

# Validación de Clusters

Bárbara Poblete

## ¿Es necesario validar los clusters?

- Por lo menos en Clasificación, la validación es parte integral del proceso
- No así en Clustering...

## aprendizaje supervisado vs. no supervisado

Supervised vs. Unsupervised Learning

Supervised	Unsupervised
<ul style="list-style-type: none"> <li>• <math>y=F(x)</math>: true function</li> <li>• D: labeled training set</li> <li>• D: <math>\{x_i, y_i\}</math></li> <li>• <math>y=G(x)</math>: model trained to predict labels D</li> <li>• Goal: <math>E \leftarrow (F(x) - G(x))^2 \approx 0</math></li> <li>• Well defined criteria: Accuracy, RMSE, ...</li> </ul>	<ul style="list-style-type: none"> <li>• Generator: true model</li> <li>• D: unlabeled data sample</li> <li>• D: <math>\{x_i\}</math></li> <li>• Learn</li> <li>• Goal: ??????????</li> <li>• Well defined criteria: ??????????</li> </ul>

source: <http://www.cs.put.poznan.pl/jstefanowski/sed/DM-7clusteringnew.pdf>

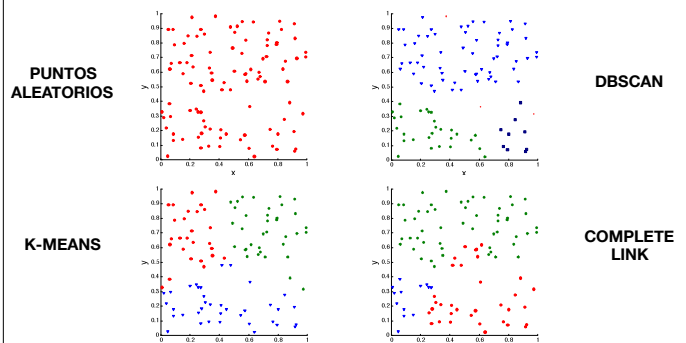
## ¿Cómo saber si nuestros clusters son buenos?

- No hay una respuesta absoluta
- Depende de la aplicación
- ¿Entonces, para qué evaluar?
  - ★ Clustering es muchas veces parte de un proceso exploratorio
  - ★ Evaluar parece innecesario en estos casos (pero no!)
  - ★ Cada algoritmo parece necesitar su propio tipo de evaluación
  - ★ k-means: SSE, pero no funciona para DBSCAN
  - ★ Es esencial! Porque siempre podemos encontrar clusters (hasta en datos aleatorios).

## Evaluamos para:

- Evitar encontrar patrones en el ruido
- Para comparar algoritmos de clustering diferentes
- Para comparar conjuntos de clusters diferentes
- Para comprar dos clusters

## CLUSTERS EN DATOS ALEATORIOS

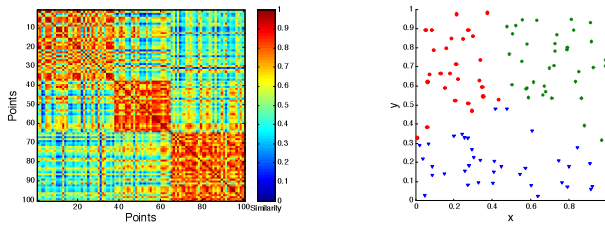
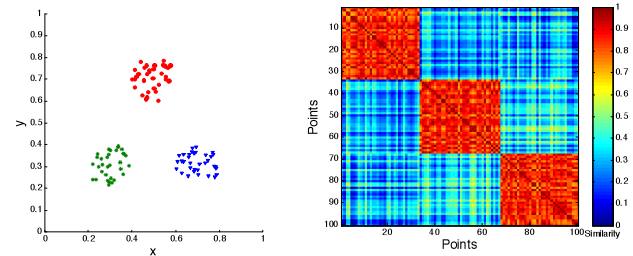


$$\text{Corr} = -0.5810$$

## Enfoque visual

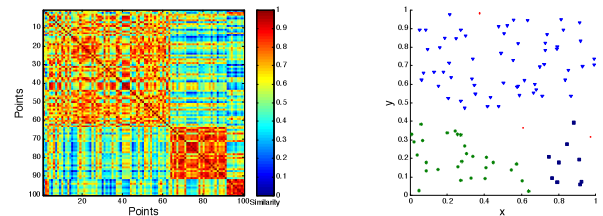
- Ordenar la matriz de similitud con respecto a etiquetas de clusters e inspeccionar visualmente

