

# Data Mining



Análisis de valores atípicos



# Análisis de Outliers

“Un *outlier* es una observación que se desvía tanto de las otras observaciones como para despertar sospechas sobre que fue generado por un mecanismo diferente”

*D. Hawkins. Identification of Outliers*

- ❑ Análisis de outliers.
  - ❑ Tipos de outliers
  - ❑ Revisión de métodos univariados: IQR, Boxplot, z-score....
  - ❑ Métodos multivariados: Local Outlier Factor (LOF)

# Tipos de outlier

## Univariado

- Son valores **atípicos** que podemos encontrar en una simple variable.
- El problema de los enfoques univariados es que son buenos para detección de extremos pero no en otros casos.

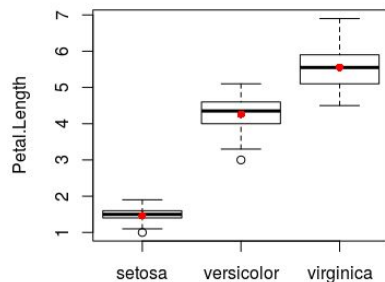
## Multivariado

- Los valores atípicos multivariados se pueden encontrar en un espacio n-dimensional.
- Para detectar valores atípicos en espacios n-dimensionales es necesario ajustar un modelo.

# Outliers

## Métodos Univariados

BoxPlot Petal.Length por Species



- ❑ IQR: Analizar los valores que están por fuera del IRQ
- ❑ Z-score.
- ❑ Identificar valores extremos a partir de 1, 2 o 3 desvíos de la media.

## Métodos Multivariados

- ❑ Análisis globales: Clustering.
  - ❑ Utilizando medidas de distancia como Mahalanobis.
- ❑ Local Outlier Factor (LOF)
  - ❑ Es un método de detección de outliers basado en distancias.
  - ❑ Calcula un score de *outlier* a partir de una distancia que se normaliza por densidad.
- ❑ Métodos basados en ensambles: IsolationForest

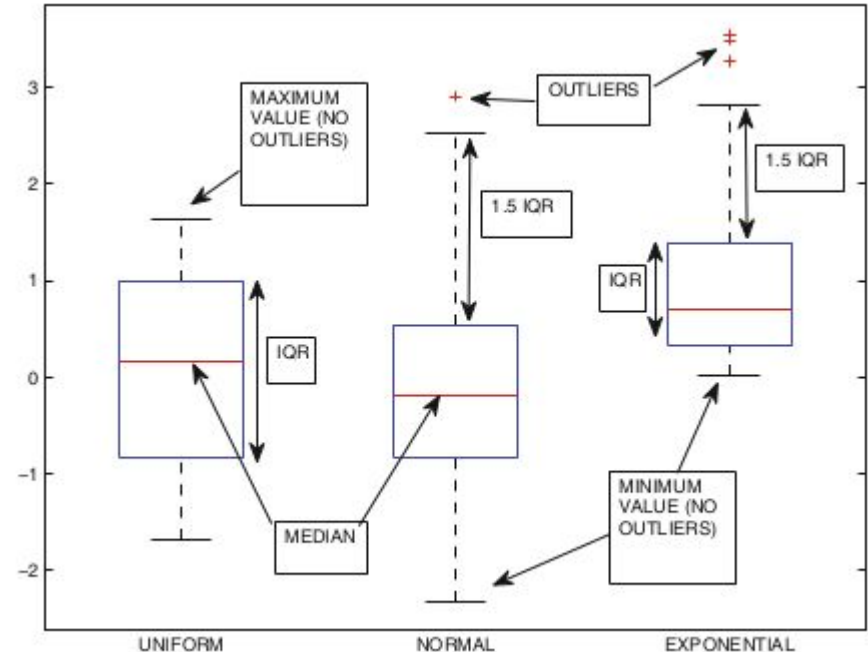
# Métodos Univariados

# Análisis de Box-Plot

Los Box-Plots permiten visualizar valores extremos univariados.

