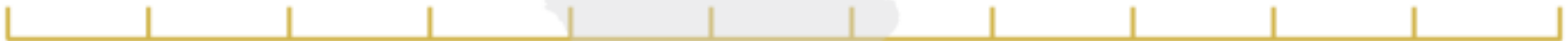




# Aplicación del Análisis de Datos Simbólico a la Detección de Inusualidades

*Dr. Oldemar Rodríguez Rojas*



# Inusualidad o atipicidad



# Objetos Simbólicos:

- “Real-life objects are too **complex** to be represented by points in a vectorial space” [Bock&Diday, 2000]
- “Symbolic objects overcome this limitation by representing **concepts** rather than individuals” [Bock&Diday, 2000]
- “Knowledge extraction from large data bases is our main aim as in Data Mining” [Diday, 1998]
- En análisis simbólico de datos transforma la Minería de Datos una herramienta para “Big Data”.

# Tabla clásica:

*Example 4.1.2: A classical data matrix with mixed variables*

The classical data matrix in Table 4.2 describes characteristics of  $N = 6$  students in terms of  $p = 4$  variables:

<i>Person</i> $u$	<i>Height</i> $Y_1$	<i>Weight</i> $Y_2$	<i>Subject</i> $Y_3$	<i>Sex</i> $Y_4$
<i>Anna</i>	1.70	65.9	<i>math</i>	1
<i>Berta</i>	1.62	61.4	<i>med</i>	1
<i>Claudia</i>	1.65	60.1	<i>phys</i>	1
<i>Daniel</i>	1.75	77.3	<i>math</i>	0
<i>Eric</i>	1.80	74.3	<i>med</i>	0
<i>Fred</i>	1.68	67.0	<i>econ</i>	0

**Table 4.2:** A classical data matrix  $\mathcal{X} = (x_{uj})$  with mixed variables

# Symbolic Data Table:

$$\underline{X} = \left( \begin{array}{c|c|c} \begin{matrix} (80, 100] \\ (100, 130] \\ (8, 10] \\ (10, 13] \end{matrix} & \begin{matrix} (D \ 0.4; \ C \ 0.3; \ S \ 0.2; \ N \ 0.1) \\ (D \ 0.1; \ C \ 0.3; \ S \ 0.4; \ N \ 0.2) \\ (D \ 0.3; \ C \ 0.5; \ S \ 0.1; \ N \ 0.1) \\ (D \ 0.3; \ C \ 0.3; \ S \ 0.3; \ N \ 0.1) \end{matrix} & \begin{matrix} \{CL, BNP\} \\ \{SPK, DB, BNP\} \\ \{SPK\} \\ \{CL, DB, BNP\} \end{matrix} \end{array} \right) \begin{array}{l} \leftarrow x'_1 \\ \leftarrow x'_2 \\ \leftarrow x'_3 \\ \leftarrow x'_4 \end{array}$$

where, e.g., the third row

$$x'_3 = ((8, 10], (D \ 0.3; \ C \ 0.5; \ S \ 0.1; \ N \ 0.1), \{SPK\})$$

describes the town  $a_3$  and could be equivalently displayed in a semi-graphical form such as



# Creando tablas simbólicas:

Classical description of Schools

Schools	Town	Nb of pupils	Kind	Level
Jaurès	Paris	320	Public	1
Condorcet	Paris	450	Public	3
Chevreur	Lyon	200	Public	2
St Hélène	Lyon	380	Private	3
St Sernin	Toulouse	290	Public	1
St Hilaire	Toulouse	210	Private	2

Symbolic description of the towns by the schools variables

Town	Nb of pupils	Kind	Level
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}

# ¿Cómo se construyen las tablas simbólicas?

Millones...

