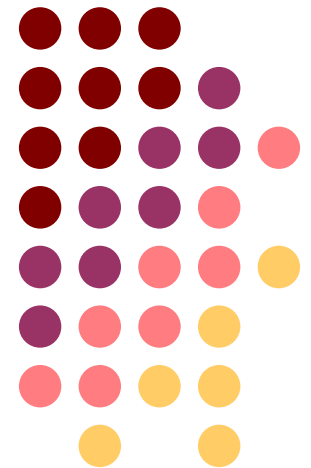


# Symbolic Data Analysis

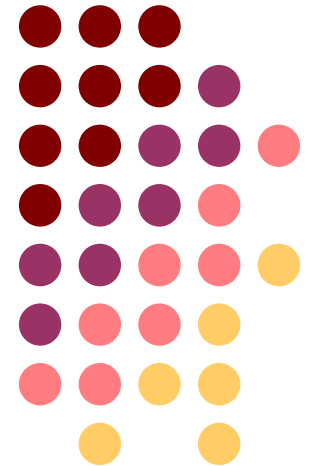
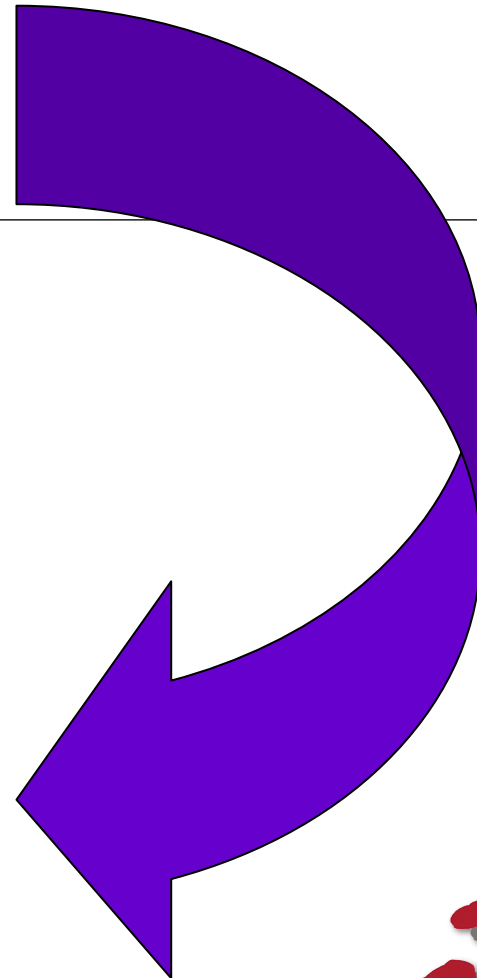
Universidade Federal de  
Pernambuco

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





# Análise de Dados Simbólicos



# Análise de Dados Simbólicos



- Surgiu em 1988 (E. Diday) e recebe influência de três grandes áreas: Análise Exploratória de Dados, Inteligência Artificial e Taxonomia Numérica.
- Nova abordagem na área de Descoberta de Conhecimento (KDD) que visa estender as técnicas estatísticas e os métodos da análise exploratória de dados para dados mais complexos chamados de Dados Simbólicos.

# Análise de Dados Simbólicos



- Novas estruturas de dados
  - Células multivaloradas;
  - Intervalos numéricos;
  - Distribuições empíricas ou de probabilidade.
- Soda's Project – Software p/ análise de dados simbólicos
  - <http://www.ceremade.dauphine.fr/~touati/sodas-pagegarde.htm>
- Livros:
  - Analysis of Symbolic Data, H.-H Bock and E. Diday, Springer-Verlag, 2000
  - Symbolic Data Analysis: Conceptual Statistics and Data Mining.y, L. Billard and E. Diday, John Wiley, 2007
  - Symbolic Data Analysis and The SODAS Software E. Diday and Monique Noirhomme-Fraiture, John Wiley, 2008

# Tabela de Dados Simbólicos



<i>Produto</i>	<i>Preço</i>	<i>Cidade</i>	<i>Cor</i>
P1	[5; 9]	Londres	{(0.1) R ; (0.3) G ; (0.6) B }
P2	[ 12 ; 15 ]	{ Paris ; Londres }	{(0.4) Y, (0.6) G}
P3	[ 3 ; 9 ]	{Bruxelas, Paris}	{ (0.3) W ; (0.7) B }
P4	[ 1 ; 8 ]	{Lisboa, Madri}	{ (0.5) W ; (0.5) B }

- **Cidade = Variável Multivalorada**
- **Cor= Variável Multivalorada Ponderada (Modal)**
- **Preço = Variável Intervalar**

# Tabela de Dados do tipo Intervalo



Temperaturas Mensais, Mínimio e Máximo, registradas em 60 estações metereológicas da China

Estações	Temperatura mensais ([min : max]) – ano 1988				
	Januay	February		November	December
AnQing	[1.8,7.1]	[2.1,7.2]		[7.8,17.9]	[4.3,11.8]
	...	...	...	...	...
ZhoJiang	[2.7,8.4]	[2.7,8.7]		[8.2;20]	[5.1,13.3]

# Alguns Métodos para Dados Simbólicos



- Análise Componente Principal e Análise Fatorial.
- Estatísticas Descritivas.
- Análise de Cluster (classificação não supervisionada).
- Análise Discriminante (classificação supervisionada).
- Modelos de Regressão
- Redes Neurais MLP

# A Idéia Básica



- **Unidades de primeira ordem** (como um cavalo ou uma pessoa qualquer), cada qual correspondendo a um único indivíduo do mundo;
- **Unidades de segunda ordem** (como o cavalo ou a pessoa, de forma geral), correspondendo a uma classe de indivíduos do mundo.



# Exemplo de Indivíduos de Primeira Ordem



Indivíduo	Classes	Variáveis Descritivas das Unidades		
<i>ID</i>	<i>Região</i>	<i>Qtd de Camas</i>	<i>Qtd de Salas de Jantar</i>	<i>Classe Social</i>
<b>1</b>	<b>Norte</b>	<b>2</b>	<b>1</b>	<b>1</b>
<b>2</b>	<b>Norte</b>	<b>2</b>	<b>1</b>	<b>3</b>
<b>3</b>	<b>Norte</b>	<b>1</b>	<b>3</b>	<b>3</b>
<b>4</b>	<b>Leste</b>	<b>1</b>	<b>3</b>	<b>3</b>
<b>5</b>	<b>Leste</b>	<b>2</b>	<b>2</b>	<b>1</b>
<b>6</b>	<b>Leste</b>	<b>1</b>	<b>2</b>	<b>3</b>

# Obtendo os Indivíduos de Segunda Ordem



Classes	Variáveis Descritivas das Unidades		
<i>Regiões</i>	<i>Qtd de Camas</i>	<i>Qtd de Salas de Jantar</i>	<i>Classe Social</i>
Norte	2	1	1
Norte	2	1	3
Norte	1	3	3
Leste	1	3	3
Leste	2	2	1
Leste	1	2	3

Classes	Variáveis Descritivas das Unidades		
<i>Regiões</i>	<i>Qtd de Camas</i>	<i>Qtd de Salas de Jantar</i>	<i>Classe Social</i>
Norte	(1/3) 1, (2/3) 2	(2/3) 1, (1/3) 3	(1/3) 1, (2/3) 3
Leste	(2/3) 1, (1/3) 2	(2/3) 2, (1/3) 3	(1/3) 1, (2/3) 3

