



ESCALAMIENTO MULTIDIMENSIONAL (MDS)



Contenido

I. Introducción

1.1 Definición

1.2 Objetivo

1.3 Modelos de MDS

II. Escalamiento Multidimensional Clásico (cMDS)

2.1 Formalización

III. Escalamiento Multidimensional Métrico (mMDS)

3.1 Mínimos cuadrados ordinarios

IV. Escalamiento Multidimensional No Métrico (nMDS)



1.1 Definición

¿Qué es escalamiento multidimensional (MDS)?

MDS es una técnica que trabaja con proximidades entre persona, objetos o estímulos usados para producir una representación espacial de estos ítems.

La **proximidad** expresa la *similaridad* o *disimilaridad* entre los valores que representan a los objetos.



Su objetivo

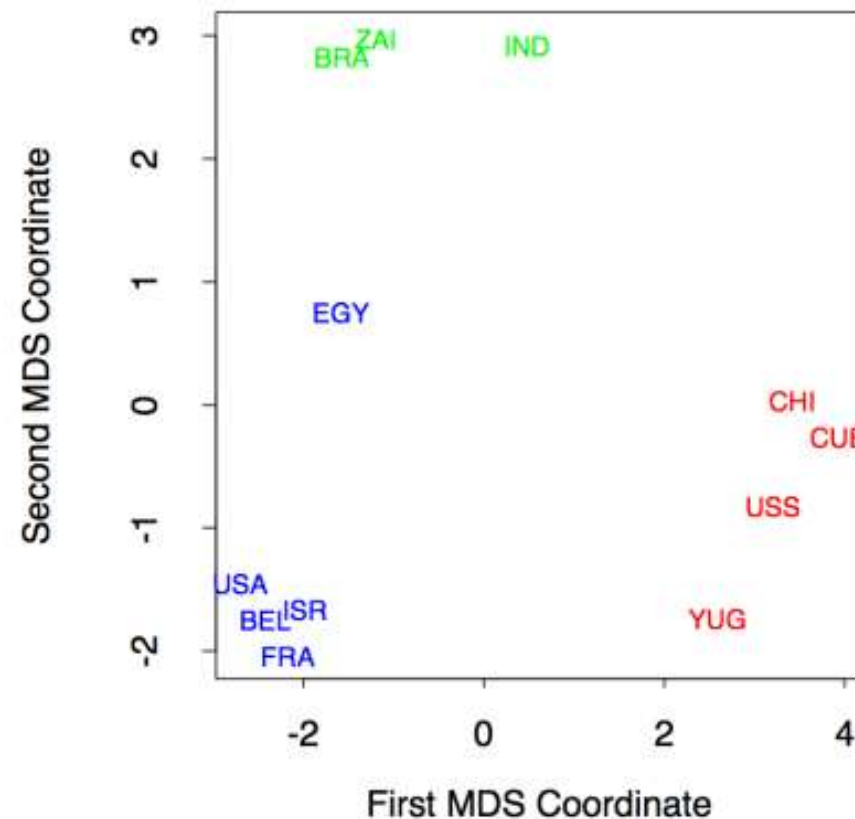
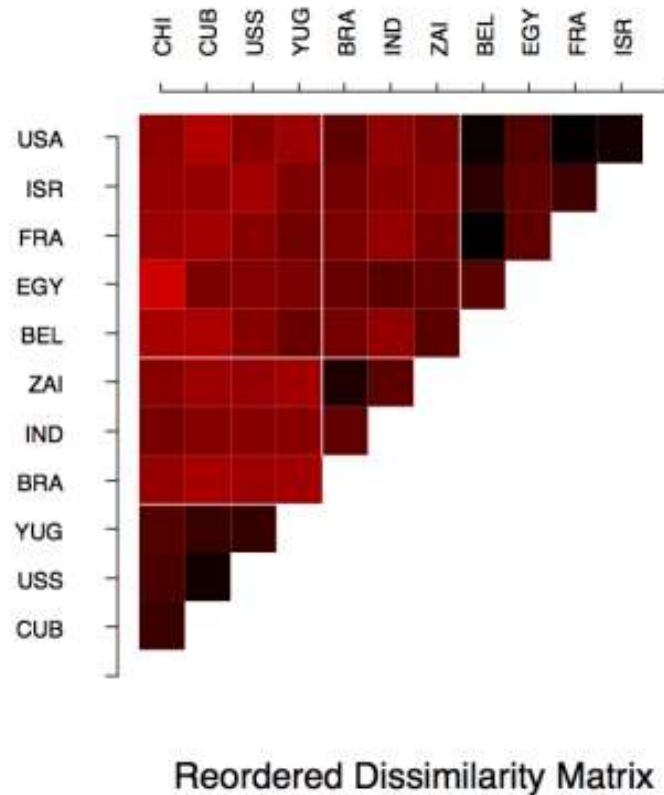
Es la reducción de la dimensionalidad porque su interés radica en encontrar un conjunto de puntos en un espacio de menor dimensión, por lo general euclidiano, que represente la configuración de los datos de una dimensión superior pero desconocida.

La configuración en alta dimensión es representada por la distancia (d) o matriz de disimilaridad

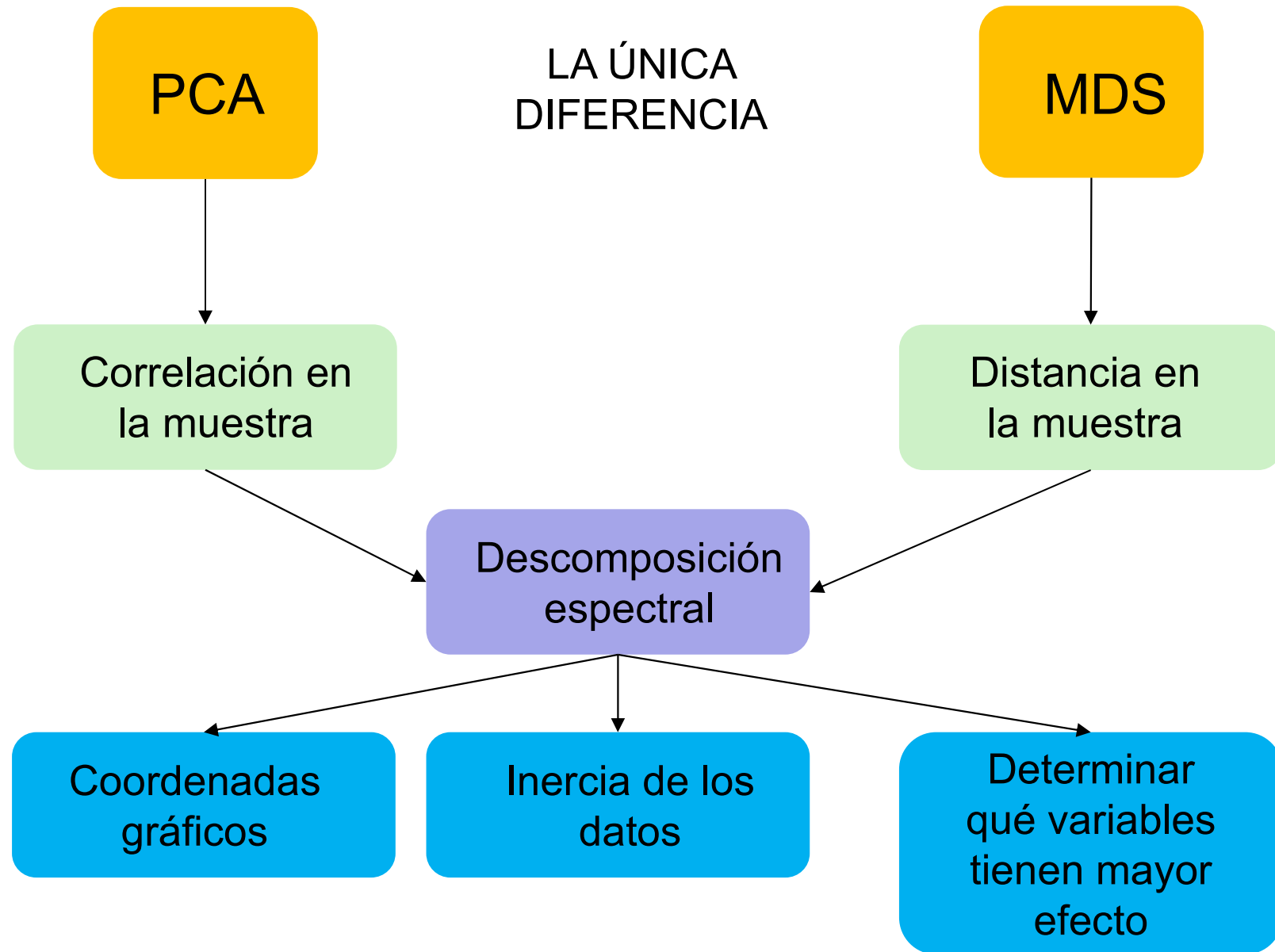
“Dado una disimilaridad (no necesariamente una métrica), reconstruir un mapa que preserve distancias”