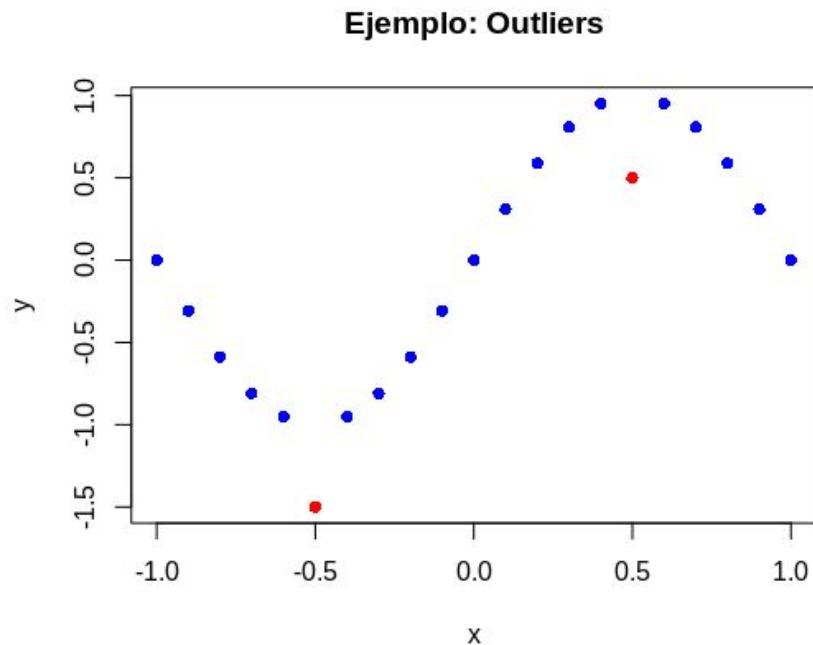
Las estadísticas de una distribución univariada se resumen en términos de cinco cantidades:

- Mínimo/máximo (bigotes)
- Primer y tercer cuantil (caja)
- Mediana (línea media de la caja)
- $IQR = Q3 - Q1$
- Generalmente la regla de decisión es  $\pm 1.5 \cdot IQR$



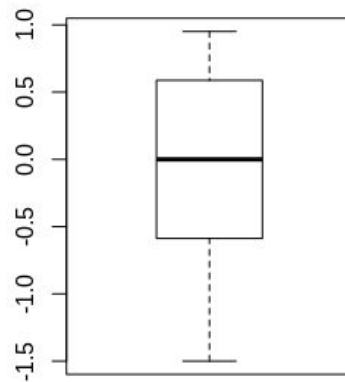
# Limitaciones del Box-Plot

En el scatter se observan dos valores atípicos.



¿Qué pasa con el box-plot?