# Exemplo de Indivíduos de Primeira Ordem – outro exemplo



✦ Em uma ilha há 600 pássaros juntos: 400 Swallows, 100 Ostriches e 100 Penguins

<i>Birds</i>	<i>Species</i>	<i>Flying</i>	<i>Size</i>
<b>1</b>	<b>Penguin</b>	<b>No</b>	<b>80</b>
.	...	...	...
.	...	...	...
.	...	...	...
<b>599</b>	<b>Swallow</b>	<b>Yes</b>	<b>70</b>
<b>600</b>	<b>Ostrich</b>	<b>No</b>	<b>125</b>

# Exemplo de Indivíduos de Primeira Ordem – outro exemplo



- Swallow bird



- Ostrich



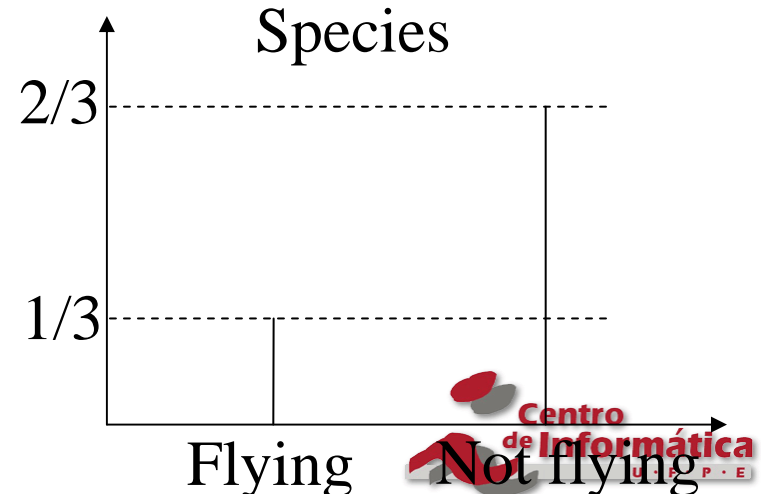
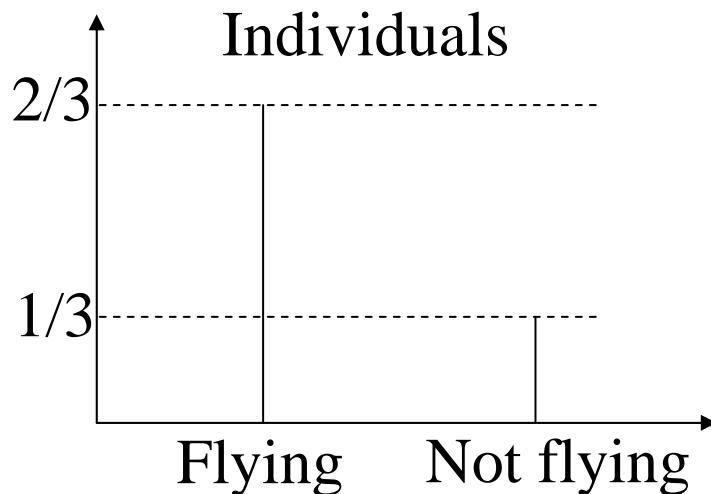
- Penguin



# Obtendo os Indivíduos de Segunda Ordem



<i>Species</i>	<i>Flying</i>	<i>Size</i>	<i>Migration</i>
<b>Swallow</b>	{Yes}	[60,85]	{ (0.1) N ; (0.9) Y }
<b>Ostrich</b>	{No}	[85,160]	{ (1.0) N ; (0.0) Y }
<b>Penguin</b>	{No}	[70,95]	{ (0.0) N ; (1.0) Y }



# Base de dados descrevendo jogadores de futebol



<i>Player</i>	<i>Team</i>	<i>Age</i>	<i>Weight</i>	<i>Height</i>	<i>Nationality</i>
<b>Fernandes</b>	<b>Spain</b>	<b>29</b>	<b>85</b>	<b>1.84</b>	<b>Spanish</b>
<b>Rodríguez</b>	<b>Spain</b>	<b>23</b>	<b>90</b>	<b>1.92</b>	<b>Brazilian</b>
<b>Mballo</b>	<b>France</b>	<b>25</b>	<b>82</b>	<b>1.9</b>	<b>Senegalese</b>
<b>Zidane</b>	<b>France</b>	<b>27</b>	<b>78</b>	<b>1.85</b>	<b>French</b>

<i>Team</i>	<i>AGE</i>	<i>Weigh</i>	<i>Height</i>	<i>Nationality</i>	<i>Number of goals at the wolrd cup 1998</i>
<b>Spain</b>	<b>[23,29]</b>	<b>[85,90]</b>	<b>[1.84,1.92]</b>	<b>(0.5 Sp, 0.5 Br)</b>	<b>18</b>
<b>France</b>	<b>[21,28]</b>	<b>[85,90]</b>	<b>[1.84,1.92]</b>	<b>(0.5 Fr, 0.5 Se)</b>	<b>24</b>



# Tabela de Dados Simbólicos

- As células podem conter dados complexos
  - Valores numéricos
    - Ex. height (Tom) = 1.80
  - Intervalos
    - Ex. age (Spain) = [23,29]
  - Categóricos
    - Ex. Nationality (Deco) = {brasileira}
    - Ex. Nationality (Spain) = {brasilian, spanish, french}
  - Modal
    - Ex. Nationality (Spain) = {(0.1)brasilian, (0.8)spanish, (0.1)french}

# Mais um exemplo de Tabela de Dados Simbólicos



<i>Produto</i>	<i>Altura</i>	<i>Cidade</i>	<i>Cor</i>
P1	3.5	Londres	{ R , G , B }
P2	[ 3 , 8 ]	{ Paris , Londres }	
P3	{ P , M , G , GG }		{ (0.3) W , (0.7) B }
P4	[ (1/3) [2,3] , (2/3) [4,5] ]		

- Outros tipos de dados simbólicos

- Taxonomia (induced rules)
- Dependência hierárquica (mother-daughter variable)
- Dependência lógica

cor = "fria" ← "violeta", "azul" e "verde"

"flying" (mother variable), "speed of flying" (daughter variable), "

If age(w) <= 2 months → (height < 80 cm)





# Fontes de Dados Simbólicos

- A partir de variáveis categóricas
  - Como tipo do empregado
  - Obtido por clusterização de grandes massas de dados
- De bancos de dados
  - Consultas originando novas variáveis (*queries*)
- Do conhecimento do especialista
- De dados confidenciais
  - Para esconder informações privadas.
- De dados estocásticos
  - Distribuição de probabilidade
- De séries temporais
  - Descrevendo intervalos de tempo



# O processo de generalização

- “O processo de generalização é aplicado a um conjunto de indivíduos para produzir uma descrição simbólica”
- Exemplo para descrição da espécie “Swallow”:
  - $d = (\{\text{yes}\}, [60,85],[90\% \text{ yes}, 10\% \text{ no}])$
- Dados simbólicos
  - Cor das espécies na tabela de dados simbólicos
  - Swallows  $\rightarrow$  [branco, preto]
  - Penguins  $\rightarrow$  {branco, preto}