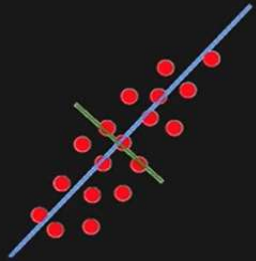


Su objetivo



“Dado una disimilaridad (no necesariamente una métrica), reconstruir un mapa que preserve distancias”



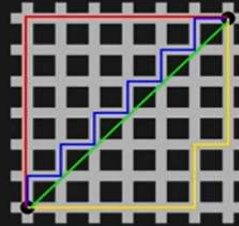


$$d_{ij} \approx \|x_i - x_j\|_2$$

Eigendecomposition

Classical MDS (PCoA)

Torgerson (1958) and Gower (1966)



$$d_{ij} \approx \delta_{ij}$$

Iterative Solution

Metric MDS (least squares MDS)

Shepard (1962)
Kruskal (1964)



$$d_{ij} \approx f(\delta_{ij})$$

Iterative Solution

Non-Metric MDS (least squares MDS)

Shepard (1962)
Kruskal (1964)

Principal Coordinates Analysis
Distancias euclideas

$$d_{r1,s1} < d_{r2,s2} < \dots < d_{rm,sm} \leftrightarrow f(d_{r1,s1}) < f(d_{r2,s2}) < \dots < f(d_{rm,sm})$$

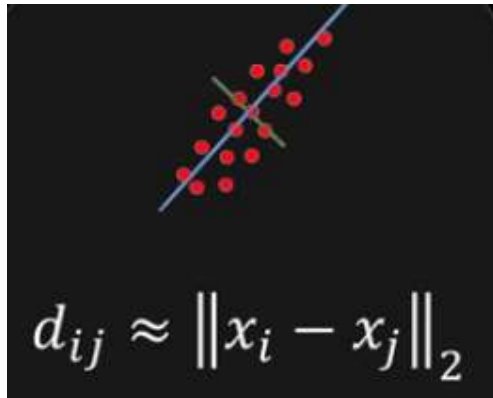


Definiciones básicas del MSD

δ_{ij} \longrightarrow Distancia observada / disimilaridad
(en el espacio original)

d_{ij} \longrightarrow Distancia en una configuración de menor
dimensión
(en el espacio reducido)

$\hat{d}_{ij} = f(\delta_{ij})$ \longrightarrow Disparidad



Implementación del MDS clásico (cMDS)

1. **Matriz B:** Calcular la matriz de doble centralidad o matriz producto interno
2. **Descomposición espectral:** Calcular los m valores estadísticos más grandes
3. **Proyectar:** Construir la matriz de coordenadas finales



1. **Matriz B:** Calcular la matriz de doble centralidad o matriz producto interno

B = HAH es la matriz de doble centralización $n \times n$ o matriz producto interno

Donde: **H** = $\mathbf{I} - n^{-1} \mathbf{11}^T$ es la matriz de centralización $n \times n$

y **A** = $\begin{pmatrix} a_{rs} \end{pmatrix}$ con $a_{rs} = -\frac{1}{2} d_{rs}^2$,

$$d_{rs}^2 = (\mathbf{x}_r - \mathbf{x}_s)^T (\mathbf{x}_r - \mathbf{x}_s).$$



2. Descomposición espectral: Calcular los m valores característicos más grandes

$$B = V \Lambda V'$$



3. **Proyección:** Construir la matriz de coordenadas finales (en un espacio de menor dimensión “m”)

$$X_{(m)} = V_m \sqrt{\Lambda_m}$$

$X_{(m)}$: matriz de puntos (coordenadas)

V_m : matriz construida con las m primeros vectores característicos de V

Λ_m : matriz m-dimensional de valores característicos asociados a V



Definición (Mardia, pp.397)

Una matriz de distancia **D** se le dice euclideana si existe una configuración de puntos en algún espacio euclideo cuya distancia interpuntos están dadas por **D**; esto es, si para algún **p**, existen puntos $x_1, \dots, x_n \in R^p$ tal que:

$$d_{rs}^2 = (\mathbf{x}_r - \mathbf{x}_s)^T (\mathbf{x}_r - \mathbf{x}_s).$$

Teorema 14.2.1 (Mardia)

Dado una matriz de distancia \mathbf{D} y de doble centralización \mathbf{B} , entonces, \mathbf{D} es euclidea si y solo si $\mathbf{B} \geq 0$. En particular, los siguientes resultados se confirman:

Nota.-

Si \mathbf{D} es la matriz de distancias euclideas interpuntos para una configuración $\mathbf{z} = (x_1, x_2, \dots, x_n)$, entonces,

$$b_{rs} = (x_r - \bar{x})(x_s - \bar{x})$$

En notación matricial $\mathbf{B} = (\mathbf{H}\mathbf{X})(\mathbf{H}\mathbf{X})^t$ tal que $\mathbf{B} \geq 0$.



Prueba:

Los elementos de la matriz $\mathbf{D} = d_{rs}$ son distancias euclideas, si y solo si, la matriz de doble centralidad $\mathbf{B} = \mathbf{H}\mathbf{A}\mathbf{H}$ es definida no negativa.

⇒) Hipótesis d_{rs}^2 es una distancia euclidea

$$\mathbf{A} = \left[a_{rs} = -\frac{1}{2}d_{rs}^2 \right] \quad \text{y} \quad \mathbf{B} = \mathbf{H}\mathbf{A}\mathbf{H}$$

Tesis \mathbf{B} es semidefinida positiva

⇒ Como $\mathbf{B} = \mathbf{H}\mathbf{A}\mathbf{H} = (\mathbf{I}_n - 1/n \mathbf{J}_n)\mathbf{A}(\mathbf{I}_n - 1/n \mathbf{J}_n)$

$$b_{rs} = a_{rs} - \bar{a}_{.s} - \bar{a}_{r.} + \bar{a}_{..} \quad \text{..... (i)}$$

$$\text{Pero p.h. } a_{rs} = -\frac{1}{2}d_{rs}^2 = -1/2 (\mathbf{x}_r - \mathbf{x}_s)^t (\mathbf{x}_r - \mathbf{x}_s)$$

$$= - (1/2)\mathbf{x}_r^t \mathbf{x}_r - (1/2)\mathbf{x}_s^t \mathbf{x}_s + \mathbf{x}_r^t \mathbf{x}_s \quad \text{..... en (i)}$$



$$\Rightarrow b_{rs} = (x_r - \bar{x})(x_s - \bar{x})$$

$$\Rightarrow \mathbf{B} = \begin{bmatrix} b_{rs} \end{bmatrix} = (\mathbf{Hx})(\mathbf{Hx})^t = \mathbf{ZZ}^t \geq 0, \forall \mathbf{Z}$$

\Leftarrow) Hipótesis \mathbf{B} es semidefinida positiva

$$\mathbf{A} = \left[a_{rs} = -\frac{1}{2}d_{rs}^2 \right] \quad \text{y} \quad \mathbf{B} = \mathbf{H}\mathbf{A}\mathbf{H}$$

Tesis d_{rs}^2 es una distancia euclidea

\Rightarrow a) simetría

$$\text{p.h. } \mathbf{B} = \mathbf{ZZ}^t \quad \Rightarrow \quad \mathbf{B}^t = \mathbf{ZZ}^t = \mathbf{B} \quad \text{luego, } \mathbf{B} \text{ es simétrica}$$



p.h. $\mathbf{B} = \mathbf{H}\mathbf{A}\mathbf{H}$, entonces, \mathbf{A} es simétrica.

$$\text{Además, } a_{rs} = -\frac{1}{2}d_{rs}^2, \quad \Rightarrow \quad d_{rs}^2 = -2a_{rs} = -2a_{sr} = d_{sr}^2$$

b) no negatividad

$$d_{rs}^2 = \mathbf{d}^t \mathbf{d} > 0 \quad \text{solo si } d \neq 0. \text{ En particular, } d = x_r - x_s \neq 0$$



$$x_r \neq x_s$$

c) identidad

$$d_{rs}^2 = 0 \quad \text{pero } d_{rs}^2 = \mathbf{d}^t \mathbf{d} \text{ es igual a cero solo si } \mathbf{d} = 0$$



$$x_r = x_s$$

De a), b) y c) concluimos que d_{rs}^2 es una distancia euclidea (...)