**Box-Plot de la variable Y**

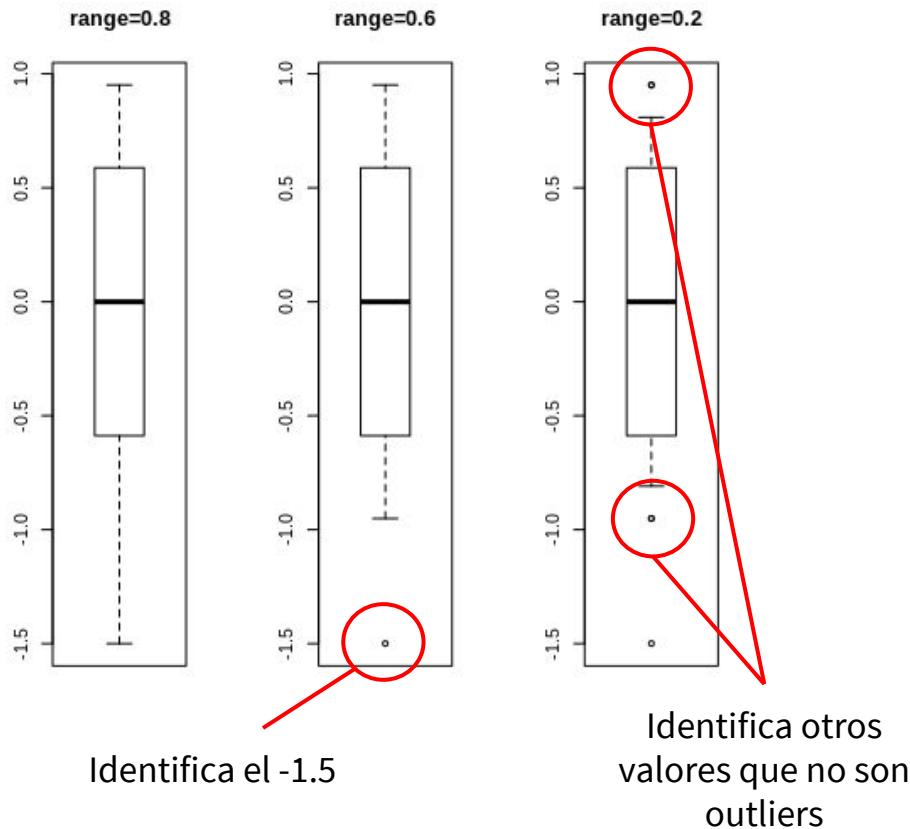


# Limitaciones del Box-Plot

Podemos intervenir modificando el largo de los bigotes del gráfico.

En R podemos hacerlo con el parámetro *range*.

También es posible modificar los cuantiles para regular el tamaño de las cajas.



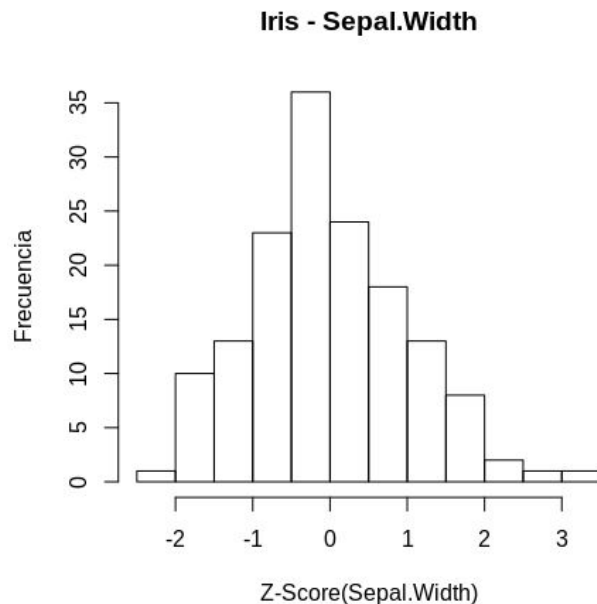


# Z-Score

Z-Score es una métrica que indica cuántas desviaciones estándar tiene una observación de la media muestral, asumiendo una distribución gaussiana.

$$Z_i = \frac{x_i - \mu}{\sigma}$$

- Cuando calculamos Z-Score para cada muestra debemos fijar un umbral.
- Un buen umbral puede ser: 2.5, 3, 3.5 o más desviaciones estándar.



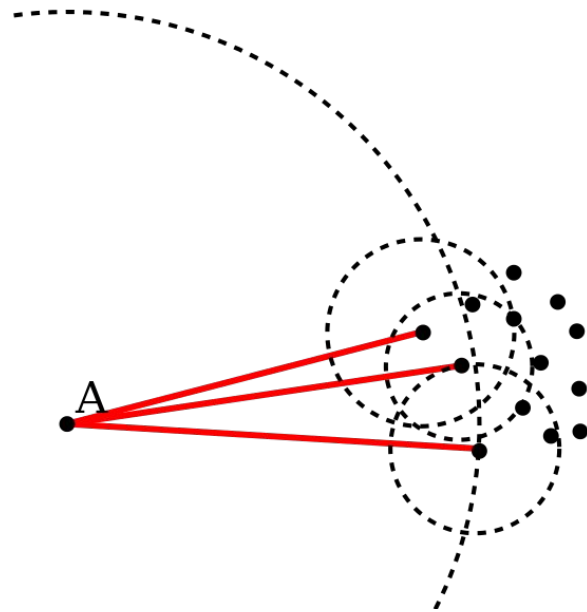
# Métodos Multivariados

# Local Outlier Factor (LOF)

El método del LOF valora puntos en un conjunto de datos multivariados cuyas filas se supone que se generan de forma independiente.

LOF es un método basado en densidad que utiliza la búsqueda de vecinos más cercanos.

