

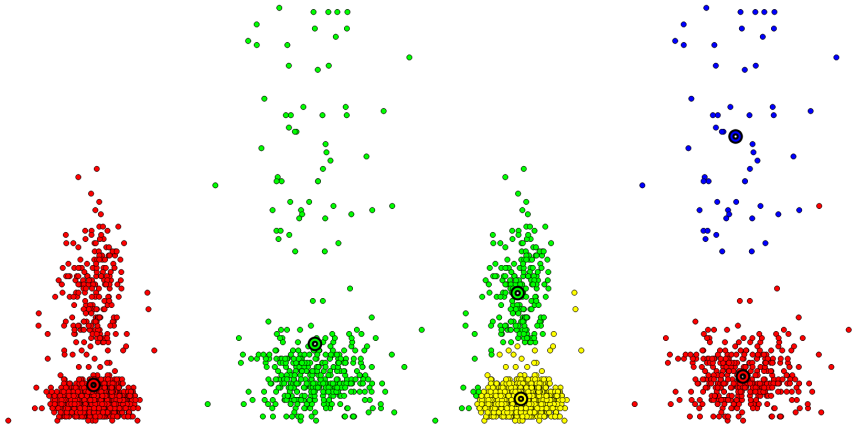
Aprendizaje no supervisado

1.3. Medidas de evaluación de agrupamientos

Javier Sevilla

Evaluación

Extrínseca vs. Intrínseca



Evaluación extrínseca

- ▶ Se conoce el agrupamiento real (**verdad básica**)
- ▶ Se compara el resultado del algoritmo con la verdad básica

Evaluación extrínseca

- ▶ Se conoce el agrupamiento real (**verdad básica**)
- ▶ Se compara el resultado del algoritmo con la verdad básica
- ▶ No existe el concepto de etiqueta:
búsqueda de la correspondencia entre clúster real y predicho

Evaluación intrínseca

- ▶ No se conoce el agrupamiento real (**verdad básica**), ni se sabe si existe
- ▶ Se mide la congruencia del agrupamiento
- ▶ Diferentes criterios posibles

$\{B_l\}_{l=1}^{K'}$: Verdad básica

$\{C_k\}_{k=1}^K$: Agrupamiento resultante de un algoritmo de *clustering*

$\{\mathbf{c}_k\}_{k=1}^K$: Centro(ide)s de los clústeres resultantes

$n_l = |B_l|$: Tamaño de un clúster verdadero

$n_k = |C_k|$: Tamaño de un clúster resultante

$n_{kl} = |C_k \cap B_l|$: Número de ejemplos que comparten un clúster resultante y otro verdadero

Error:

$$E = 1 - \frac{1}{n} \max_{\sigma} \sum_{l=1}^{K'} n_{\sigma(l)l}$$

donde σ es una función de $\sigma : \{1, \dots, K'\} \rightarrow \{1, \dots, K\}$

Error:

$$E = 1 - \frac{1}{n} \max_{\sigma} \sum_{l=1}^{K'} n_{\sigma(l)l}$$

donde σ es una función de $\sigma : \{1, \dots, K'\} \rightarrow \{1, \dots, K\}$

- ▶ Recorrido sobre los clústeres reales
- ▶ Máximo (optimista) para identificar la correspondencia $C-B$

Evaluación

Extrínseca

Precisión:

$$P_{kl} = \frac{n_{kl}}{n_{k\cdot}}$$

Recall:

$$R_{lk} = \frac{n_{kl}}{n_{\cdot l}}$$

Precisión:

$$P_{kl} = \frac{n_{kl}}{n_k}$$

Recall:

$$R_{lk} = \frac{n_{kl}}{n_l}$$

- ▶ Medidas entre un clúster real y otro resultante
- ▶ Precisión: ¿Cuántos de los elementos del clúster resultante k lo son también del clúster real l ?
- ▶ *Recall*: ¿Cuántos de los elementos del clúster real l lo son también del clúster resultante k ?