Resultado a) $\mathbf{b}_{rs} = (\mathbf{x}_r - \bar{\mathbf{x}})^t (\mathbf{x}_s - \bar{\mathbf{x}})$

Aquí asumimos que
la configuración es \mathbf{X}

$$\Rightarrow \mathbf{B} = (\mathbf{X} - \mathbf{1}_n \bar{\mathbf{x}})^t (\mathbf{X} - \mathbf{1}_n \bar{\mathbf{x}})$$

$$\mathbf{B} = (\mathbf{H}\mathbf{X})(\mathbf{H}\mathbf{X})^T = \mathbf{W}\mathbf{W}^t \geq 0$$

Resultado b) $\mathbf{\Lambda}_1 = \text{diag}(\lambda_1, \dots, \lambda_p), \quad \mathbf{V}_1 = [\mathbf{v}_1, \dots, \mathbf{v}_p]$ dado que: $r(\mathbf{B})=p$

$$\Rightarrow \mathbf{\Lambda}_1^{\frac{1}{2}} = \text{diag}(\lambda_1^{\frac{1}{2}}, \dots, \lambda_p^{\frac{1}{2}})$$

$$\Rightarrow \mathbf{X} = \mathbf{V}_1 \mathbf{\Lambda}_1^{\frac{1}{2}}$$

$$\mathbf{B} = \mathbf{V}_1 \mathbf{\Lambda}_1 \mathbf{V}_1^T \text{ es la descomposición espectral de } \mathbf{B}$$

$$\Rightarrow \mathbf{B} = \mathbf{X}\mathbf{X}^T$$

O también: $[\mathbf{B}]_{rs} = b_{rs} = \mathbf{x}_r^T \mathbf{x}_s.$

Matriz producto
interno de la
configuración \mathbf{X}



Se sabe que \mathbf{d} es una distancia euclídea:

$$d_{rs}^2 = \mathbf{x}_r^T \mathbf{x}_r + \mathbf{x}_s^T \mathbf{x}_s - 2\mathbf{x}_r^T \mathbf{x}_s.$$

Aplicando sumatorias:

$$\begin{aligned} \Rightarrow \quad \frac{1}{n} \sum_{r=1}^n d_{rs}^2 &= \frac{1}{n} \sum_{r=1}^n \mathbf{x}_r^T \mathbf{x}_r + \mathbf{x}_s^T \mathbf{x}_s, \\ \frac{1}{n} \sum_{s=1}^n d_{rs}^2 &= \mathbf{x}_r^T \mathbf{x}_r + \frac{1}{n} \sum_{s=1}^n \mathbf{x}_s^T \mathbf{x}_s, \\ \frac{1}{n^2} \sum_{r=1}^n \sum_{s=1}^n d_{rs}^2 &= \frac{2}{n} \sum_{r=1}^n \mathbf{x}_r^T \mathbf{x}_r. \end{aligned}$$

Asumiendo que:

$$\sum_{r=1}^n x_{ri} = 0 \quad (i = 1, \dots, p)$$

Porque la configuración
X está centrada en **0**



Se sabe que:

$$\begin{aligned}b_{rs} &= \mathbf{x}_r^T \mathbf{x}_s, \\&= -\frac{1}{2} \left(d_{rs}^2 - \frac{1}{n} \sum_{r=1}^n d_{rs}^2 - \frac{1}{n} \sum_{s=1}^n d_{rs}^2 + \frac{1}{n^2} \sum_{r=1}^n \sum_{s=1}^n d_{rs}^2 \right) \\&= a_{rs} - \bar{a}_{r.} - \bar{a}_{.s} + \bar{a}_{..},\end{aligned}$$

Donde:

$$\begin{aligned}a_{rs} &= -\frac{1}{2} d_{rs}^2, \\ \bar{a}_{r.} &= n^{-1} \sum_s a_{rs}, \quad \bar{a}_{.s} = n^{-1} \sum_r a_{rs}, \quad \bar{a}_{..} = n^{-2} \sum_r \sum_s a_{rs}\end{aligned}$$



Expresando **B** en términos de **A**:

$$\begin{aligned}(\mathbf{x}_r - \mathbf{x}_s)^T (\mathbf{x}_r - \mathbf{x}_s) &= \mathbf{x}_r^T \mathbf{x}_r + \mathbf{x}_s^T \mathbf{x}_s - 2\mathbf{x}_r^T \mathbf{x}_s \\&= b_{rr} + b_{ss} - 2b_{rs} \\&= a_{rr} + a_{ss} - 2a_{rs} \\&= -2a_{rs} = \delta_{rs}^2,\end{aligned}$$

Dado que : $\mathbf{B}\mathbf{1} = \mathbf{H}\mathbf{A}\mathbf{H}\mathbf{1} = \mathbf{0}$, entonces, **1** es un vector característico cuyo valor característico asociado es **0**.



El vector **1** es ortogonal
a las columnas de **X**

dado que $\mathbf{B} = (\mathbf{H}\mathbf{X})(\mathbf{H}\mathbf{X})^T$ (...)

EJEMPLO.- Δ
Disparidad (d) entre 5 objetos

	O_1	O_2	O_3	O_4	O_5
O_1	0.00	2.37	1.05	2.34	2.34
O_2	2.37	0.00	2.34	2.81	2.81
O_3	1.05	2.34	0.00	2.42	2.42
O_4	3.54	3.63	3.50	0.00	4.24
O_5	2.34	2.81	2.42	4.24	0.00

- ✓ No se conocen los datos originales
- ✓ No se sabe cómo se determinaron las disparidades (δ_{ij})
- ✓ Aplicamos el teorema 14.2.1 (calculamos **B**)
- ✓ Si son datos métricos, entonces,
- ✓ Aplicamos algún procedimiento MDS



Para calcular la matriz **B** hacemos $\delta^2_{ij} \forall i,j$ y se obtiene la siguiente matriz:

$$\begin{pmatrix} 0.00 & 5.62 & 1.11 & 12.59 & 5.48 \\ 5.62 & 0.00 & 5.50 & 13.20 & 7.94 \\ 1.10 & 5.50 & 0.00 & 12.31 & 5.87 \\ 12.59 & 13.20 & 12.31 & 0.00 & 18.03 \\ 5.48 & 7.94 & 5.87 & 18.03 & 0.00 \end{pmatrix}$$

Dado que $a_{ij} = -\frac{1}{2} \delta^2_{ij}$

$$A = -\frac{1}{2} \begin{pmatrix} 0.00 & 5.62 & 1.11 & 12.59 & 5.48 \\ 5.62 & 0.00 & 5.50 & 13.20 & 7.94 \\ 1.10 & 5.50 & 0.00 & 12.31 & 5.87 \\ 12.59 & 13.20 & 12.31 & 0.00 & 18.03 \\ 5.48 & 7.94 & 5.87 & 18.03 & 0.00 \end{pmatrix}$$

Y dado que: $b_{ij} = a_{ij} - \bar{a}_{i.} - \bar{a}_{.j} + \bar{a}_{..}$

Obtenemos:

$$\mathbf{B} = \begin{pmatrix} 1.454 & -0.610 & 0.898 & -1.708 & -0.034 \\ -0.610 & 2.946 & -0.551 & -1.267 & -0.518 \\ 0.898 & -0.551 & 1.452 & -1.569 & -0.230 \\ -1.708 & -1.267 & -1.569 & 7.720 & -3.176 \\ -0.034 & -0.518 & -0.230 & -3.176 & 3.958 \end{pmatrix}$$

Cuyos valores característicos son:

10.0895 3.6366 3.2546 0.5493 0.0000

Entonces, $\mathbf{B} \geq 0$, lo que significa que la matriz original de disimilaridades es euclideana



Datos

$$X = \begin{pmatrix} 3 & 4 & 4 & 6 & 1 \\ 5 & 1 & 1 & 7 & 3 \\ 6 & 2 & 0 & 2 & 6 \\ 1 & 1 & 1 & 0 & 3 \\ 4 & 7 & 3 & 6 & 2 \\ 2 & 2 & 5 & 1 & 0 \\ 0 & 4 & 1 & 1 & 1 \\ 0 & 6 & 4 & 3 & 5 \\ 7 & 6 & 5 & 1 & 4 \\ 2 & 1 & 4 & 3 & 1 \end{pmatrix}$$

```
# Cargar datos
> X = read.delim("clipboard")
> dim(X)
[1] 10 5
```

```
> D <- dist(X)
> D
```