<i>Id-trx</i>	<i>Causal</i>	<i>Sucursal</i>	<i>Monto</i>	<i># Tarjeta</i>
3457	36	Curridabat	2,500.00	1000
1251	28	San Pedro	1,750.00	1001
3245	39	Grecia	2,400.00	1000
7635	35	San Pedro	1,900.00	1001
3245	35	Alajuela	1,850.00	1001
5367	27	Alajuela	1,900.00	1002
6486	34	Heredia	1,600.00	1002

Análisis  
Clásico

Datos



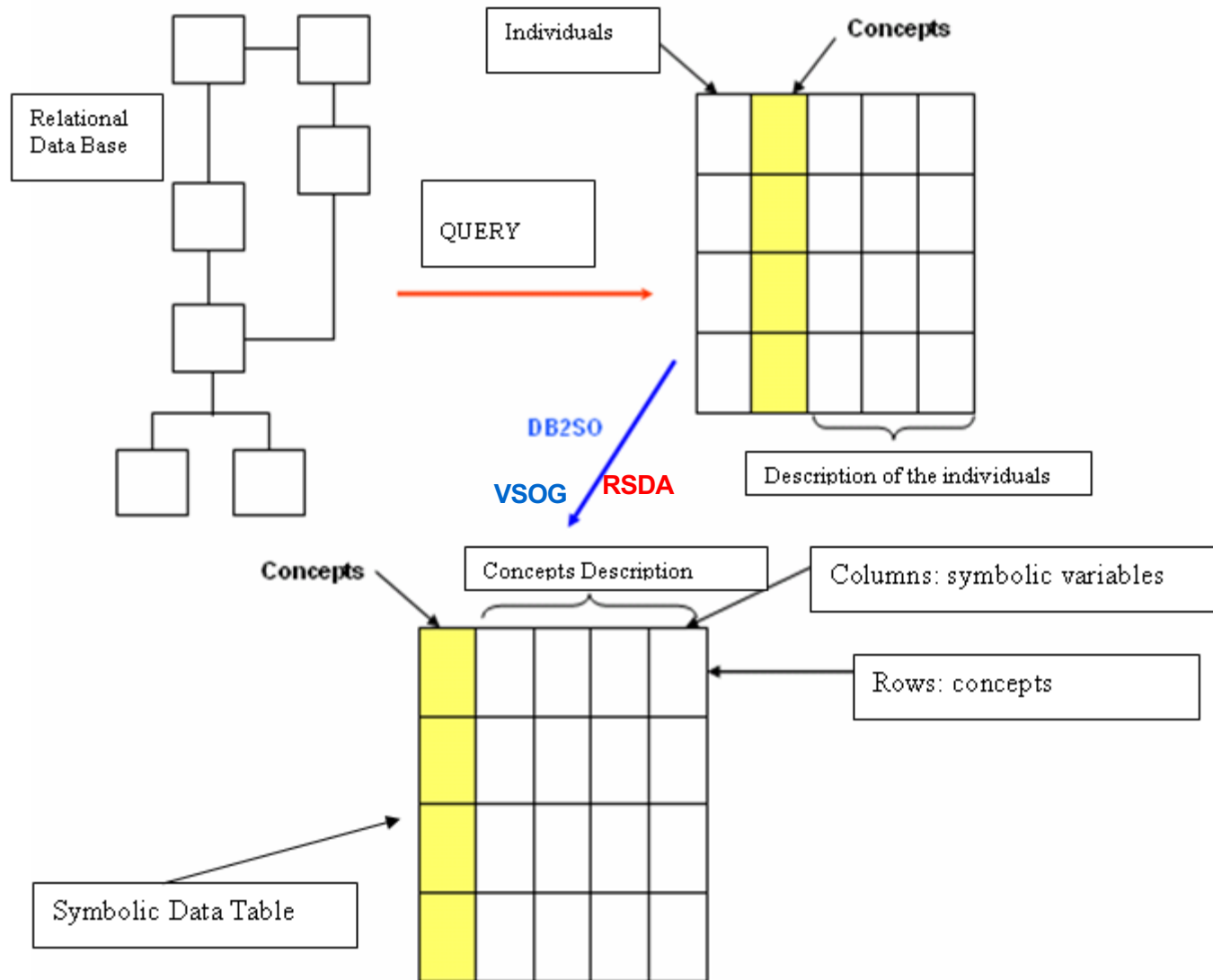
Análisis  
Multivariado  
Simbólico

Conceptos

Cientos...