# Generalizando dados simbólicos a partir de dados fuzzy, imprecisos ou conjuntivos



- Dados fuzzy
  - Ex. variável numérica *altura (homem)* = 1,60 m
  - Pode ser associada ao valor
    - “pequeno” com peso 0,9
    - “médio” com peso 0,1
    - “alto” com peso 0,0
- Dados imprecisos
  - Ocorrem quando não é possível obter uma medida exata
  - Ex: é possível dizer que uma árvore tem  $10\text{m} \pm 1$
  - Significa que a altura da árvore está no intervalo [9,11]
- Dados conjuntivos
  - Ocorrem quando muitas categorias aparecem simultaneamente.
  - Ex. uma maçã pode ser vermelha e verde, ou amarela, as três cores ao mesmo tempo

# Incerteza e dados simbólicos



- Ex1. Vamos considerar que não sabemos a altura do tenista Rafael Nadal;
- Considerando que ao jogar com o tenista Roger Federer, o Nadal parece ser apenas um pouco mais baixo, e sabendo que Federer mede 1.85m.
- É plausível considerar que
  - Altura (Nadal) = [1.80, 1.85] (incerteza)
- Ex2. Considerando a altura dos dez melhores tenistas ranqueados da ATP, podemos considerar a altura como dado simbólico:
  - Altura (top 10) = [1.80, 1.85] (variabilidade)
- Os dois intervalos são iguais mas representam semânticas completamente diferentes

# Dados simbólicos a partir de dados estruturados



- Dados estruturados ocorrem quando há variáveis do tipo mother/daughter ou taxonômicas ou ainda em tabelas associadas em um SGDB
- É possível unir tabelas com poucas variáveis em comum usando dados simbólicos

# Dados simbólicos a partir de dados estruturados



<i>Escola</i>	<i>Cidade</i>	<i>Nr alunos</i>	<i>Tipo</i>	<i>Nível</i>
Jaurès	Paris	320	Public	1
Condorcet	Paris	450	Public	3
Chevreur	Lyon	200	Public	2
St Helene	Lyon	380	Private	3
St Sernin	Toulouse	290	Public	1
St Hilare	Toulouse	210	Private	2

Descrição clássica de Escolas

<i>Cidade</i>	<i>Nr Alunos</i>	<i>Tipo</i>	<i>Nível</i>
Paris	[320,450]	(100) public	{1,3}
Lyon	[200,380]	(50%) public, (50%) private	{2,3}
Toulouse	[210,290]	(50%) public, (50%) private	{1,2}

Descrição simbólica de cidades pela variável escola - generalização

# Dados simbólicos a partir de dados estruturados



<i>Hospital</i>	<i>Cidade</i>	<i>Nr Leitos</i>	<i>Código da especialidade</i>
Lariboisiere	Paris	750	5
St Louis	Paris	1200	3
Herriot	Lyon	650	3
Besgenettes	Lyon	720	2
Purpan	Toulouse	520	6
Marchant	Toulouse	450	2

Descrição clássica de Hospitais

<i>Cidade</i>	<i>Nr Leitos</i>	<i>Código da especialidade</i>
Paris	[750,1200]	{3,5}
Lyon	[650,720]	{2,3}
Toulouse	[450,520]	{2,6}

Descrição simbólica de cidades pela variável hospital - generalização

# Dados simbólicos a partir de dados estruturados



<i>Cidade</i>	<i>Nr Alunos</i>	<i>Tipo</i>	<i>Nível</i>
Paris	[320,450]	(100) public	{1,3}
Lyon	[200,380]	(50%) public, (50%) private	{2,3}
Toulouse	[210,290]	(50%) public, (50%) private	{1,2}

<i>Cidade</i>	<i>Nr Leitos</i>	<i>Código da especialidade</i>
Paris	[750,1200]	{3,5}
Lyon	[650,720]	{2,3}
Toulouse	[450,520]	{2,6}

<i>Cidade</i>	<i>Nr alunos</i>	<i>Tipo</i>	<i>Nível</i>	<i>Nr Leitos</i>	<i>Código da especialidade</i>
Paris	[320,450]	(100) public	{1,3}	[750,1200]	{3,5}
Lyon	[200,380]	(50%) public, (50%) private	{2,3}	[650,720]	{2,3}
Toulouse	[210,290]	(50%) public, (50%) private	{1,2}	[450,520]	{2,6}



# Os quatro tipos de estatísticas e “*data mining*”



Análise Clássica

Análise Simbólica

Dados Clássicos

Caso 1

Caso 2

<i>Cidade</i>	<i>Min. Alunos</i>	<i>Máx. alunos</i>	<i>Public</i>	<i>Private</i>	<i>Nível 1</i>	<i>Nível 2</i>	<i>Nível 3</i>
Paris	320	450	100	0	1	0	1
Lyon	200	380	50	50	0	1	1
Toulouse	210	290	50	50	1	1	0

Formação de tabelas de dados clássicos  
Análise de dados simbólicos sobre dados simbólicos.

Aplicação de estatística clássica.  
Muita informação é perdida

# SODAS – Symbolic Official Data Analysis System



- Protótipo disponível gratuitamente
  - <http://www.info.fundp.ac.be/asso/>
  - Funcionalidades
  - Construção de tabelas de dados simbólicos a partir de BD's tradicionais
  - Descrição de regras e hierarquias
  - Análise dos dados através de métodos de análise de dados simbólicos
    - Estatística descritiva
    - Análise Fatorial
    - Agrupamento
    - Árvore de Decisão
    - ...

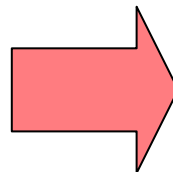
# CLASSIFICADOR SIMBÓLICO APLICADO A IMAGENS SAR



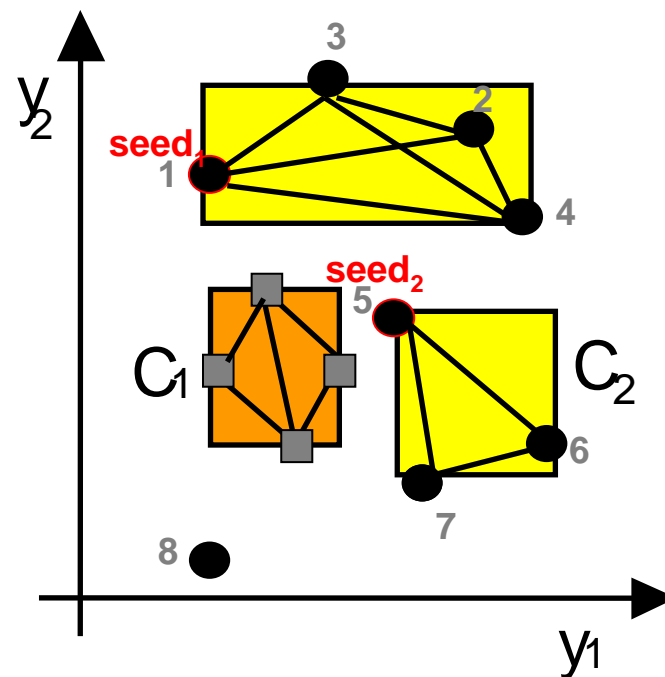
## Etapa de Aprendizagem

Entrada: Dados usuais

segmento	níveis de cinza	
	media ( $y_1$ )	desvio padrão ( $y_2$ )
seg <sub>1</sub>	44.50	12.50
seg <sub>2</sub>	83.60	3.60
seg <sub>3</sub>	120.30	6.45
•	•	•
•	•	•