Pureza:

$$P_u = \sum_{k=1}^K \frac{n_k}{n} \max_{l \in \{1, \dots, K'\}} P_{kl}$$

Pureza:

$$P_u = \sum_{k=1}^K \frac{n_k}{n} \max_{l \in \{1, \dots, K'\}} P_{kl}$$

- ▶ Media ponderada de la precisión
- ▶ Recorrido sobre los clústeres resultantes
- ▶ Máximo (optimista) para identificar la correspondencia $C-B$

Medida F:

$$F1 = \sum_{l=1}^{K'} \frac{n_{.l}}{n} \max_{k \in \{1, \dots, K\}} \left(\frac{2P_{kl}R_{lk}}{P_{kl} + R_{lk}} \right)$$

Medida F:

$$F1 = \sum_{l=1}^{K'} \frac{n_{.l}}{n} \max_{k \in \{1, \dots, K\}} \left(\frac{2P_{kl}R_{lk}}{P_{kl} + R_{lk}} \right)$$

- ▶ Media ponderada de la media armónica de la precisión y el *recall*
- ▶ Recorrido sobre los clústeres reales
- ▶ Máximo (optimista) para identificar la correspondencia $C-B$

Entropía:

$$H = - \sum_{k=1}^K \frac{n_{k\cdot}}{n} \sum_{l=1}^{K'} \frac{n_{kl}}{n_{k\cdot}} \log \frac{n_{kl}}{n_{k\cdot}}$$

Entropía:

$$H = - \sum_{k=1}^K \frac{n_{k\cdot}}{n} \sum_{l=1}^{K'} \frac{n_{kl}}{n_{k\cdot}} \log \frac{n_{kl}}{n_{k\cdot}}$$

- ▶ Media ponderada de la entropía de cada clúster resultante
- ▶ Entropía: mide cómo se distribuyen los ejemplos de un clúster resultante entre los clústeres reales (crece a mayor desorden)
- ▶ Recorrido (principal) los clústeres resultantes

Información mutua:

$$I = \sum_{k=1}^K \sum_{l=1}^{K'} \frac{n_{kl}}{n} \log \frac{n \cdot n_{kl}}{n_{k\cdot} \cdot n_{\cdot l}}$$

Información mutua:

$$I = \sum_{k=1}^K \sum_{l=1}^{K'} \frac{n_{kl}}{n} \log \frac{n \cdot n_{kl}}{n_{k\cdot} \cdot n_{\cdot l}}$$

- ▶ Información mutua entre dos agrupamientos (real y resultante)
- ▶ I: mide cómo se explican mutuamente ambos agrupamientos

¿Es razonable asumir la existencia de la verdad básica?

¿Es razonable asumir la existencia de la verdad básica?

¿Para qué queremos entonces un algoritmo de *clustering*?

La raíz del cuadrado de la media de la desviación típica:

$$RMSSTD = \sqrt{\frac{\sum_{k=1}^K \sum_{\mathbf{x}_i \in C_k} \|\mathbf{x}_i - \mathbf{c}_k\|^2}{v \cdot \sum_{k=1}^K (|C_k| - 1)}}$$

La raíz del cuadrado de la media de la desviación típica:

$$RMSSTD = \sqrt{\frac{\sum_{k=1}^K \sum_{\mathbf{x}_i \in C_k} \|\mathbf{x}_i - \mathbf{c}_k\|^2}{v \cdot \sum_{k=1}^K (|C_k| - 1)}}$$

- Mide la heterogeneidad de los clústeres
- Se reduce fácilmente aumentando el número de clústeres resultantes K

Medida R -cuadrado

$$R^2 = \frac{\sum_{\mathbf{x}_i} \|\mathbf{x}_i - \mathbf{c}\|^2 - \sum_{k=1}^K \sum_{\mathbf{x}_{i'} \in C_k} \|\mathbf{x}_{i'} - \mathbf{c}_k\|^2}{\sum_{\mathbf{x}_i} \|\mathbf{x}_i - \mathbf{c}\|^2}$$

donde \mathbf{c} es el centro de todo el *dataset*.