## Visualizando la matriz de similitud (clusters reales)

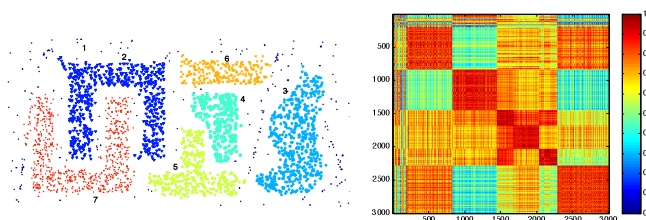


K-means

## Visualizando clusters sobre datos aleatorios



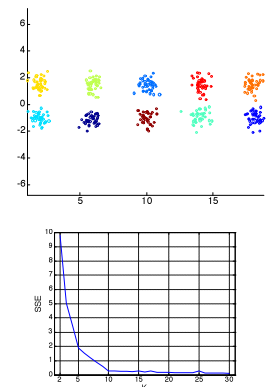
DBSCAN

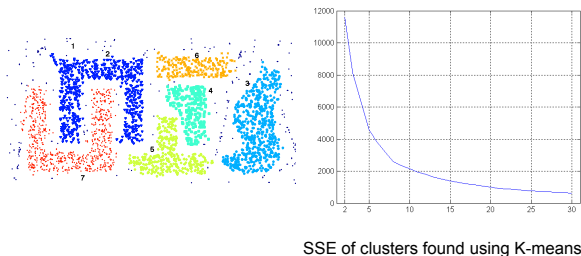


DBSCAN

## Medidas internas: SSE

- Clusters en figuras más complicadas no están bien separados
- Índice interno: SSE
- Permite comparar 2 clusters, o 2 soluciones de clustering
- Permite estimar el número de clusters





Curva SSE para un dataset más complicado

## Metodología para validar clusters

- Necesidad de contar con una metodología para interpretar cualquier medida (¿qué es bueno? ¿qué no?)
- Usamos la estadística para crear una metodología

## Metodología para validar clusters

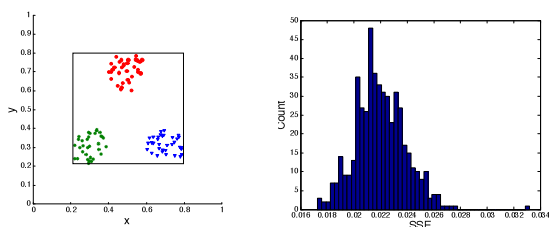
- Mientras más **atípico** es un resultado, más probable que sea reflejo de estructuras válidas
- Podemos comparar índices que resultan de datos aleatorios, con los de nuestros datos
- valores poco probables indican resultados válidos

## Metodología para validar clusters

- Al comparar resultados de dos clustering (dos cluster sets), no es muy necesario usar una metodología
- Pero en este caso la pregunta es si la diferencia es **significativa** (estadísticamente - repetible y en magnitud)

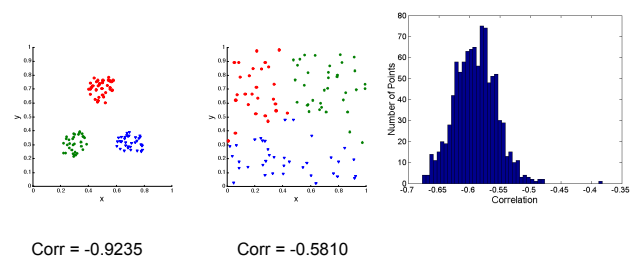
## Metodología: Ejemplo SSE

- Comparar SSE = 0.005 contra 3 clusters de datos aleatorios
- Histograma muestra distribución SSE para 500 sets de datos aleatorios (100 puntos), en el mismo rango



## Otro ejemplo: Correlación

- Correlación entre matrices de incidencia y proximidad para 2 sets de datos



## Medidas internas: Cohesión y separación

- **Cohesión de clusters:** mide qué tan cercanos son los objetos en un cluster (ej: SSE)
- **Separación de clusters:** mide qué tan diferente o bien separado es un cluster de otros