	1	2	3	4	5	6	7	8	9
2	5.196152								
3	8.366600	6.082763							
4	7.874008	8.062258	6.324555						
5	3.464102	6.557439	8.366600	9.273618					
6	5.656854	8.426150	8.831761	5.291503	7.874008				
7	6.557439	8.602325	8.185353	3.872983	7.416198	5.000000			
8	6.164414	8.888194	8.366600	6.928203	6.000000	7.071068	5.744563		
9	7.416198	9.055385	6.855655	8.888194	6.557439	7.549834	8.831761	7.416198	
10	4.358899	6.164414	7.681146	4.795832	7.141428	2.645751	5.099020	6.708204	8.000000



cmdscale()

cmdscale(dist(X), k=5)

cmdscale(D,k=5)

k = 2, 3

```
> cmdscale(dist(X),k=5)
      [,1]      [,2]      [,3]      [,4]
[1,] -1.6038325  2.38060903 -2.2301092 -0.3656856
[2,] -2.8246377 -2.30937202 -3.9523782  0.3419185
[3,] -1.6908272 -5.13970089  1.2880306  0.6503227
[4,]  3.9527719 -2.43233961  0.3833746  0.6863995
[5,] -3.5984894  2.75538195 -0.2551393  1.0783741
[6,]  2.9520356  1.35475175 -0.1899027 -2.8211220
[7,]  3.4689928  0.76411068  0.3016531  1.6369166
[8,]  0.3545235  2.31408566  2.2161772  2.9240116
[9,] -2.9362323 -0.01279597  4.3117385 -2.5122743
[10,]  1.9256952  0.32526941 -1.8734445 -1.6188611
      [,5]
[1,]  0.11536476
[2,]  0.33169405
[3,] -0.05133897
[4,] -0.03460933
[5,] -1.26125237
[6,]  0.12385813
[7,] -1.94209512
[8,]  2.00450379
[9,] -0.18911558
[10,]  0.90299062
```



minimizar $\phi = \sum_{r,s}^n (d_{rs}^2 - \hat{d}_{rs}^2)$

distancia MDS

distancia euclídea

```
> dist(X)-dist(cmdscale(dist(X),k=5))
```

O también

```
> D-dist(cmdscale(D,k=5))
```

donde **D** es la matriz
de distancias
euclídeas

⇒ $\Phi \approx 0$

Comprobación haciendo $r=1$ y $s=2$ en la configuración de MDS \hat{X} :

	\hat{O}_1	1.6	2.38	2.23	-0.37	0.12
	\hat{O}_2	2.82	-2.31	3.95	0.34	0.33
		↓	↓	↓	↓	↓
$\hat{x}_{1j} - \hat{x}_{2j}$	→	-1.22	4.69	-1.72	-0.71	-0.21
$(\hat{x}_{1j} - \hat{x}_{2j})^2$	→	1.4884	21.9961	2.9584	0.5041	0.0441

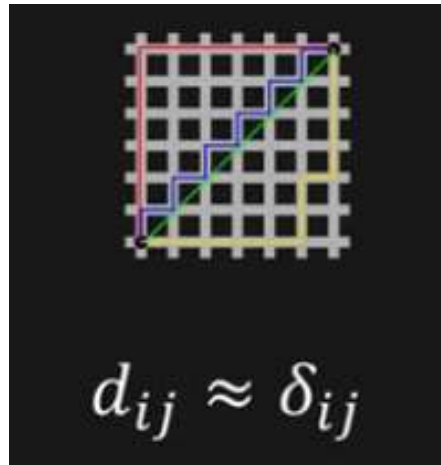


$$\hat{d}_{12}^2 = 5.19529595$$

distancia euclídea entre O_1 y O_2 en el espacio de menor dimensión

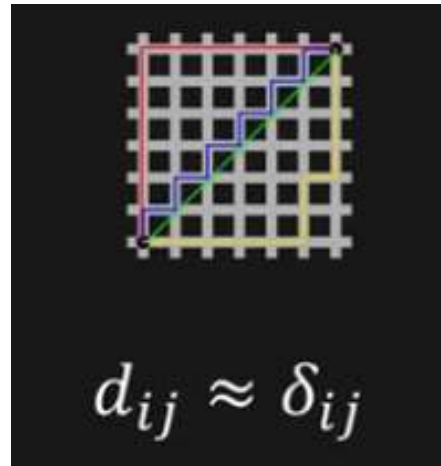
$$d_{12}^2 = 5.196152$$

distancia euclídea real entre O_1 y O_2



Implementación del MDS métrico (mMDS)

Si existen n objetos con disimilaridades $\{\delta_{ij}\}$ el MDS métrico ($mMDS$) intenta encontrar un conjunto de puntos en un espacio donde cada punto representa a un objeto y las distancias entre puntos $\{d_{rs}\}$ son tal que: $d_{rs} \approx f(\delta_{rs})$ donde f es una función paramétrica monótonica continua.

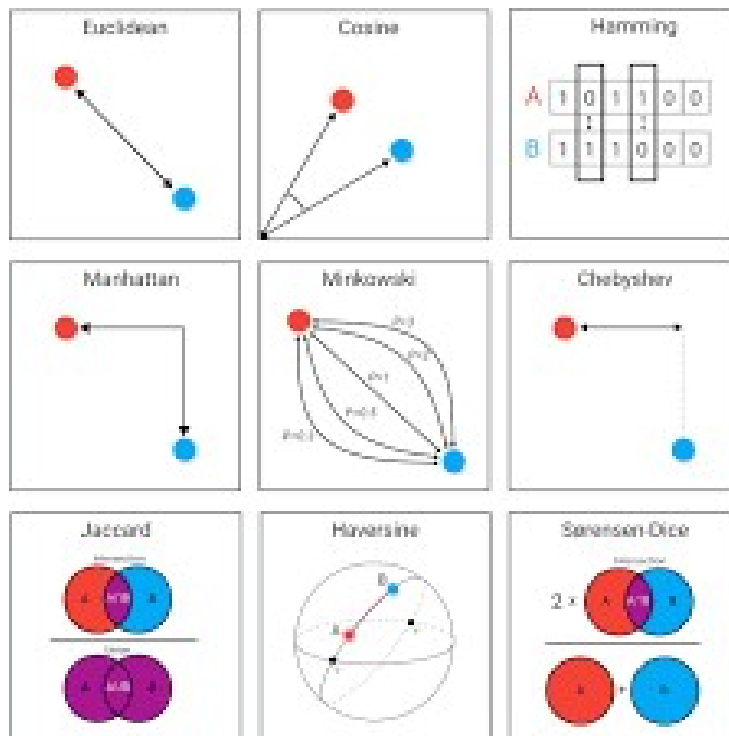


Implementación del MSD
métrico
(mMDS)

1. **Matriz de distancias:** La d_{ij} es una métrica (escala de proporción)
2. **Solución iterativa:** inicialización, cálculo de distancias, calculo de la perdida y optimización para minimizar *stress*
3. **Pérdida:** Standardized Residual Sum of Squares (*STRESS*)



1. Matriz de distancias: Tratamiento más general de distancias. La métrica.



Positivity

$$d_{ij} \geq 0, \text{ for } i \neq j$$
$$d_{ii} = 0$$

Symmetry

$$d_{ij} = d_{ji}$$

Triangle inequality
→ Metric

$$d_{ik} \leq d_{ij} + d_{jk}$$





Medidas de disimilaridad para datos cuantitativos

- ① $d(x, y) \geq 0$,
- ② $d(x, y) = 0$ if and only if $x = y$,
- ③ $d(x, y) = d(y, x)$,
- ④ $d(x, z) \leq d(x, y) + d(y, z)$.