Aproximação do Grafo de  
Vizinhos Mútuos





## Saída: Dados Simbólicos que descrevem as regiões (grupos de segmentos)

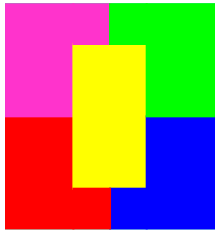
Grupo	Sub-Grupo	níveis de cinza	
		media ( $y_1$ )	desvio padrão ( $y_2$ )
Região 1	$G_{11}$	[37.35,57.70]	[0,20,5.62]
Região 2	$G_{21}$	[132.56,160.79]	[0.73,6.84]
•	$G_{22}$	[167.12,196.67]	[1.30,10.66]
•	•	•	•
•	•	•	•

## Etapa de Alocação: Funções de Proximidade

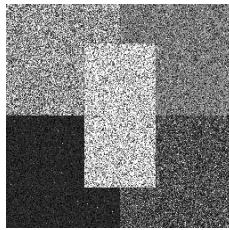
1 - Distância (De Carvalho et al (1998))  $d_r(s,s') = \frac{1}{p} \left[ \sum_{i=1}^p \left\{ \Phi_{jk\gamma} (s_i, s'_i) \right\}^r \right]^{1/r}$

2 - Palumbo et al (1996)  $q_{ss'} = \frac{\pi(s \oplus s') - \pi(s)}{\pi(s)}, \pi(s) \neq 0$

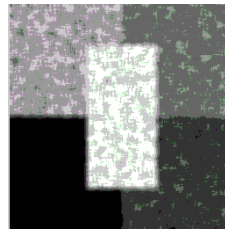
# EXPERIMENTO MONTE CARLO



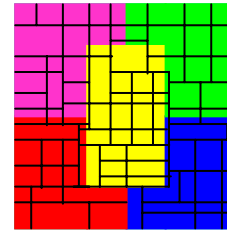
PHANTOM



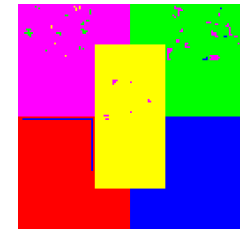
MULTIPLICATIVE  
MODEL  
(Frery et al 1997)



LEE  
FILTER

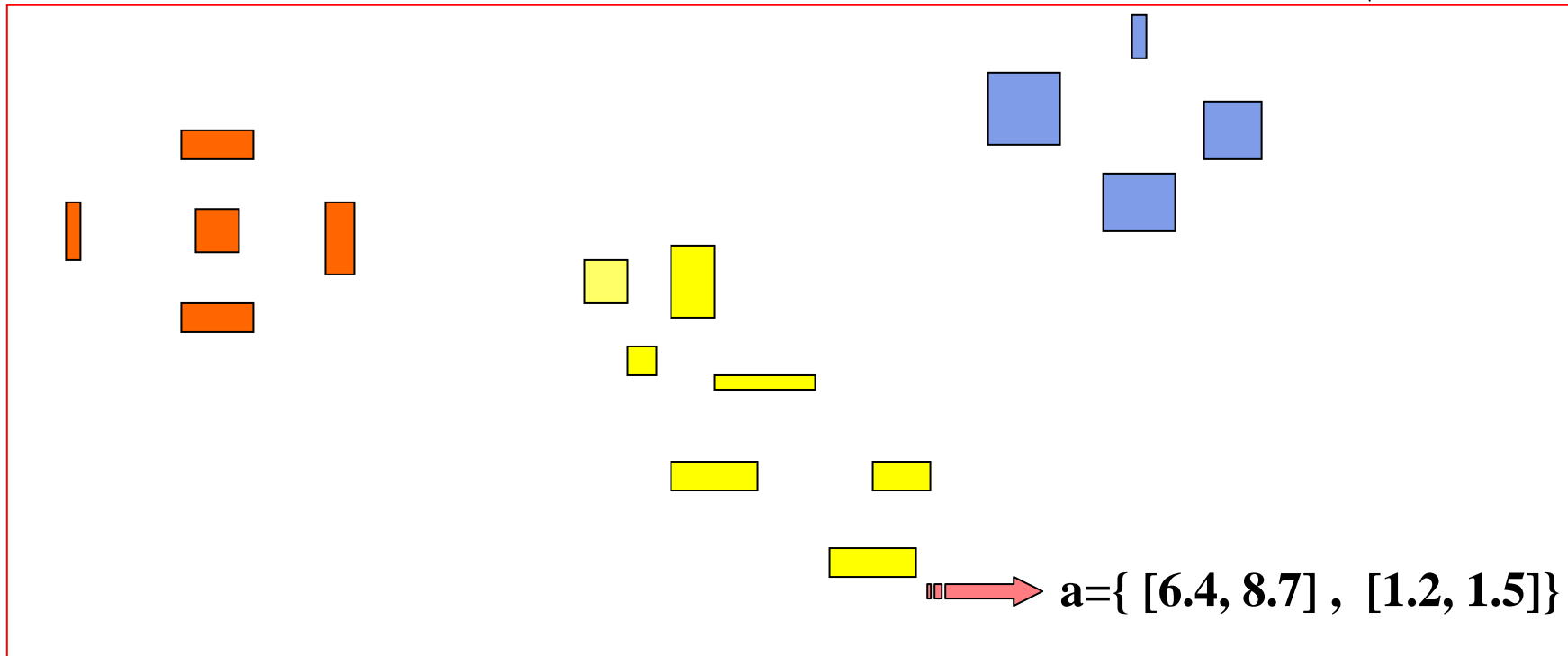


SEGMENTATION  
REGION GROWTH



SYMBOLIC  
CLASSIFIER

# Análise de Cluster

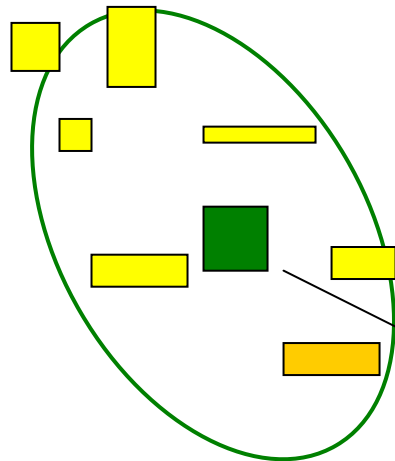
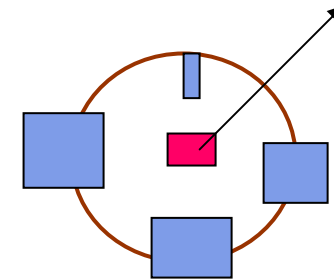


O algoritmo de **nuvens dinâmicas** visa encontrar uma partição e um conjunto de representantes das classes otimizando um critério de ajustamento entre classes e seus representantes. **As distâncias adaptativas** permitem encontrar clusters de formas e tamanhos diferentes.