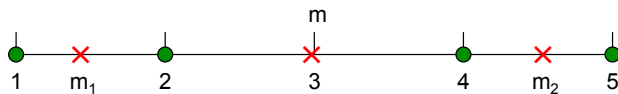
## EJ. (SSE) Cohesión y Separación

- Cohesión se mide como **within cluster sum of squares (SSE)**

$$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

- Separación se mide como **between cluster sum of squares (BSS)**

$$BSS = \sum_i |C_i| (m - m_i)^2$$



K=1 cluster:

$$WSS = (1 - 3)^2 + (2 - 3)^2 + (4 - 3)^2 + (5 - 3)^2 = 10$$

$$BSS = 4 \times (3 - 3)^2 = 0$$

$$Total = 10 + 0 = 10$$

K=2 clusters:

$$WSS = (1 - 1.5)^2 + (2 - 1.5)^2 + (4 - 4.5)^2 + (5 - 4.5)^2 = 1$$

$$BSS = 2 \times (3 - 1.5)^2 + 2 \times (4.5 - 3)^2 = 9$$

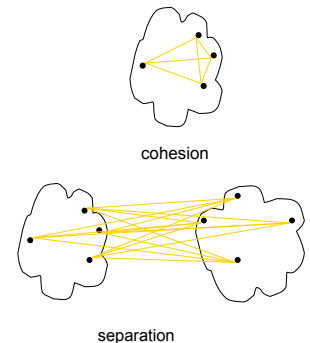
$$Total = 1 + 9 = 10$$

## Medidas internas: Cohesión y separación

- Enfoque basado en grafos de proximidad

- Cohesión: suma de los pesos de todos los arcos en un cluster

- Separación: suma de los pesos entre nodos del cluster y de otros clusters



## Medidas internas: Coeficiente de Silhouette

- Combina ideas de cohesión y separación, pero para puntos individuales, como también para clusters y clusterings (estos últimos son promedios)
- Para un punto individual,  $i$ 
  - Calcular  $a$  = distancia promedio de  $i$  a los puntos de su cluster
  - Calcular  $b$  =  $\min$ (distancia promedio de  $i$  a puntos de otro cluster)
  - $s = (b - a) / \max(a, b)$
  - valores entre -1 y 1, mientras más cerca a 1 mejor (y a más cerca de 0)

## Pureza y Entropía

- Pureza: Nivel en que un cluster contiene elementos de una sola clase (se usa la clase predominante)
- Entropía: Cantidad de clases diferentes que contiene un cluster

## Medidas externas: Pureza y entropía

**Table 5.9.** K-means Clustering Results for LA Document Data Set

Cluster	Entertainment	Financial	Foreign	Metro	National	Sports	Entropy	Purity
1	3	5	40	506	96	27	1.2270	0.7474
2	4	7	280	29	39	2	1.1472	0.7756
3	1	1	1	7	4	671	0.1813	0.9796
4	10	162	3	119	73	2	1.7487	0.4390
5	331	22	5	70	13	23	1.3976	0.7134
6	5	358	12	212	48	13	1.5523	0.5525
Total	354	555	341	943	273	738	1.1450	0.7203

**entropy** For each cluster, the class distribution of the data is calculated first, i.e., for cluster  $j$  we compute  $p_{ij}$ , the 'probability' that a member of cluster  $j$  belongs to class  $i$  as follows:  $p_{ij} = m_{ij}/m_j$ , where  $m_j$  is the number of values in cluster  $j$  and  $m_{ij}$  is the number of values of class  $i$  in cluster  $j$ . Then using this class distribution, the entropy of each cluster  $j$  is calculated using the standard formula  $e_j = -\sum_{i=1}^L p_{ij} \log_2 p_{ij}$ , where the  $L$  is the number of classes. The total entropy for a set of clusters is calculated as the sum of the entropies of each cluster weighted by the size of each cluster, i.e.,  $e = \sum_{j=1}^K \frac{m_j}{m} e_j$ , where  $m_j$  is the size of cluster  $j$ ,  $K$  is the number of clusters, and  $m$  is the total number of data points.

**purity** Using the terminology derived for entropy, the purity of cluster  $j$ , is given by  $\text{purity}_j = \max_i p_{ij}$  and the overall purity of a clustering by  $\text{purity} = \sum_{j=1}^K \frac{m_j}{m} \text{purity}_j$ .

## Validación con Expertos

- Se pueden evaluar los clusters para ver si producen el resultado esperado y comparar con otras soluciones
- Se puede generar una clasificación de validación

## Comentario final

- La etapa de validación es la parte más difícil y frustrante del análisis de clusters
- Sin embargo es necesario
- Idealmente se deben combinar medidas externas e internas