Euclidean distance

$$\delta_{rs} = \left\{ \sum_i (x_{ri} - x_{si})^2 \right\}^{\frac{1}{2}}$$

Weighted Euclidean

$$\delta_{rs} = \left\{ \sum_i w_i (x_{ri} - x_{si})^2 \right\}^{\frac{1}{2}}$$

Mahalanobis distance

$$\delta_{rs} = \{(\mathbf{x}_r - \mathbf{x}_s)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_r - \mathbf{x}_s)\}^{\frac{1}{2}}$$

City block metric

$$\delta_{rs} = \sum_i |x_{ri} - x_{si}|$$

Minkowski metric

$$\delta_{rs} = \left\{ \sum_i w_i |x_{ri} - x_{si}|^\lambda \right\}^{\frac{1}{\lambda}} \quad \lambda \geq 1$$

Canberra metric

$$\delta_{rs} = \sum_i |x_{ri} - x_{si}| / (x_{ri} + x_{si})$$

Divergence

$$\delta_{rs} = \frac{1}{p} \sum_i (x_{ri} - x_{si})^2 / (x_{ri} + x_{si})^2$$

Bray-Curtis

$$\delta_{rs} = \frac{1}{p} \frac{\sum_i |x_{ri} - x_{si}|}{\sum_i (x_{ri} + x_{si})}$$

Soergel

$$\delta_{rs} = \frac{\sum_i |x_{ri} - x_{si}|}{\sum_i \max(x_{ri}, x_{si})}$$

Bhattacharyya distance

$$\delta_{rs} = \left\{ \sum_i (x_{ri}^{\frac{1}{2}} - x_{si}^{\frac{1}{2}})^2 \right\}^{\frac{1}{2}}$$

Wave-Hedges

$$\delta_{rs} = \frac{1}{p} \sum_i \left(1 - \frac{\min(x_{ri}, x_{si})}{\max(x_{ri}, x_{si})} \right)$$

Angular separation

$$\delta_{rs} = 1 - \frac{\sum_i x_{ri} x_{si}}{[\sum_i x_{ri}^2 \sum_i x_{si}^2]^{\frac{1}{2}}}$$

Correlation

$$\delta_{rs} = 1 - \frac{\sum_i (x_{ri} - \bar{x}_r)(x_{si} - \bar{x}_s)}{\left\{ \sum_i (x_{ri} - \bar{x}_r)^2 \sum_i (x_{si} - \bar{x}_s)^2 \right\}^{\frac{1}{2}}}$$

Datos binarios

	Object <i>s</i>		
	1	0	
Object <i>r</i>	1	$\begin{array}{ c c } \hline a & b \\ \hline \end{array}$	$a + b$
	0	$\begin{array}{ c c } \hline c & d \\ \hline \end{array}$	$c + d$
	$a + c$	$b + d$	$p = a + b + c + d$

Coeficientes de
similitud para
datos binarios

Braun, Blanque

$$s_{rs} = \frac{a}{\max\{(a+b), (a+c)\}}$$

Czekanowski, Sørensen, Dice

$$s_{rs} = \frac{2a}{2a + b + c}$$

Hamman

$$s_{rs} = \frac{a - (b + c) + d}{a + b + c + d}$$

Jaccard coefficient

$$s_{rs} = \frac{a}{a + b + c}$$

Kulczynski

$$s_{rs} = \frac{a}{b + c}$$

Kulczynski

$$s_{rs} = \frac{1}{2} \left(\frac{a}{a+b} + \frac{a}{a+c} \right)$$

Michael

$$s_{rs} = \frac{4(ad - bc)}{\{(a+d)^2 + (b+c)^2\}}$$

Mountford

$$s_{rs} = \frac{2a}{a(b+c) + 2bc}$$

Mozley, Margalef

$$s_{rs} = \frac{a(a+b+c+d)}{(a+b)(a+c)}$$

Ochiai

$$s_{rs} = \frac{a}{[(a+b)(a+c)]^{\frac{1}{2}}}$$

Phi

$$s_{rs} = \frac{ad - bc}{[(a+b)(a+c)(b+d)(c+d)]^{\frac{1}{2}}}$$

Rogers, Tanimoto

$$s_{rs} = \frac{a+d}{a+2b+2c+d}$$

Russell, Rao

$$s_{rs} = \frac{a}{a+b+c+d}$$

Simple matching coefficient

$$s_{rs} = \frac{a+d}{a+b+c+d}$$

Simpson

$$s_{rs} = \frac{a}{\min\{(a+b), (a+c)\}}$$

Sokal, Sneath, Anderberg

$$s_{rs} = \frac{a}{a+2(b+c)}$$

Yule

$$s_{rs} = \frac{ad - bc}{ad + bc}$$

2. **Solución iterativa:** familia de algoritmos diseñados para alcanzar una configuración óptima de baja dimensión.

Pasos

1. **Inicialización:** Se inicializa con puntos en posiciones aleatorias
2. **Cálculo de distancias:** Se obtiene la matriz de distancias para la configuración
3. **Cálculo de la pérdida (*loss*):** Evaluar la función *stress*
4. **Optimizar:** “descenso por gradiente” para actualizar el *stress* minimizado

2. Cálculo de distancias: Se obtiene la matriz de distancias para la configuración

Una matriz $\mathbf{D}(n \times n)$ se llama matriz distancia si es simétrica y

$$d_{rr} = 0, \quad d_{rs} > 0, \quad \text{si } r \neq s$$

tal que si d_{rs} denota una distancia euclídea entre \mathbf{P}_r y \mathbf{P}_s ,



3. Cálculo de la pérdida (*loss*): Evaluar la función *stress*

Standardized Residual Sum of Squares (*STRESS*)

Dado una dimensión p y una función monótona $f()$, entonces, $mMDS$ trata de encontrar una configuración óptima $X \subset \mathbb{R}^p$ tal que:

$$f(d_{ij}) \approx \hat{d}_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$$

y

$$f(d_{ij}) = Ad_{ij} + b$$

De esta forma, $\mathbf{P}_1, \dots, \mathbf{P}_n$ con coordenadas $\mathbf{x}_i' = (x_{i1}, \dots, x_{ip})$, $i=1, \dots, n$ representan una solución $mMDS$ p -dimensional



Es óptima en la medida que explícitamente se cumple:

$$\text{stress} = \mathcal{L}(\hat{d}_{ij}) = \left(\sum_{i < j} (\hat{d}_{ij} - f(d_{ij}))^2 / \sum d_{ij}^2 \right)^{\frac{1}{2}}$$

y la métrica MDS minimiza $\mathcal{L}(\hat{d}_{ij})$ sobre todos los \hat{d}_{ij} y \mathbf{A}, \mathbf{b} .

La métrica usual del *mMDS* es el caso especial donde $f(d_{ij}) = d_{ij}$

“La solución usual *mMDS* (por optimización) \neq a la solución *cMDS*”



Diagrama de Shepard

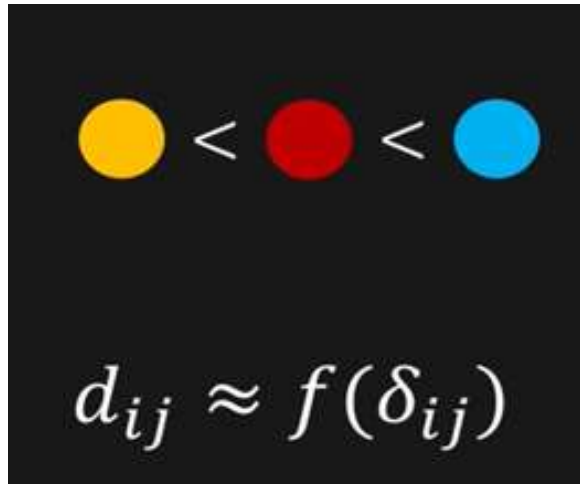




4. **Optimizar:** “descenso por gradiente” o para actualizar el *stress* minimizado