El método calcula los *scores* para cada uno de los puntos a partir de la tasa promedio de densidad de los puntos vecinos con respecto a si mismo.



# Componentes del cálculo de LOF

K-Distancia de una observación  $\mathbf{o}$ ,  $\mathbf{dist}_k(\mathbf{o})$ : Es la distancia entre  $\mathbf{o}$  y su  $k$ -ésimo vecino más cercano (kNN)

Vecindad de  $\mathbf{o}$ ,  $\mathbf{N}_k(\mathbf{o}) = \{\mathbf{o}' \mid \mathbf{o}' \text{ en } \mathbf{D}, \mathbf{dist}(\mathbf{o}, \mathbf{o}') \leq \mathbf{dist}_k(\mathbf{o})\}$

La distancia de vecindad son todas las distancias dentro del radio de  $\mathbf{dist}_k(\mathbf{o})$

Distancia esperada: de  $\mathbf{o}'$  a  $\mathbf{o}$ :  $\mathbf{reachdist}(\mathbf{o} \leftarrow \mathbf{o}') = \mathbf{max}(\mathbf{dist}_k(\mathbf{o}), \mathbf{dist}(\mathbf{o}, \mathbf{o}'))$

- Densidad Local Esperada de  $\mathbf{o}$ :  
(Local Reachability Density)

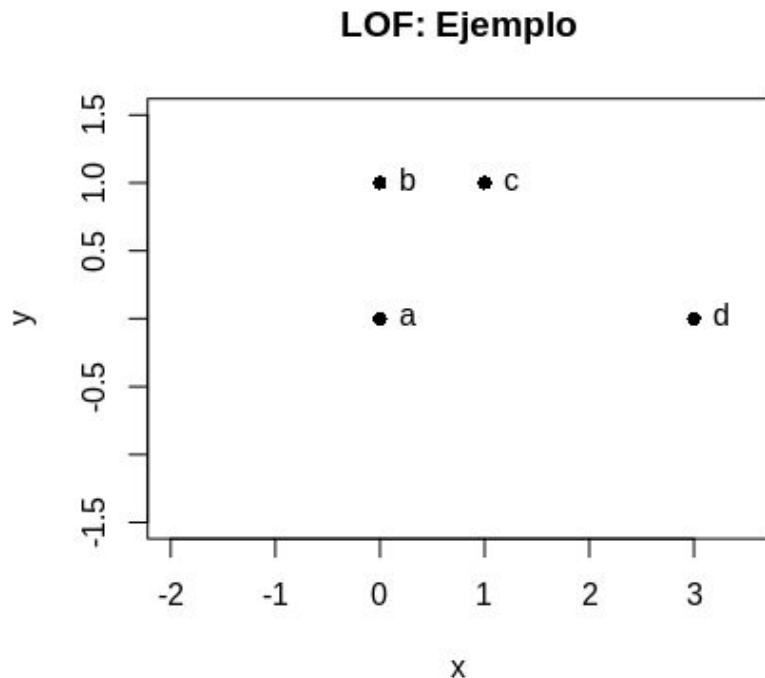
$$lrd_k(o) = \frac{\|N_k(o)\|}{\sum_{o' \in N_k(o)} reachdist_k(o' \leftarrow o)}$$

$$LOF_k(o) = \frac{\sum_{o' \in N_k(o)} \frac{lrd_k(o')}{lrd_k(o)}}{\|N_k(o)\|} = \sum_{o' \in N_k(o)} lrd_k(o') \cdot \sum_{o' \in N_k(o)} reachdist_k(o' \leftarrow o)$$

# LOF: Ejemplo

Dado el siguiente conjunto de datos:

	<b>x</b>	<b>y</b>
<b>a</b>	0	0
<b>b</b>	0	1
<b>c</b>	1	1
<b>d</b>	3	0



**Paso 1:** Calcular la matriz de distancias.

	<b>a</b>	<b>b</b>	<b>c</b>
<b>b</b>	1		
<b>c</b>	2	1	
<b>d</b>	3	4	3

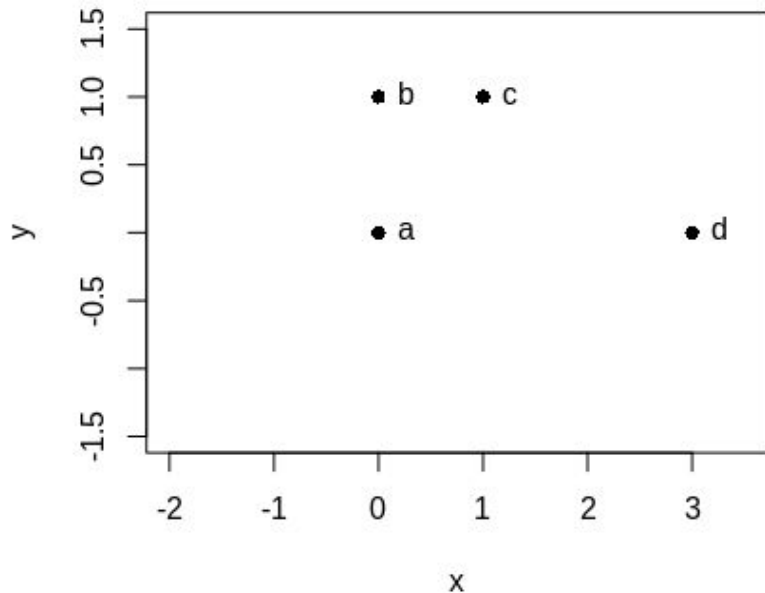
Para el ejemplo utilizamos distancia de Manhattan:

$$d_1(\mathbf{p}, \mathbf{q}) = \|\mathbf{p} - \mathbf{q}\|_1 = \sum_{i=1}^n |p_i - q_i|$$

# LOF: Ejemplo

**Paso 2:** Calculamos las distancias al k-ésimo vecino más cercano. Para este ejemplo **k=2**.

LOF: Ejemplo



$$\text{Dist}_2(a) = \text{dist}(a, c) = 2.$$

La distancia del k=2 vecino más cercano es la distancia de  $a \rightarrow c$

Recordemos la matriz de distancias para a:

El 2do NN es **c**

	a	b	c
b	1		
c	2	1	
d	3	4	3

$$\text{Dist}_2(b) = \text{dist}(b, a) = 1$$

También puede estar c

$$\text{Dist}_2(c) = \text{dist}(c, a) = 2$$

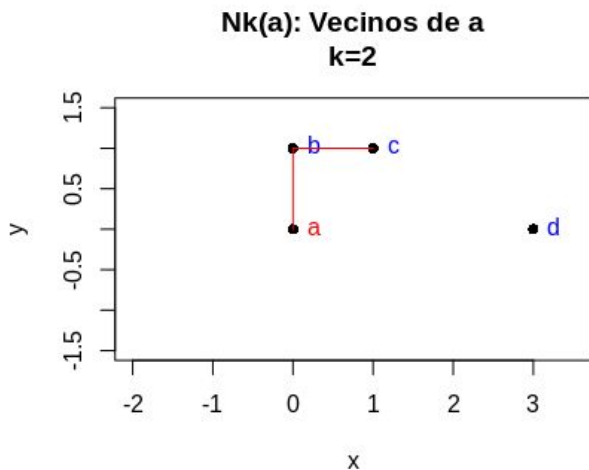
$$\text{Dist}_2(d) = \text{dist}(d, a) = 3$$

También puede estar c

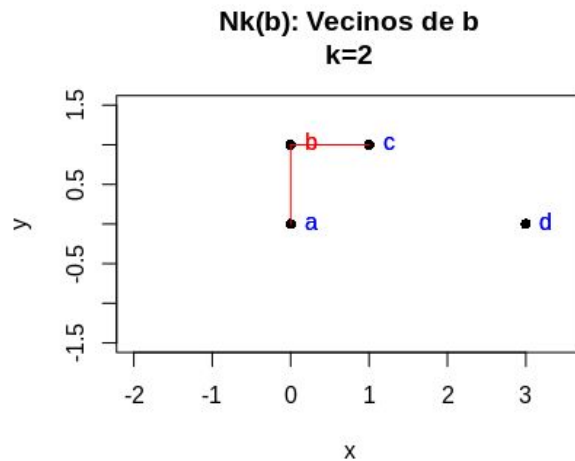
# LOF: Ejemplo

**Paso 3:** Calculamos todas las distancias de la vecindad:  $N_k(o)$

$$N_k(o) = \{o' \mid o' \text{ in } D, \text{dist}(o, o') \leq \text{dist } k(o)\}$$