<i># Tarjeta</i>	<i>Causal</i>	<i>Sucursal</i>	<i>Monto</i>	<i>Salario</i>
1000	36(1/2),39(1/2)	{Curr-50%,Gre-50%}	[2.4,2.5]	255.4
1001	28(1/3),35(2/3)	{SP-66%,Al-33%}	[1.75,1.9]	122,2
1002	27(1/2),34(1/2)	{Al-50%,Her-50%}	[1.6,1.9]	534,5

# From a Relational Database to a symbolic data table:





# Symbolic Data Base:

**Relational** Data Base

100% knowledge

15 Gigabyte



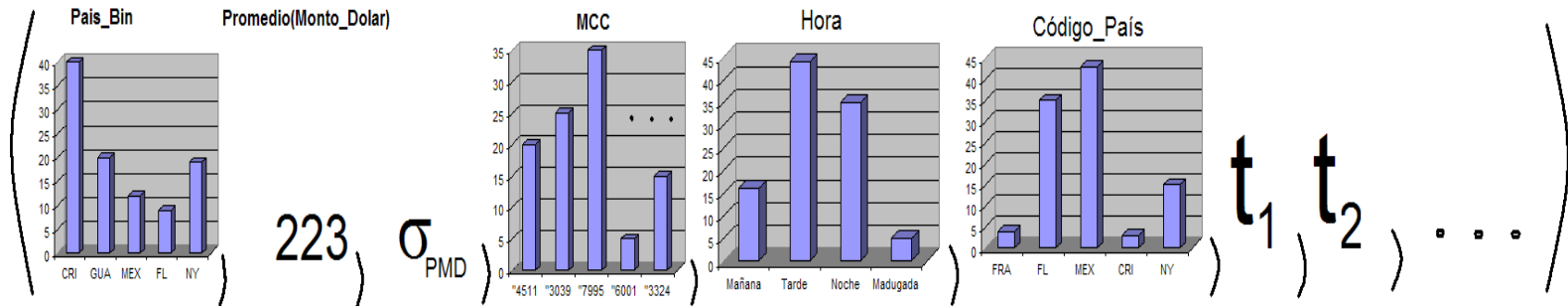
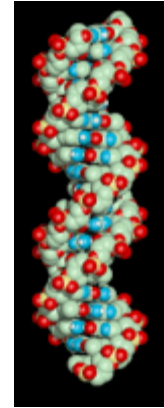
**Symbolic** Data Base

90 % knowledge

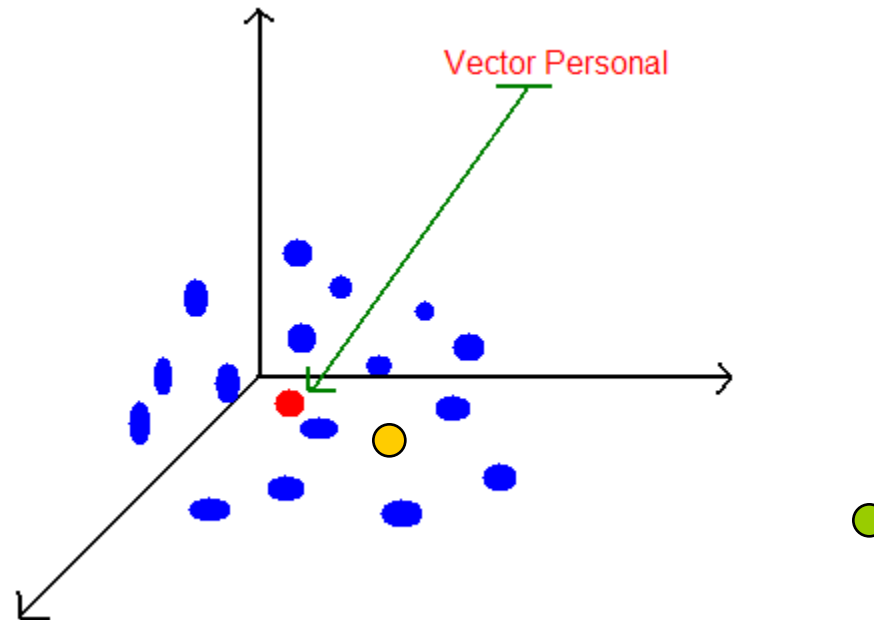
10.3 Megabyte



# Perfil simbólico

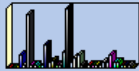


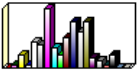
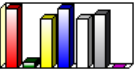







# Espacio con el Perfil Simbólico y sus Transacciones

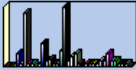


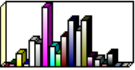
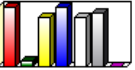