# Saída: A partição e as descrições das classes



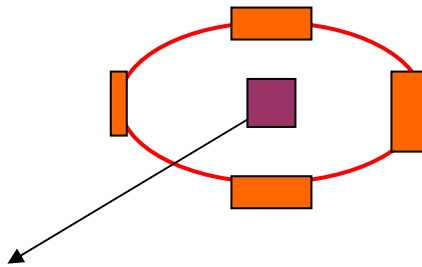
$C_3: ([18.4, 19.6], [9.7, 10.5])$



$C_2: ([13.4, 15.3], [3.7, 5.5])$

protótipo

$C_1: ([3.4, 4.3], [6.7, 7.5])$



Critério a ser otimizado:

$$W = \sum_{i=1}^k \sum_{x \in C_i} d_i(x, G_i)$$

Distâncias Adaptativas

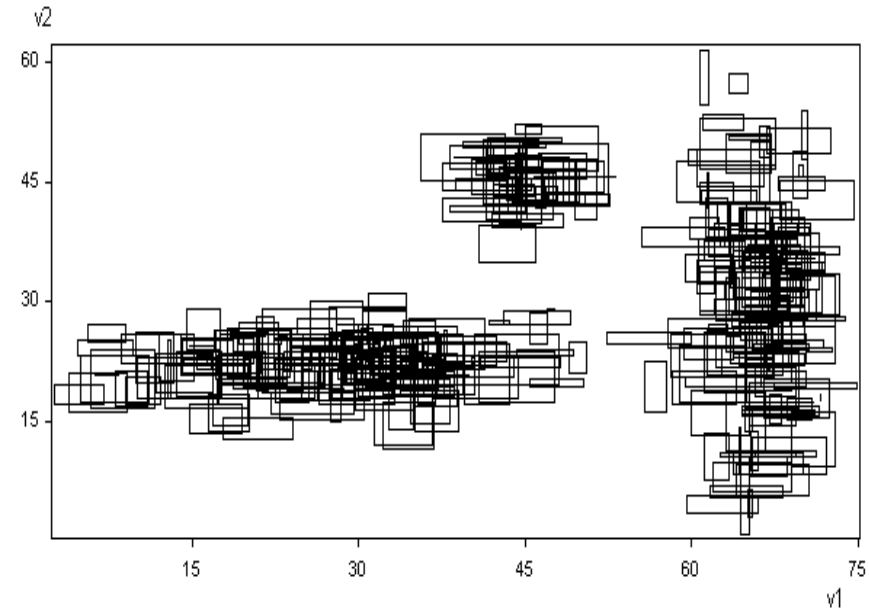
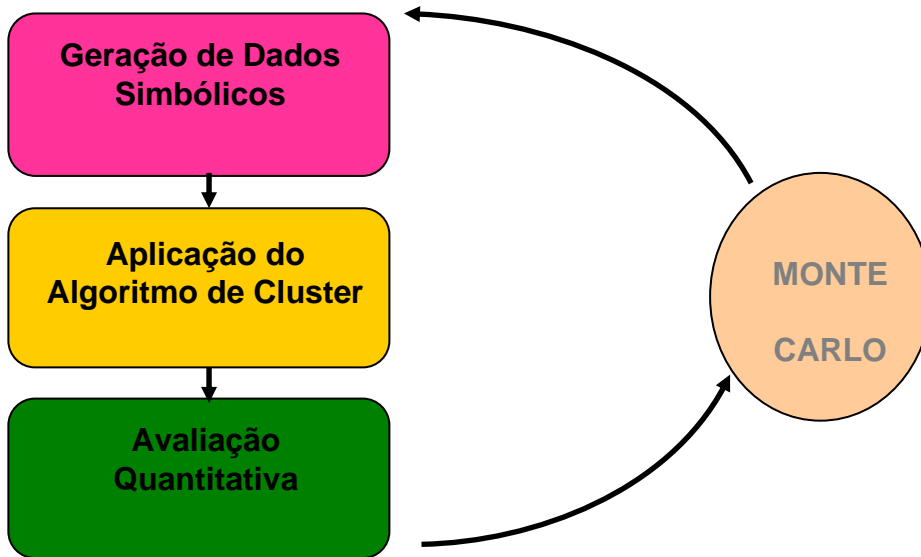
City-Block

Euclideana

Chebyshev

Mahalanobis

# Método Monte Carlo



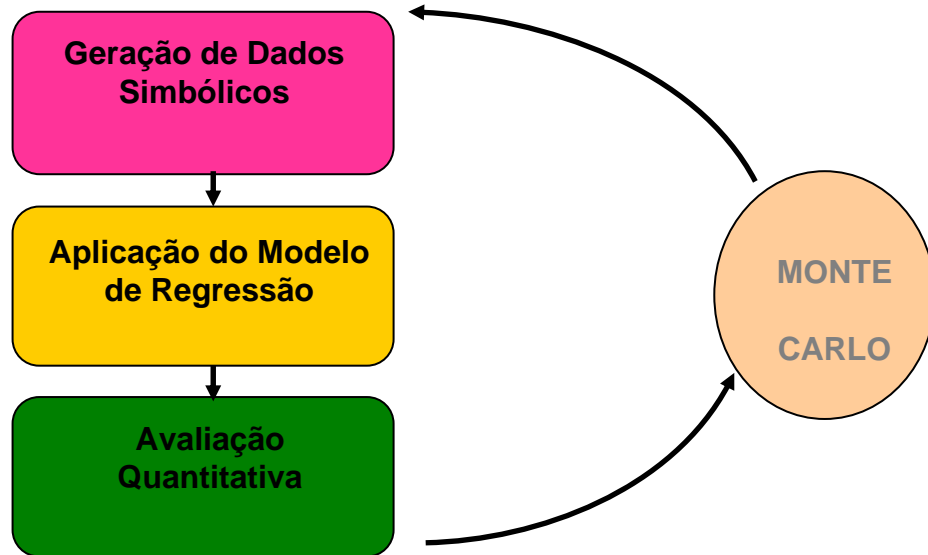
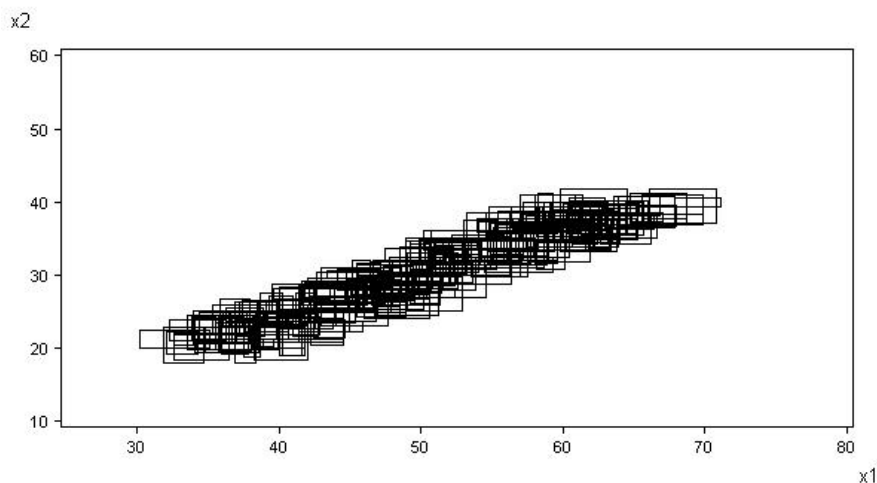