$$N_2(a) = \{b, c\}$$

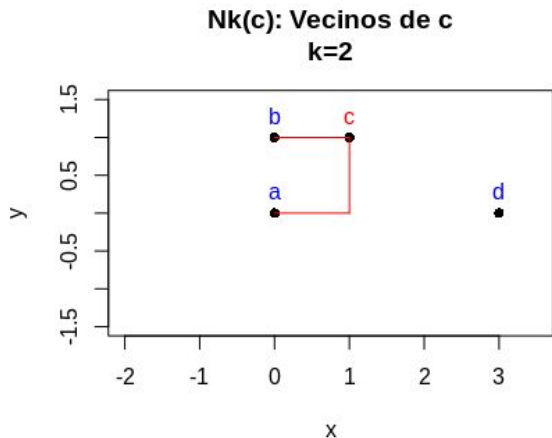


$$N_2(b) = \{a, c\}$$

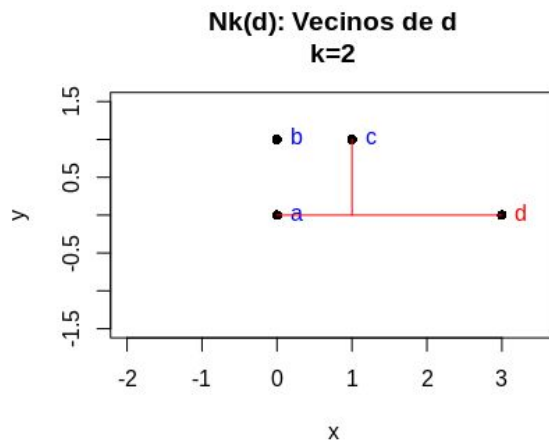
# LOF: Ejemplo

**Paso 3:** Calculamos todas las distancias de la vecindad:  $N_k(o)$

$$N_k(o) = \{o' \mid o' \text{ in } D, \text{dist}(o, o') \leq \text{dist } k(o)\}$$



$$N_2(c) = \{b, a\}$$



$$N_2(d) = \{a, c\}$$



# LOF: Ejemplo

**Paso 4:** Calculamos todos los  $\text{LRD}_k(\mathbf{o})$

$\| \mathbf{N}_k(\mathbf{o}) \|$  es la cantidad de objetos en  $\mathbf{N}_k(\mathbf{o})$ ,

por ejemplo:  $\| \mathbf{N}_2(\mathbf{a}) \| = \| \{b, c\} \| = 2$

$$\text{LDR}_k(\mathbf{a}) = \frac{\| \mathbf{N}_2(\mathbf{a}) \|}{\text{reachdist}_2(b \leftarrow a) + \text{reachdist}_2(c \leftarrow a)}$$

$$\text{LDR}_2(\mathbf{a}) = \frac{\| \mathbf{N}_2(\mathbf{a}) \|}{\text{reachdist}_2(b \leftarrow a) + \text{reachdist}_2(c \leftarrow a)} = \frac{2}{1+2} = 0.667$$

$$\text{LDR}_2(\mathbf{c}) = \frac{\| \mathbf{N}_2(\mathbf{c}) \|}{\text{reachdist}_2(a \leftarrow c) + \text{reachdist}_2(b \leftarrow c)} = \frac{2}{1+2} = 0.667$$

$$\text{lrd}_k(o) = \frac{\| \mathbf{N}_k(o) \|}{\sum_{o' \in \mathbf{N}_k(o)} \text{reachdist}_k(o' \leftarrow o)}$$

$$\text{reachdist}_k(o \leftarrow o') = \max\{\text{dist}_k(o), \text{dist}(o, o')\}$$

$$\text{reachdist}_2(b \leftarrow a) = \max\{\text{dist}_2(b), \text{dist}(b, a)\} = \max\{1, 1\} = 1$$