Gradient Descent

SMACOF
Scaling by Majorizing a
Complicated Fuction

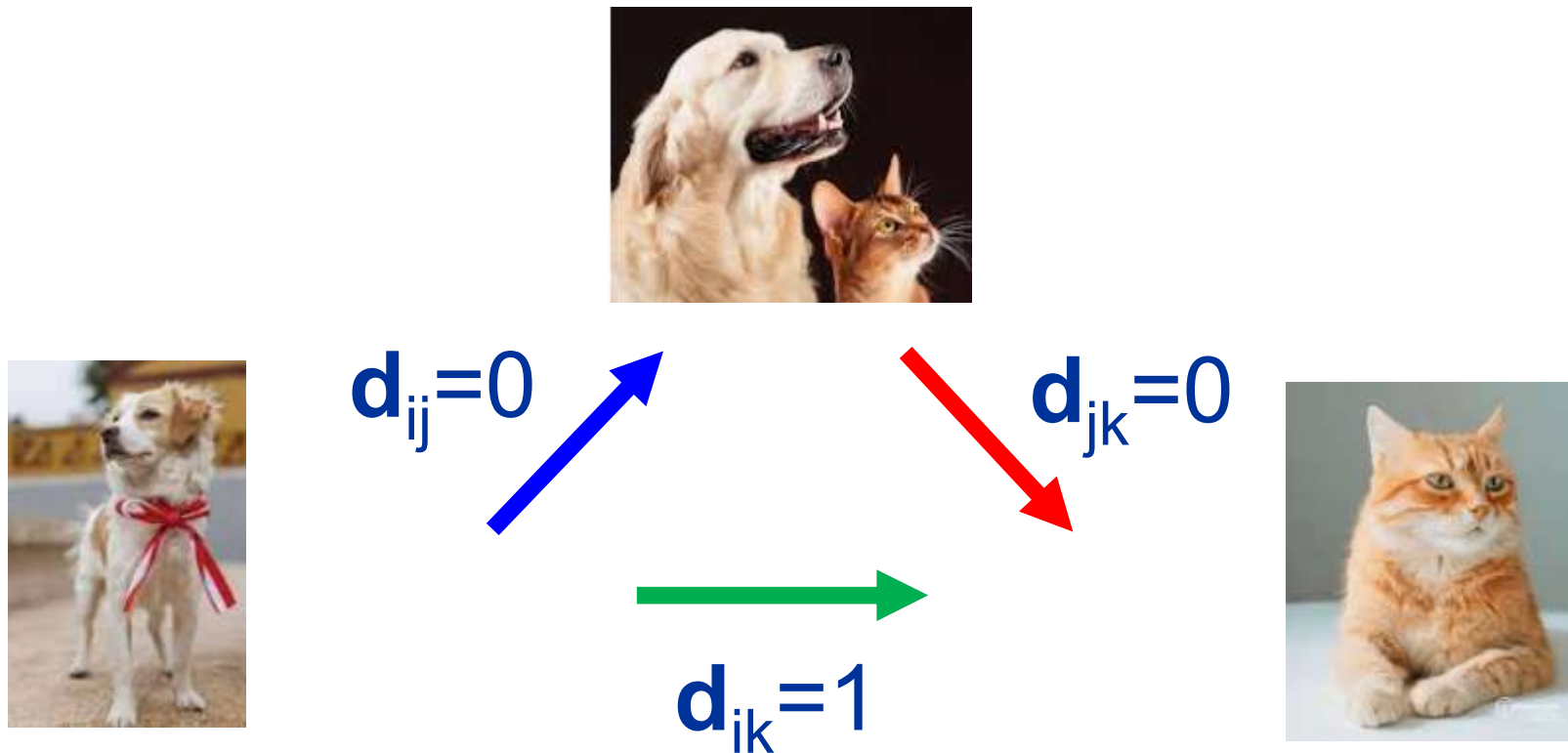


Implementación del MSD
no métrico
(nMDS)

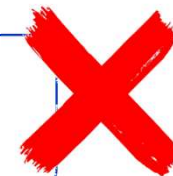
*Ordinal – función monotónica – disparidad – preserva
distancia – regresión isotónica – algoritmo iterativo*



No métrico

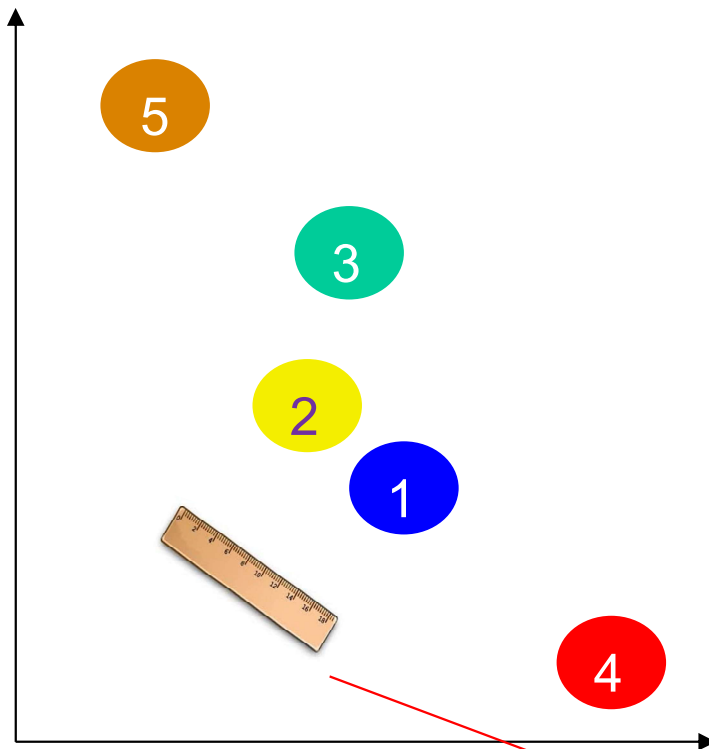


$$d_{ik} \leq d_{ij} + d_{jk}$$





Ordinal MDS

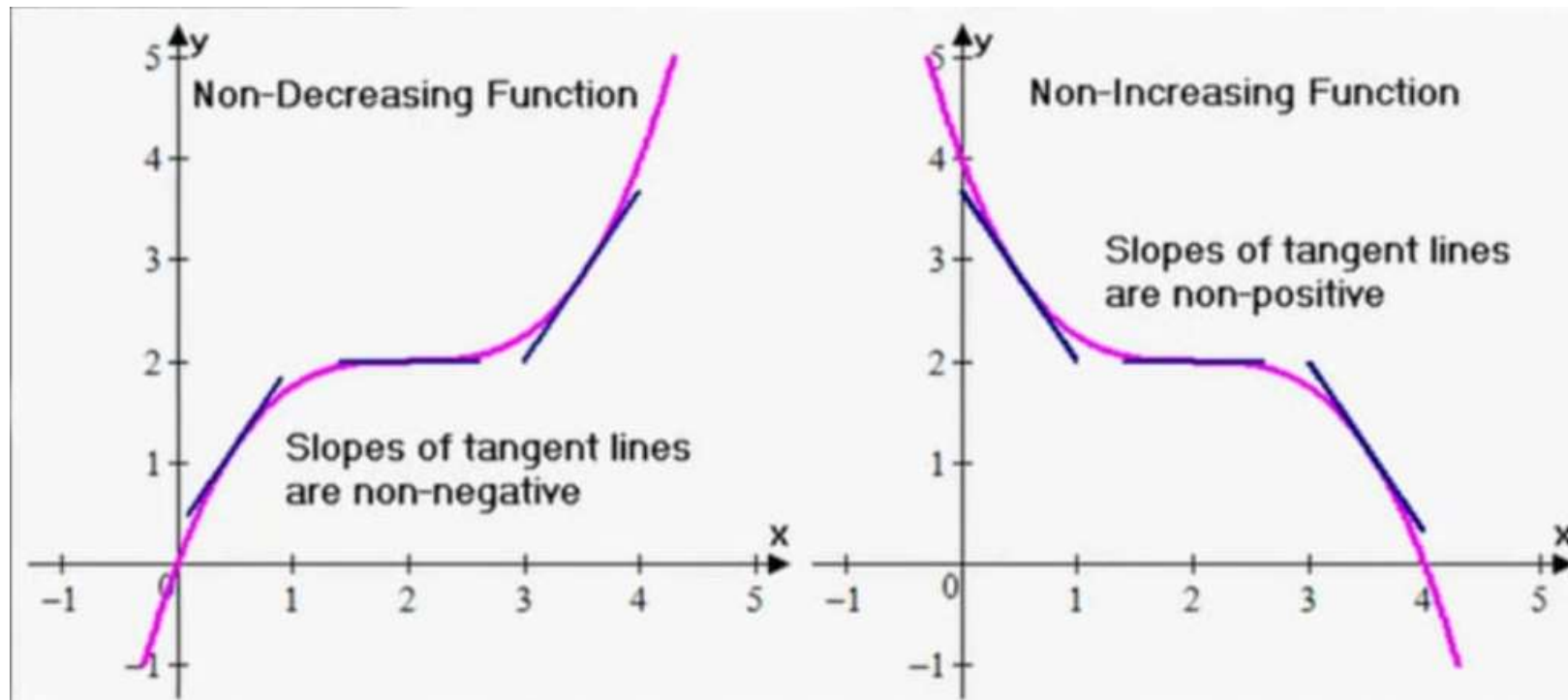


	695	760	6264	10160
		545	5563	10538
			5567	9997
				7420

$$\delta_{12} < \delta_{23} < \delta_{13} < \dots$$

¿Cómo aseguramos que en el espacio de menos dimensión se preservan las distancias reales?