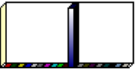



# Transacción Inusual

Comparación Vector Personal vs Transacción

Concept	TablaAprendizajeAbril.BNumber	TablaAprendizajeAbril.duration	Interval.TablaAprendizajeAbril.duration	TablaAprendizajeAbril.indInternacional	TablaAprendizajeAbril.hora	TablaAprendizajeAbril.diaSemana	Distancia Maxima	Desviacion Estandar Distancia Hausdorff	Promedio Distancias	Afinidad Ma
8102148			[0 , 240]				1.0396	0.0639	0.7624	0.8397
La Transaccion			[10 , 10]				1.0326	0	1.0326	0.8665

# Transacción Usual

Concept	TablaAprendizajeAbril.BNumber	TablaAprendizajeAbril.duration	Interval.TablaAprendizajeAbril.duration	TablaAprendizajeAbril.indInternacional	TablaAprendizajeAbril.hora	TablaAprendizajeAbril.diaSemana	Distancia Maxima	Desviacion Estandar Distancia Hausdorff	Promedio Distancias	Afinidad Media
8102148			[0 , 240]				1.0396	0.0639	0.7624	0.8397
La Transaccion			[39 , 39]				0.7206	0	0.7206	0.5066

# Distancias entre Objetos Simbólicos



FIG. 4.3 – *Comparaison d'intervalles*

# Distancias entre Objetos Simbólicos

Notation	Intervalle	Qualitatif
$A_k$	$[a_l, a_u]$	$a_1, \dots$
$B_k$	$[b_l, b_u]$	$b_1, \dots$
inters ou $ A_k - -B_k $	longueur $(A_k \cap B_k)$	$\text{card}(A_k \cap B_k)$
$l_s$ ou $ A_k + +B_k $	$ \max(a_u, b_u) - \min(a_l, b_l) $	$l_a + l_b - \text{inters}$
$l_a$ ou $ A_k $	$ a_u - a_l $	$\text{card}(A_k)$
$l_b$ ou $ B_k $	$ b_u - b_l $	$\text{card}(B_k)$
$ U_k $	longueur du domaine	cardinal du domaine

## Mesure proposée par Ichino [Ichino88]

En reprenant les notations du tableau 4.1 (page 146), la mesure s'écrit de la façon suivante :

$$\Phi(A_k, B_k) = |A_k + +B_k| - |A_k - -B_k| + \gamma(2|A_k - -B_k| - |A_k| - |B_k|)$$

# Distancias entre Objetos Simbólicos

## Distance de Hausdorff [Bandemer et Nather92]

Entre deux ensembles A et B, la distance s'écrit :

$$D_h(A, B) = \max\left\{\sup_{a \in A} \inf_{b \in B} d(a, b), \sup_{b \in B} \inf_{a \in A} d(a, b)\right\}$$

La mesure du Khi2 est la suivante:

$$d(1, 2) = \sum_{k=1}^p \left( \sum_{i=1}^q \frac{1}{f_{ik}^1 + f_{ik}^2} \left( \frac{f_{ik}^1}{\sum_{i=1}^q f_{ik}^1} - \frac{f_{ik}^2}{\sum_{i=1}^q f_{ik}^2} \right)^2 \right)^{\frac{1}{2}}$$