$$\text{reachdist}_2(c \leftarrow a) = \max\{\text{dist}_2(c), \text{dist}(c, a)\} = \max\{2, 2\} = 2$$

$$\text{LDR}_2(\mathbf{b}) = \frac{\| \mathbf{N}_2(\mathbf{b}) \|}{\text{reachdist}_2(a \leftarrow b) + \text{reachdist}_2(c \leftarrow b)} = \frac{2}{2+2} = 0.5$$

$$\text{LDR}_2(\mathbf{d}) = \frac{\| \mathbf{N}_2(\mathbf{d}) \|}{\text{reachdist}_2(a \leftarrow d) + \text{reachdist}_2(c \leftarrow d)} = \frac{2}{3+3} = 0.33$$

# LOF: Ejemplo

**Paso 5:** Calcular  $LOF_k(o)$

$$LOF_k(o) = \frac{\sum_{o' \in N_k(o)} \frac{lrd_k(o')}{lrd_k(o)}}{\|N_k(o)\|} = \sum_{o' \in N_k(o)} lrd_k(o') \cdot \sum_{o' \in N_k(o)} reachdist_k(o' \leftarrow o)$$

$$LOF_2(a) = (lrd_2(b) + lrd_2(c)) * (reachdist_2(b \leftarrow a) + reachdist_2(c \leftarrow a)) = (0.5 + 0.667) * (1 + 2) = 3.501$$

$$LOF_2(b) = (lrd_2(a) + lrd_2(c)) * (reachdist_2(a \leftarrow b) + reachdist_2(c \leftarrow b)) = (0.667 + 0.667) * (2 + 2) = 5.336$$

$$LOF_2(c) = (lrd_2(b) + lrd_2(a)) * (reachdist_2(b \leftarrow c) + reachdist_2(a \leftarrow c)) = (0.5 + 0.667) * (1 + 2) = 3.501$$

$$LOF_2(d) = (lrd_2(a) + lrd_2(c)) * (reachdist_2(a \leftarrow d) + reachdist_2(c \leftarrow d)) = (0.667 + 0.667) * (3 + 3) = 8.004$$

# LOF: Ejemplo

**Paso 6:** Ordenamos todos los valores de  $\text{LOF}_k(o)$

$$\text{LOF}_2(d) = \mathbf{8.004} \leftarrow$$

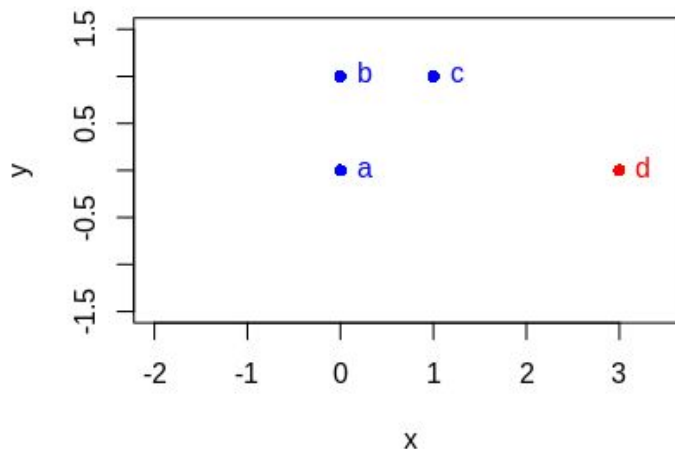
$$\text{LOF}_2(b) = 5.336$$

$$\text{LOF}_2(a) = 3.501$$

$$\text{LOF}_2(c) = 3.501$$

Una vez que ordenamos tenemos que definir un umbral de corte.

En este caso el punto **d** es un claro valor atípico.



# Bibliografía

- ❑ Jiawei Han, Micheline Kamber, Jian Pei. 2012. Tercera edición. Data Mining: Concepts and Techniques. Cap. 2 y Cap. 3
- ❑ [http://www.cse.ust.hk/~leichen/courses/comp5331/lectures/LOF\\_Example.pdf](http://www.cse.ust.hk/~leichen/courses/comp5331/lectures/LOF_Example.pdf)