Función monotónica creciente
o decreciente





$$d_{ij} \approx f(\delta_{12}) = \hat{d}_{ij}$$

f: función monotónica que
mapea las disparidades

\hat{d} : disparidades

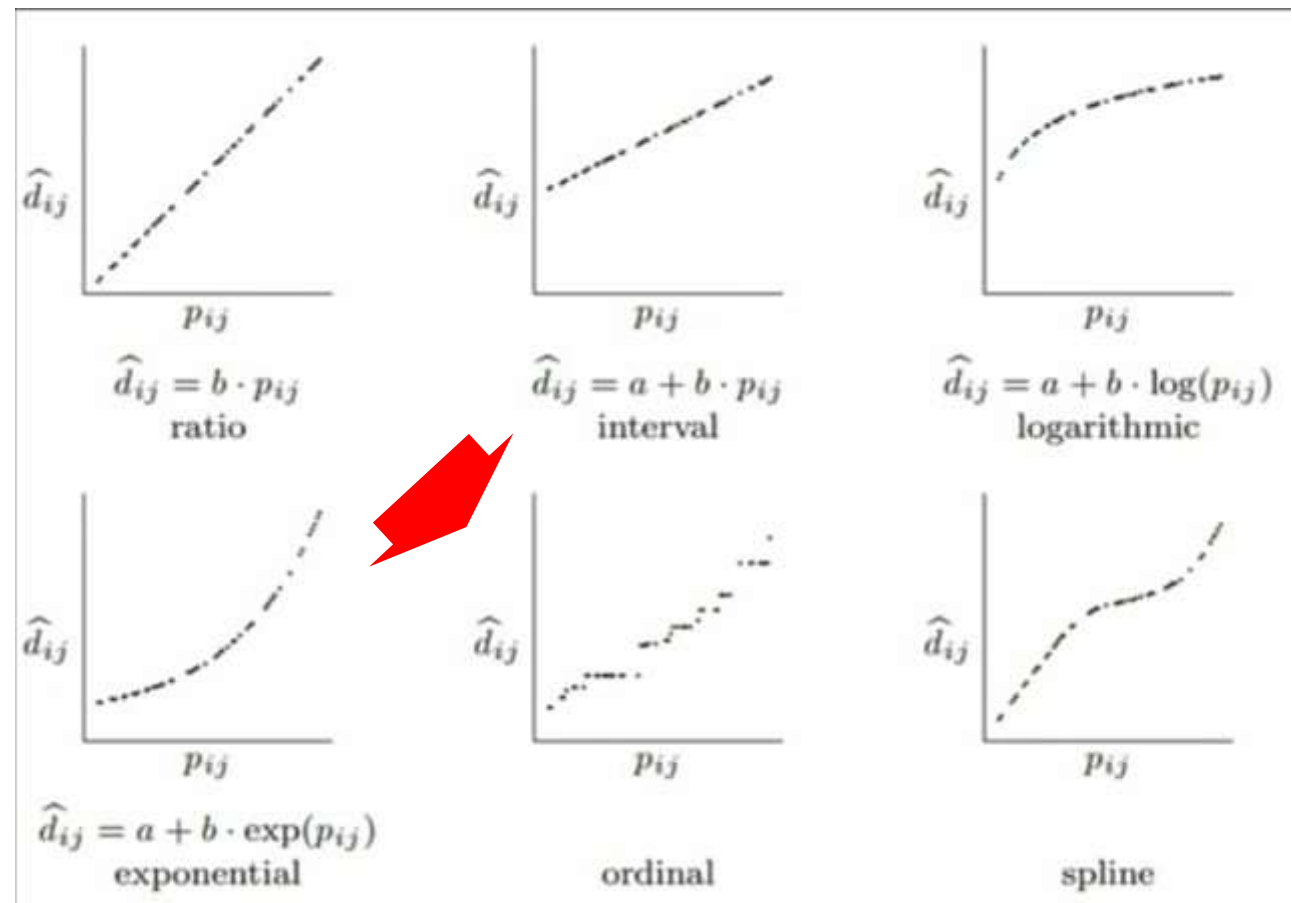
$$\delta_{ij} < \delta_{kl}$$

Preserva
el orden



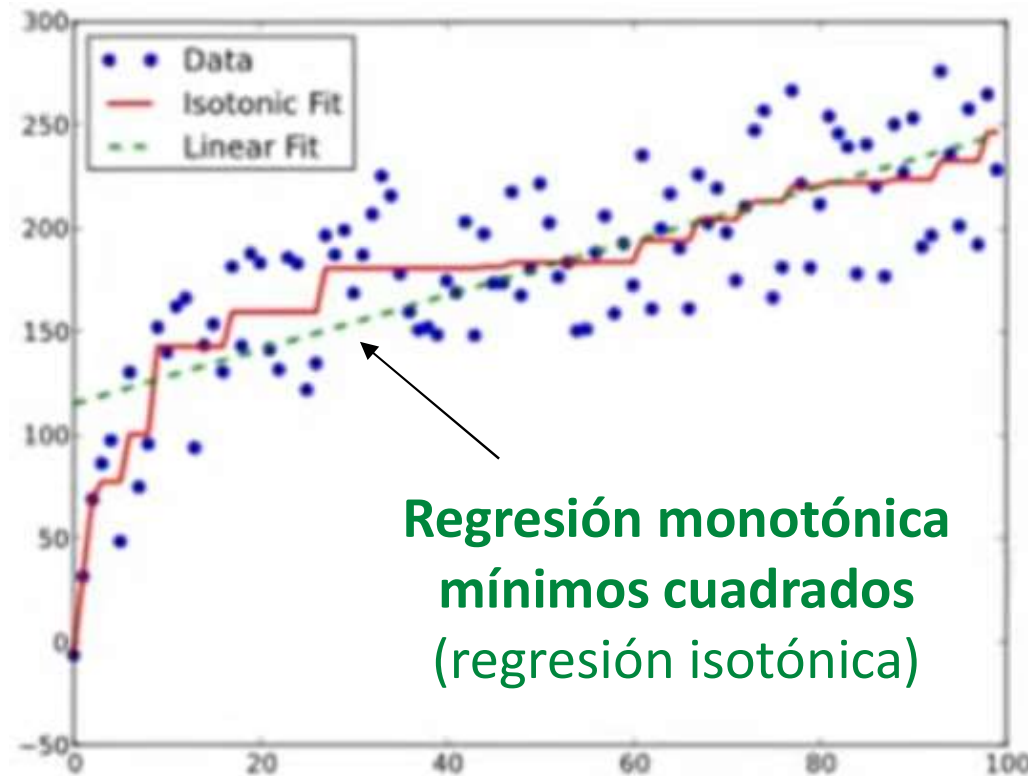
$$f(\delta_{ij}) < f(\delta_{kl})$$

Funciones monotónicas





Función monotónica



Algoritmo iterativo:
algoritmo de infractores
adyacentes al grupo

**Regresión monotónica
mínimos cuadrados
(regresión isotónica)**

Ajuste isotónico / ajuste lineal



Predicho por la
regresión
monotónica

$$Stress = \left(\frac{\sum_{i < j} (d_{ij} - \hat{d}_{ij})^2}{\sum_{i < j} d_{ij}^2} \right)^{1/2}$$

Caso: disimilaridades ordinales

A un panel de personas que tienen cierto conocimiento sobre la distribución física de las principales ciudades de una región se les nombra un par de ciudades y deben indicar si este par de ciudades ¿que tan cerca o lejos están. Para ello, se les proporciona una escala de calificación del 1 (muy cercanos) al 9 (muy lejanos).