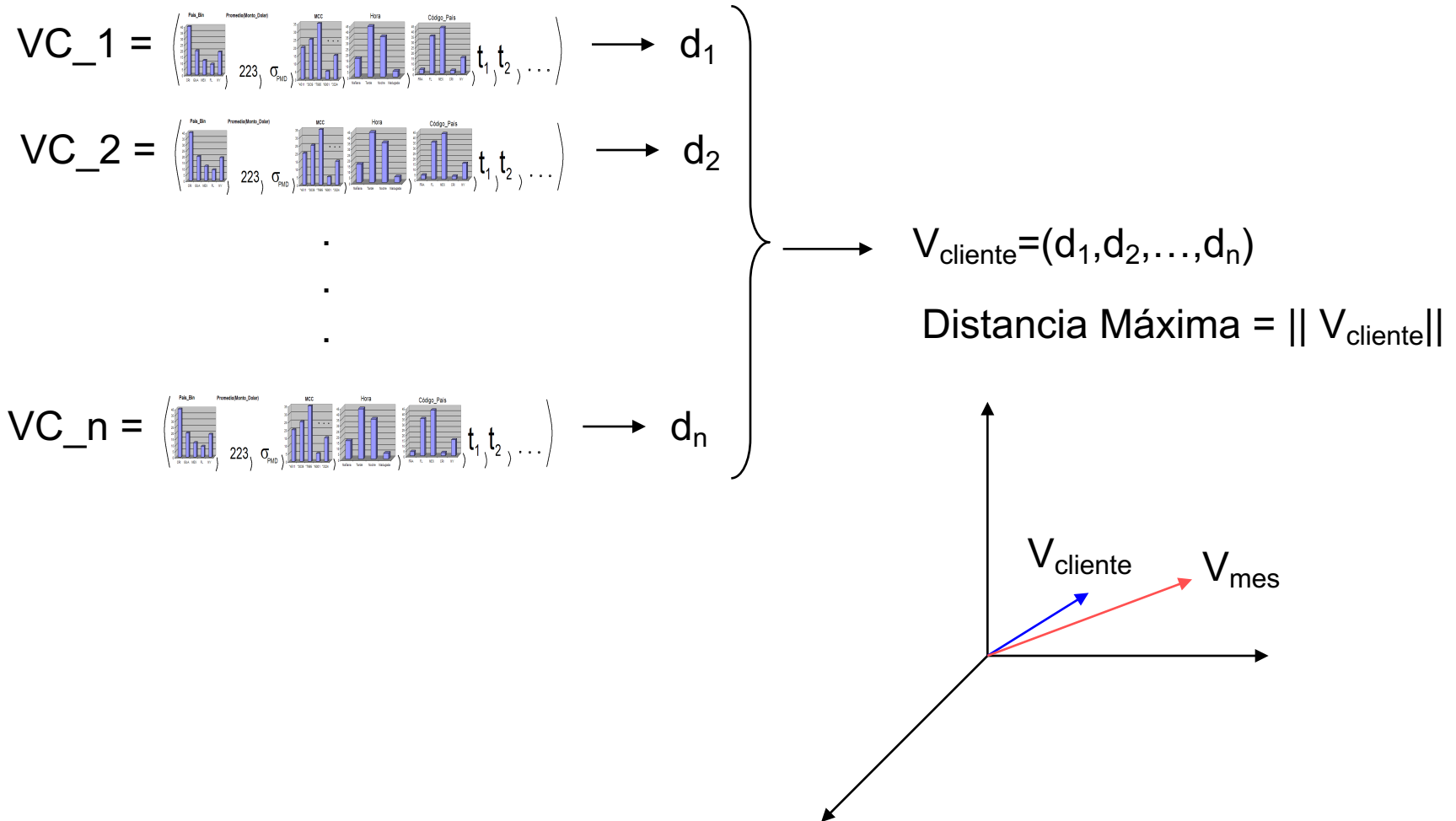


# Índice de Afinidad

Now, if we weight the importance of a variable  $Y_j$  with a weight  $w_j$  (with  $0 \leq w_j \leq 1$  and  $\sum_{j=1}^p w_j = 1$ ), we define the (weighted) **affinity similarity coefficient**  $aff(k, k') \equiv a(k, k')$  between the units (groups)  $k, k' \in E$  by the weighted average:

$$a(k, k') := \sum_{j=1}^p w_j \cdot aff(\xi_{kj}, \xi_{k'j}) = \sum_{j=1}^p w_j \cdot \sum_{l=1}^{m_j} \sqrt{\frac{n_{kjl}}{n_{kj.}} \cdot \frac{n_{k'jl}}{n_{k'j.}}}. \quad (8.68)$$

