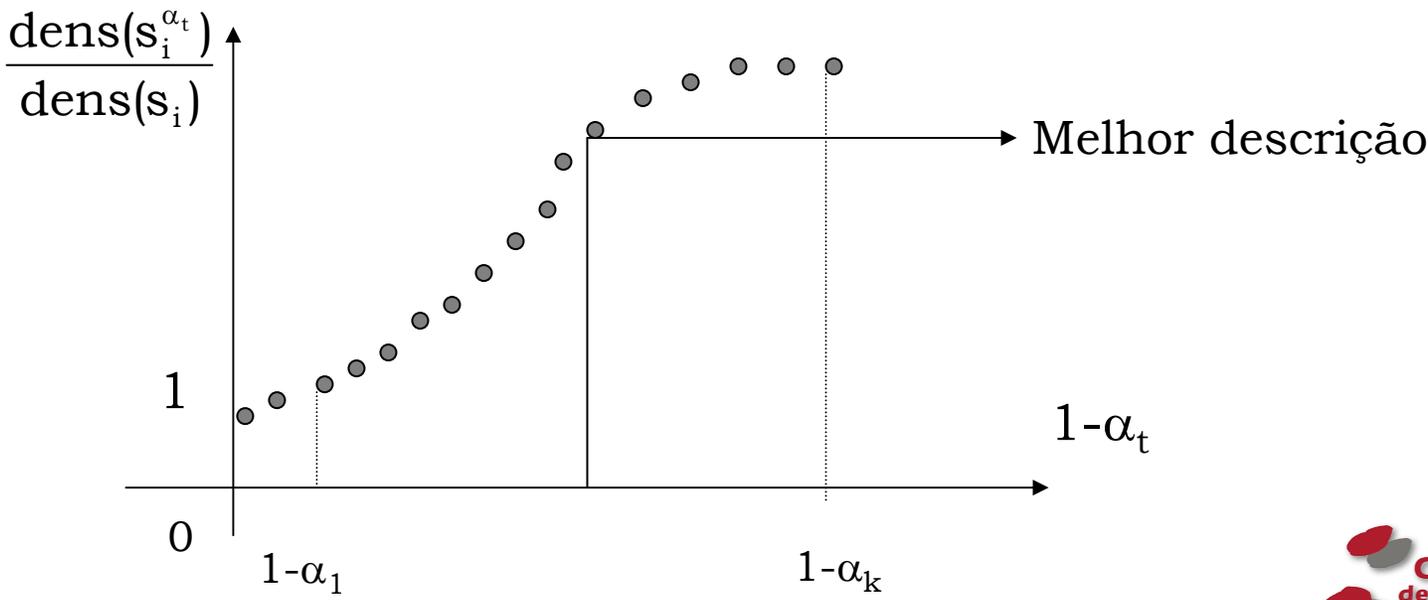
$$\text{vol}(d_i^\alpha) = \min_{d \in P(O_1) \times \dots \times P(O_p)} \left\{ \text{vol}(d) / \left| \left\{ \text{Ext}(d / C_i) \right\} \right| \geq \alpha \times |C_i| \right\}$$



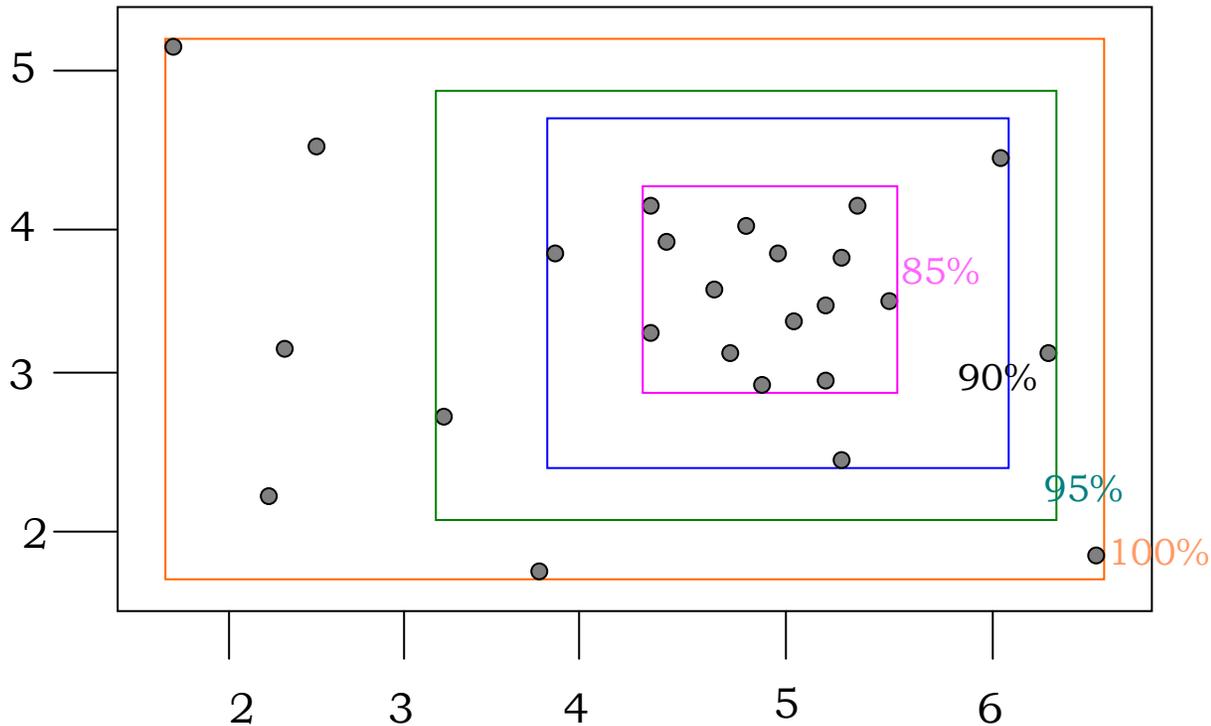
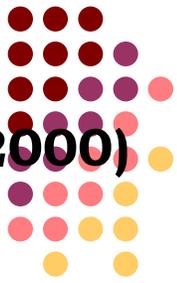
Para cada seqüência de  $\alpha_1 > \alpha_2 > \dots > \alpha_k > 0$ , pode-se definir o

$\alpha_t$  – generalização  $s_i^{\alpha_t}$

Cada  $s_i^{\alpha_t}$  é avaliado conforme a densidade relativa  $\frac{\text{dens}(s_i^{\alpha_t})}{\text{dens}(s_i)}$   
onde  $s_i$  corresponde ao objeto completo



# Exemplo do resultado do algoritmo de redução (Stephan, 2000)



$a_1 = [y_1 \in [1.5, 6.5]] \wedge [y_2 \in [1.5, 5.5]]$  **descrição inicial**

$a_1 = [y_1 \in [4.5, 5.5]] \wedge [y_2 \in [2.8, 4.3]]$  **descrição final  $\alpha = 85\%$**