Disimilaridades entre ciudades

La distancia definida no cumple las propiedades de distancia (requiere un tratamiento ordinal)

0	1	5	2	7	2	2	7	8	1	ATLANTA
1	0	3	3	7	5	2	7	6	1	CHICAGO
5	3	0	3	3	6	6	4	4	6	DENVER
2	3	3	0	5	4	5	6	7	5	HOUSTON
7	7	3	5	0	8	9	1	4	8	LOS ANGELES
2	5	6	4	8	0	4	9	9	3	MIAMI
2	2	6	5	9	4	0	9	8	1	NEW YORK
7	7	4	6	1	9	9	0	2	9	SAN FRANCISCO
8	6	4	7	4	9	8	2	0	8	SEATTLE
1	1	6	5	8	3	1	9	8	0	WASHINGTON, DC

Cantidades mayores, mayor disimilaridad y viceversa



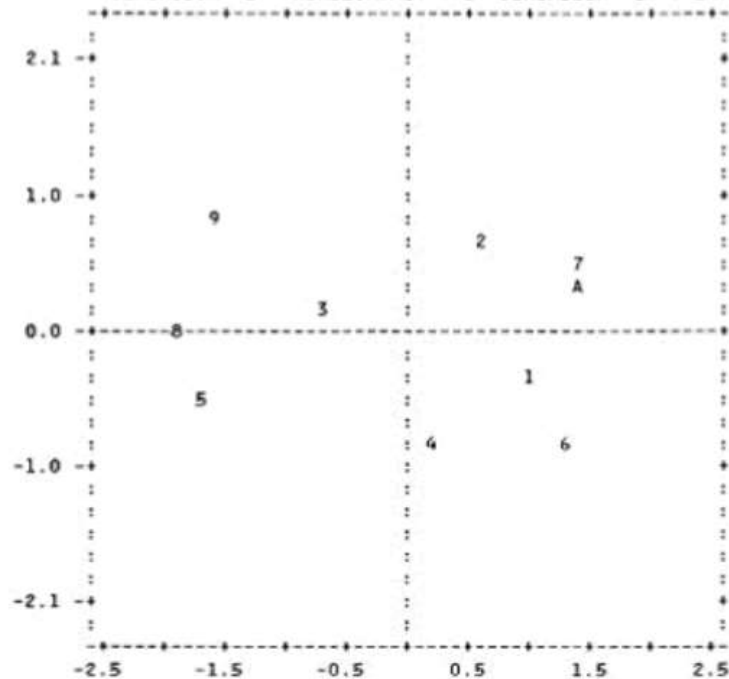
nMDS

STIMULUS COORDINATES

STIMULUS NUMBER	STIMULUS NAME	PLOT SYMBOL	DIMENSION	
			1	2
1	ATLANTA	1	0.9586	-0.3385
2	CHICAGO	2	0.6336	0.6347
3	DENVER	3	-0.7085	0.1229
4	HOUSTON	4	0.1955	-0.8345
5	LA	5	-1.6803	-0.5261
6	MIAMI	6	1.3276	-0.8150
7	NEWYORK	7	1.4289	0.5066
8	SANFRAN	8	-1.8769	-0.0782
9	SEATTLE	9	-1.6377	0.9331
10	WASHDC	A	1.3592	0.3948

S-STRESS=0.07
(Muy bueno)

DERIVED STIMULUS CONFIGURATION:
DIMENSION 1 (HORIZONTAL) VS DIMENSION 2 (VERTICAL)



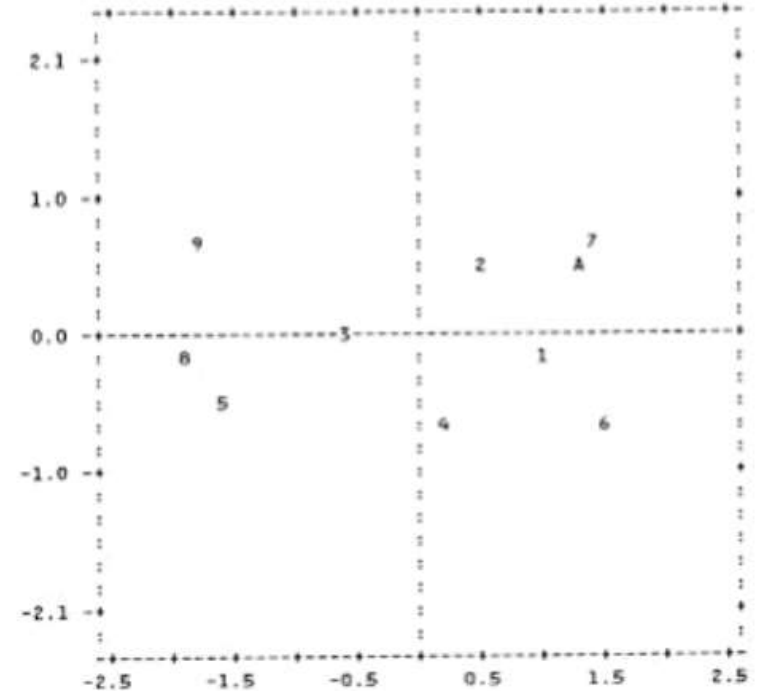
mMDS

STIMULUS COORDINATES

STIMULUS NUMBER	STIMULUS NAME	PLOT SYMBOL	DIMENSION	
			1	2
1	ATLANTA	1	0.9575	-0.1905
2	CHICAGO	2	0.5090	0.4541
3	DENVER	3	-0.6416	0.0337
4	HOUSTON	4	0.2151	-0.7631
5	LA	5	-1.6036	-0.5197
6	MIAMI	6	1.5101	-0.7752
7	NEWYORK	7	1.4284	0.6915
8	SANFRAN	8	-1.8925	-0.1500
9	SEATTLE	9	-1.7875	0.7723
10	WASHDC	A	1.3051	0.4469

S-STRESS=0.003
(Muy bueno)

DERIVED STIMULUS CONFIGURATION:
DIMENSION 1 (HORIZONTAL) VS DIMENSION 2 (VERTICAL)

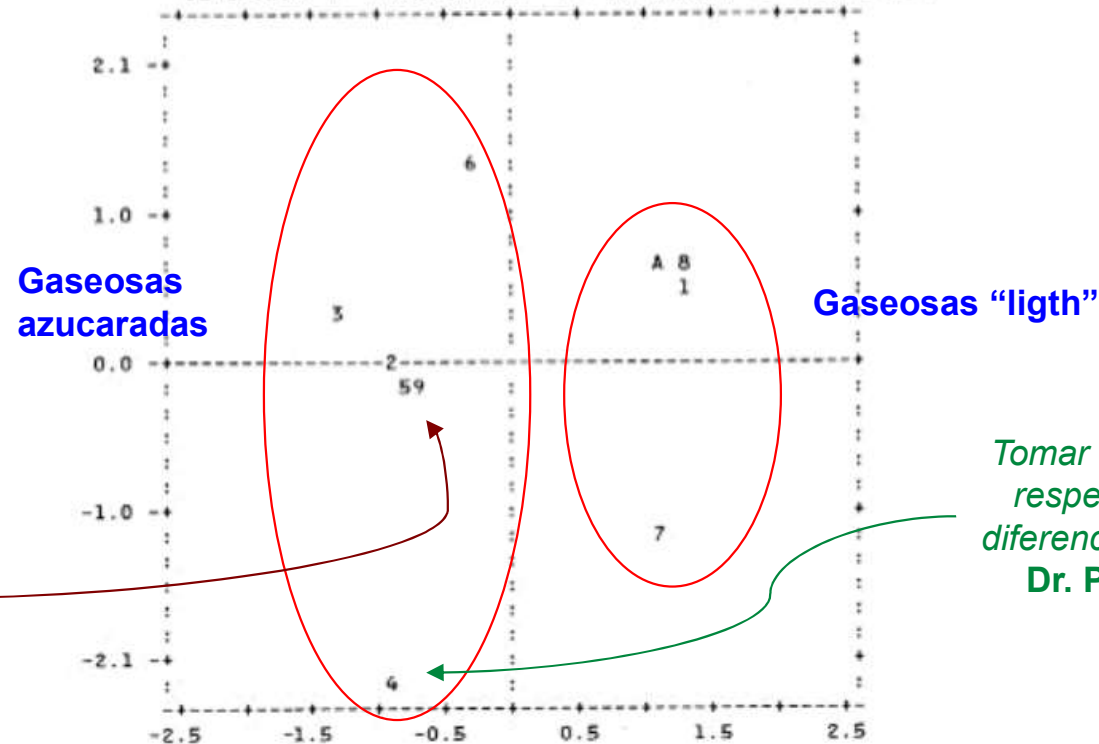




STIMULUS COORDINATES

STIMULUS NUMBER	STIMULUS NAME	PLOT SYMBOL	DIMENSION	
			1	2
1	DIETPEPS	1	1.2975	0.6838
2	RC	2	-0.8878	-0.0077
3	YUKON	3	-1.2968	0.3561
4	DRPEPPER	4	-0.8796	-2.1726
5	SHASTA	5	-0.7476	-0.1650
6	COCACOLA	6	-0.2746	1.3555
7	DIETDRPR	7	1.0675	-1.2826
8	TAB	8	1.2870	0.7553
9	PEPSI	9	-0.6862	-0.2130
10	DIETRITE	10	1.1205	0.6902

DERIVED STIMULUS CONFIGURATION:
DIMENSION 1 (HORIZONTAL) VS DIMENSION 2 (VERTICAL)





Gracias!!!

Tomado de: <https://www.quirks.com/articles/data-use-multidimensional-scaling-for-market-research>