# Perfil del Cliente

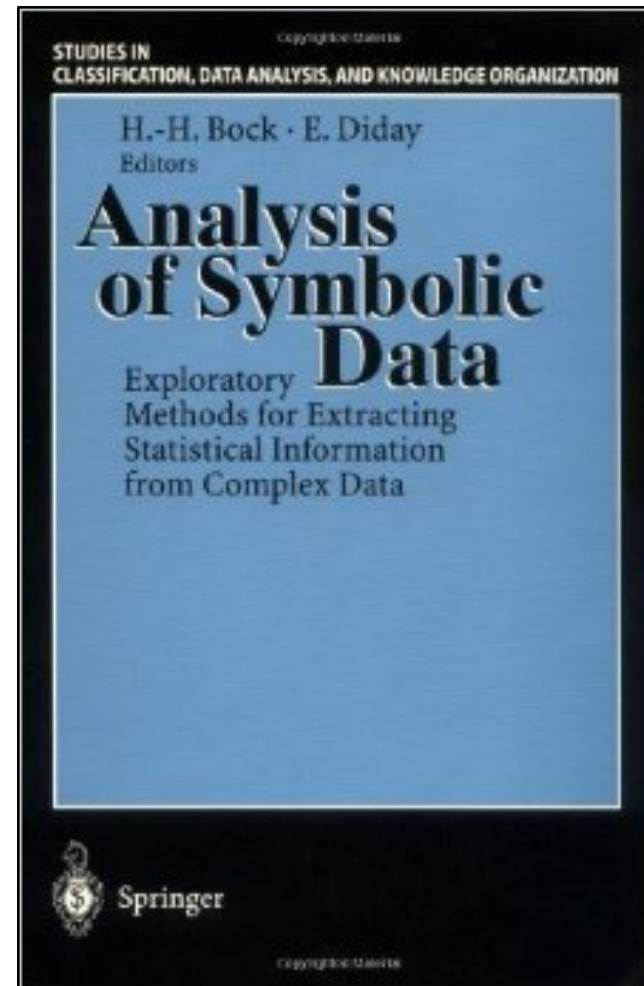
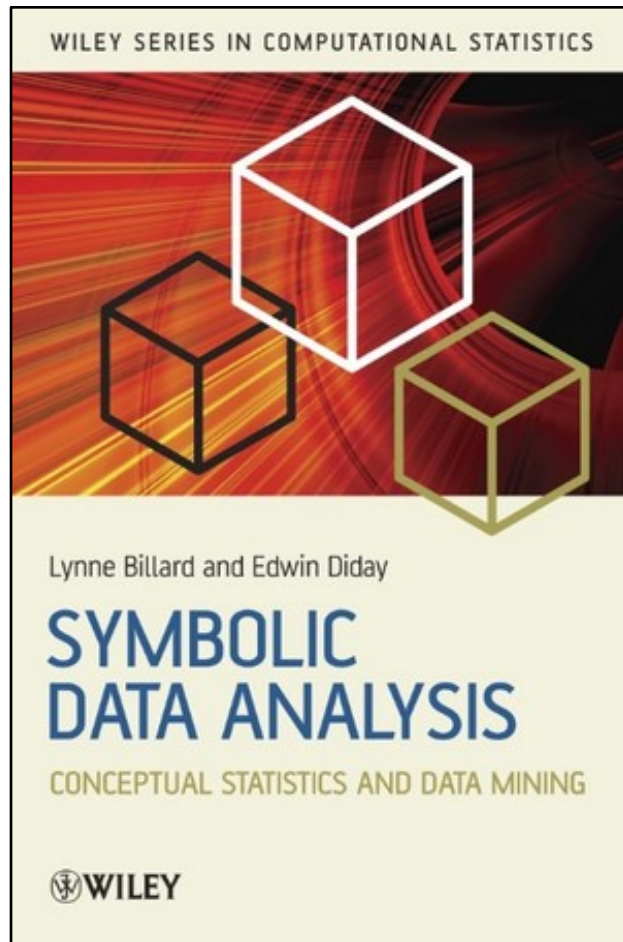


# Inusualidad o atipicidad



Concept	TablaAprendizajeAbril.BNumber	TablaAprendizajeAbril.duration	Interval.TablaAprendizajeAbril.duration	TablaAprendizajeAbril.indInternacional	TablaAprendizajeAbril.hora	TablaAprendizajeAbril.diaSemana	Distancia Maxima	Desviacion Estandar Distancia Hausdorff	Promedio Distancias	Afinidad Media
8102148			[0 , 240]				1.0396	0.0639	0.7624	0.8397
La Transaccion			[10 , 10]				1.0326	0	1.0326	0.8665







# PROMiDAT

IBEROAMERICANO

Programa Iberoamericano de  
Formación en Minería de Datos

Gracias....