# Modelos de Regressão para Dados Simbólicos Tipo-Intervalo



e	Pulso (Y)	Pressão Sistólica (X <sub>1</sub> )	Pressão Diastólica (X <sub>2</sub> )
1	[44-68]	[90-100]	[50-70]
2	[60-72]	[90-130]	[70-90]
3	[56-90]	[140-180]	[90-100]
4	[70-112]	[110-142]	[80-108]
5	[54-72]	[90-100]	[50-70]
6	[70-100]	[130-160]	[80-110]
7	[63-75]	[60-100]	[140-150]
8	[72-100]	[130-160]	[76-90]
9	[76-98]	[110-190]	[70-110]
10	[86-96]	[138-180]	[90-110]
11	[86-100]	[110-150]	[78-100]

# Problema de regressão: dados intervalares simulados



# Analizando Gestões Administrativas Municipais



- **Motivação:**
  - Agrupar, por meio de uma abordagem simbólica, cidades que apresentem gestões administrativas similares.
  - Traçar um mapa administrativo de cada município identificando, mais detalhadamente, seus defeitos e suas virtudes
  - Utilizar a opinião de seus habitantes sobre alguns serviços públicos municipais

# As cidades



BJD	STL	GRN	ARC	BEZ
PLT	AEL	IGA	SCC	RCF
JDG	OLI	PET	CAM	CPN
CAR	CSA	VSA	SLM	GOI
ARA	GVT	IPJ	PES	ESC

# As Variáveis de Interesse



<b>Serviços Públicos Municipais (variáveis)</b>	<b>Codificação</b>
Limpeza de ruas e avenidas	<b>Limpeza</b>
Recolhimento do lixo domiciliar	<b>Lixo</b>
Iluminação pública	<b>Iluminação</b>
Ensino municipal	<b>Educação</b>
Assistência à população pobre	<b>Ass_pobre</b>
Conservação de praças, parques e jardins	<b>Conserv</b>
Assistência médica municipal	<b>Saúde</b>
Promoção de festas populares	<b>Festa</b>
Abastecimento de água	<b>Agua</b>
Rede de esgoto/saneamento	<b>Saneam</b>
Segurança	<b>Segurança</b>
Pavimentação de ruas e avenidas	<b>Paviment</b>
Apoio à geração de empregos	<b>Emprego</b>
Conservação de estradas	<b>Estrada</b>
Trânsito	<b>Transito</b>
Avaliação Administrativa	<b>Aval</b>

16 variáveis

# Os Dados



- Escala de Avaliação de Desempenho das Variáveis de Interesse

Escore	1	2	3	4	5
Categoria	Péssimo	Ruim	Regular	Bom	Ótimo

# O BD clássico



ID	Município*	Limpeza	Lixo	...	Transito	Aval
1	BJD	Regular	Bom	...	Ruim	Ruim
2	BJD	Bom	Otimo	...	Bom	Bom
...	...	...	...	...	...	...
231	BJD	Bom	Otimo	...	Bom	Bom
232	STL	Pessimo	Bom	...	Ruim	Regular
...	...	...	...	...	...	...
5.241	ESC	Ruim	Bom	...	Bom	Péssimo

89.097 células

# O BD simbólico



	Limpeza	..	Trânsito	Saldo
<b>JDG</b>	Péssi(0.29), Ruim(0.15), Regul(0.26), Bom(0.26), Ótimo(0.03)	..	Péssi(0.22), Ruim(0.17), Regul(0.37), Bom(0.23), Ótimo(0.02)	n
<b>RFC</b>	Péssi(0.13), Ruim(0.10), Regul(0.34), Bom(0.35), Ótimo(0.09)	..	Péssi(0.15), Ruim(0.12), Regul(0.34), Bom(0.37), Ótimo(0.03)	n
<b>OLI</b>	Péssi(0.31), Ruim(0.13), Regul(0.27), Bom(0.24), Ótimo(0.05)	..	Péssi(0.16), Ruim(0.09), Regul(0.34), Bom(0.38), Ótimo(0.03)	n
<b>PLT</b>	Péssi(0.14), Ruim(0.07), Regul(0.24), Bom(0.41), Ótimo(0.14)	..	Péssi(0.12), Ruim(0.06), Regul(0.47), Bom(0.33), Ótimo(0.02)	p
<b>CAR</b>	Péssi(0.07), Ruim(0.04), Regul(0.24), Bom(0.49), Ótimo(0.15)	..	Péssi(0.19), Ruim(0.06), Regul(0.23), Bom(0.48), Ótimo(0.04)	n
....	....	..	....	....
<b>ESC</b>	Péssi(0.21), Ruim(0.06), Regul(0.31), Bom(0.32), Ótimo(0.09)	..	Péssi(0.21), Ruim(0.14), Regul(0.22), Bom(0.41), Ótimo(0.02)	n

400 células



# Os Dados



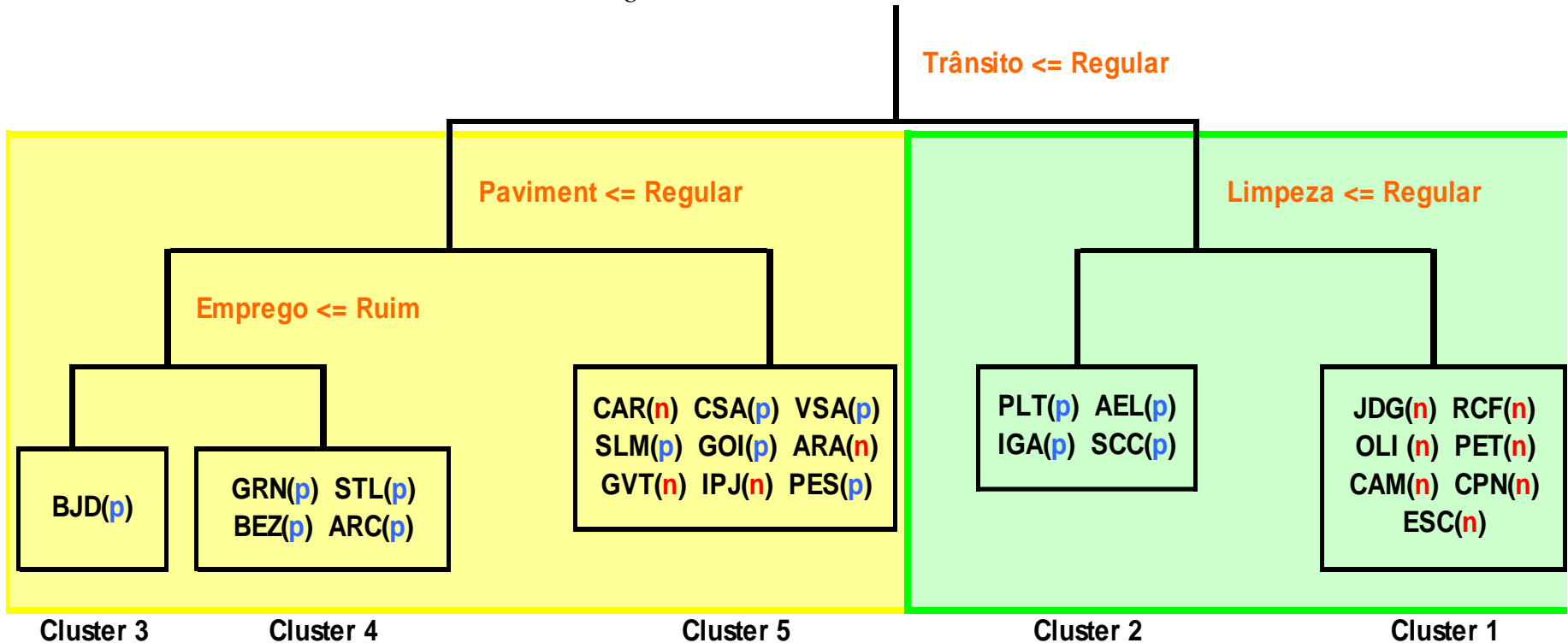
## Novas Variáveis Simbólicas

- Saldo = {n, p}
  - Avaliação Administrativa
  - Se  $[(\%Ótimo + \%Bom) - (\%Ruim + \%Péssimo)] > \%Regular$ 
    - então Saldo = p (gestão adm. positiva)
    - cc. Saldo = n (gestão adm. negativa)
- Variáveis Multivaloradas Ponderadas (Serviços Públicos)
  - Limpeza, Lixo, Iluminação, Educação, Assistência à pobreza, Conservação de praças, Saúde, Festas, Água, Saneamento, Segurança, Pavimentação, Emprego, Estrada e Trânsito.

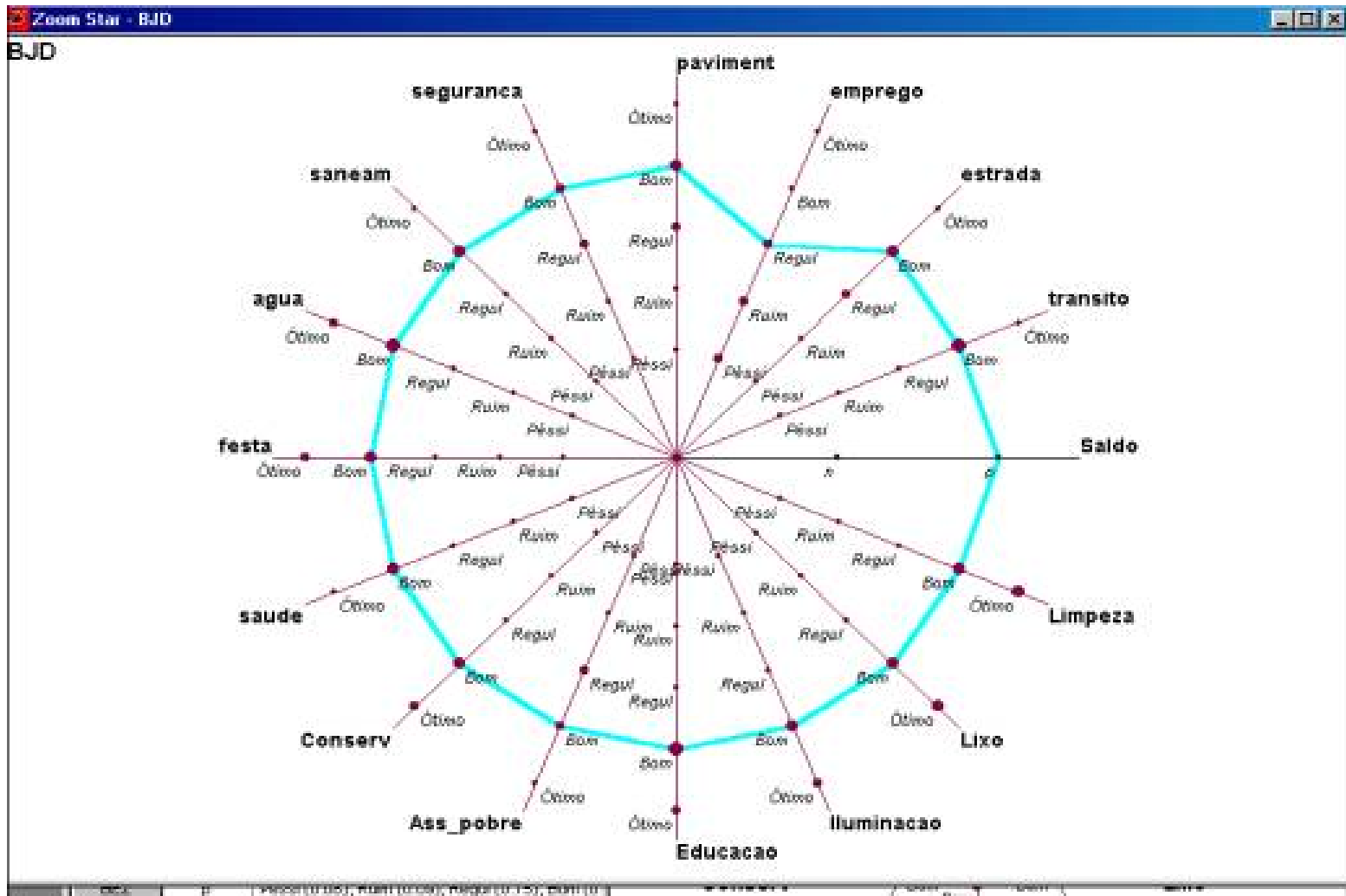
# Métodos Utilizados : A árvore de Cluster



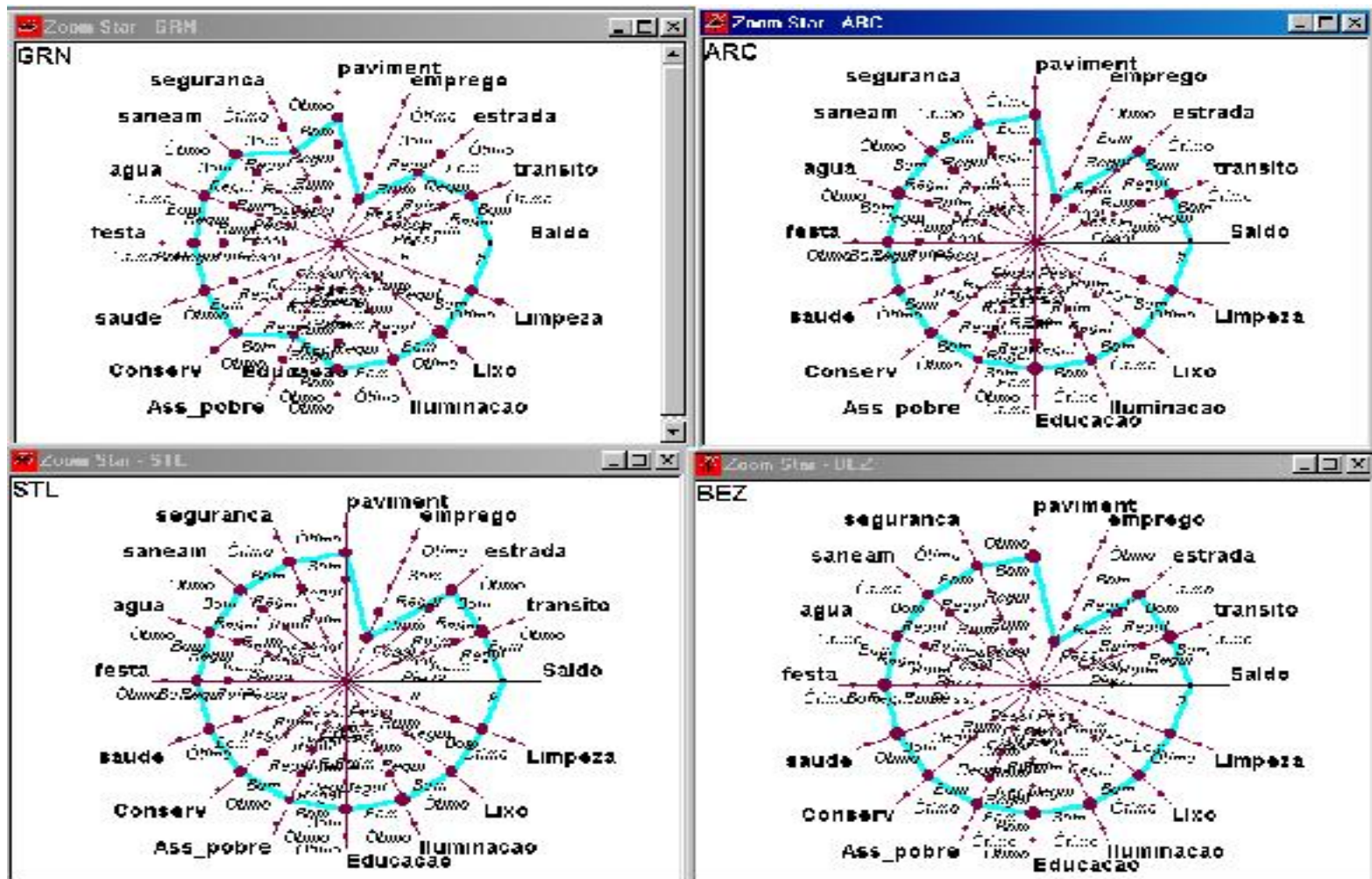
Figura 1: A Árvore de Cluster



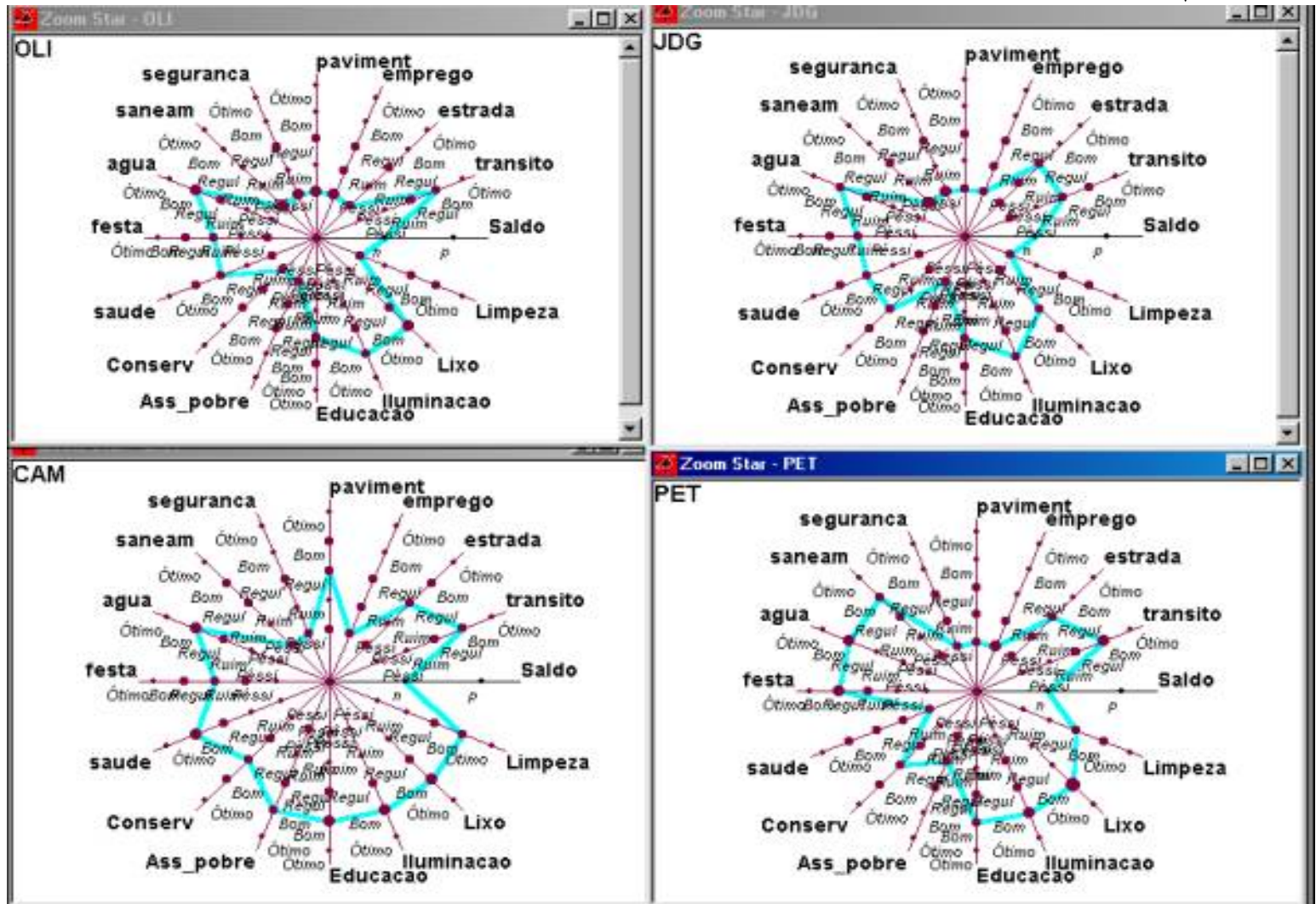
# Métodos Utilizados: Zoom Stars 2D — cluster 3



# Métodos Utilizados : Zoom Stars 2D — cluster 4



# Zoom Stars 2D — cluster 1



# Conclusões



- Utilizando a **Análise de Cluster** foi possível identificar padrões de gestões administrativas;
- Através dos gráficos **Zoom Stars 2D** identificamos os pontos positivos e negativos de cada gestão;
- As estrelas maiores e mais uniformes implicam em gestões positivas. Já estrelas menores e/ou com diversas deformidades caracterizam gestões negativas.

# Conclusões



- Uma gestão administrativa negativa foi caracterizada pelas cidades que obtiveram uma classificação *Regular*, *Ruim* ou *Péssima* em Trânsito e Limpeza de Ruas ou Avenidas;
- Uma gestão administrativa positiva foi caracterizada pelas cidades que obtiveram uma classificação *Boa* ou *Ótima* em Trânsito e Pavimentação de ruas ou avenidas;
- A gestão administrativa da cidade BJD mereceu destaque em virtude da classificação *Regular* na variável Apoio a geração de empregos.