Medida R -cuadrado

$$R^2 = \frac{\sum_{\mathbf{x}_i} \|\mathbf{x}_i - \mathbf{c}\|^2 - \sum_{k=1}^K \sum_{\mathbf{x}_{i'} \in C_k} \|\mathbf{x}_{i'} - \mathbf{c}_k\|^2}{\sum_{\mathbf{x}_i} \|\mathbf{x}_i - \mathbf{c}\|^2}$$

donde \mathbf{c} es el centro de todo el *dataset*.

- ▶ Mide la homogeneidad de los clústeres
- ▶ Acotada entre 0 (sólo un clúster) y 1 ($K = n$)
- ▶ Se incrementa fácilmente aumentando el número de clústeres resultantes K

Silueta

$$S = \frac{1}{n} \sum_{\mathbf{x}_i} \frac{b_k(\mathbf{x}_i) - a_k(\mathbf{x}_i)}{\max\{b_k(\mathbf{x}_i), a_k(\mathbf{x}_i)\}}$$

donde

$$a_k(\mathbf{x}_i) = \frac{1}{n_k - 1} \sum_{\mathbf{x}_j \in C_k: \mathbf{x}_j \neq \mathbf{x}_i} d(\mathbf{x}_i, \mathbf{x}_j)$$

y

$$b_k(\mathbf{x}_i) = \min_{h \neq k} \frac{1}{n_h} \sum_{\mathbf{x}_j \in C_h} d(\mathbf{x}_i, \mathbf{x}_j)$$

Silueta

$$S = \frac{1}{n} \sum_{\mathbf{x}_i} \frac{b_k(\mathbf{x}_i) - a_k(\mathbf{x}_i)}{\max\{b_k(\mathbf{x}_i), a_k(\mathbf{x}_i)\}}$$

donde

$$a_k(\mathbf{x}_i) = \frac{1}{n_k - 1} \sum_{\mathbf{x}_j \in C_k: \mathbf{x}_j \neq \mathbf{x}_i} d(\mathbf{x}_i, \mathbf{x}_j)$$

y

$$b_k(\mathbf{x}_i) = \min_{h \neq k} \frac{1}{n_h} \sum_{\mathbf{x}_j \in C_h} d(\mathbf{x}_i, \mathbf{x}_j)$$

- Diferencia normalizada entre la distancia intraclúster y la interclúster
- Acotada entre -1 y 1

Índice Calinski-Harabasz:

$$CH = \frac{(n - K) \sum_{k=1}^K n_k \cdot d(\mathbf{c}_k, \mathbf{c})^2}{(K - 1) \sum_{k=1}^K \sum_{\mathbf{x}_i \in C_k} d(\mathbf{x}_i, \mathbf{c}_k)^2}$$

Índice Calinski-Harabasz:

$$CH = \frac{(n - K) \sum_{k=1}^K n_k \cdot d(\mathbf{c}_k, \mathbf{c})^2}{(K - 1) \sum_{k=1}^K \sum_{\mathbf{x}_i \in C_k} d(\mathbf{x}_i, \mathbf{c}_k)^2}$$

- ▶ Suma promedio de las distancias inter e intraclúster al cuadrado
- ▶ A mayor valor, mejor agrupamiento

Índice I:

$$I = \left(\frac{\sum \mathbf{x}_i d(\mathbf{x}_i, \mathbf{c})}{K \cdot \sum_{k=1}^K \sum_{\mathbf{x}_i \in C_k} d(\mathbf{x}_i, \mathbf{c}_k)} \cdot \max_{i,j \in \{1, \dots, K\}} d(\mathbf{c}_i, \mathbf{c}_j) \right)^p$$

Índice I:

$$I = \frac{\sum_{\mathbf{x}_i} d(\mathbf{x}_i, \mathbf{c})}{K \cdot \sum_{k=1}^K \sum_{\mathbf{x}_i \in C_k} d(\mathbf{x}_i, \mathbf{c}_k)} \cdot \left(\max_{i,j \in \{1, \dots, K\}} d(\mathbf{c}_i, \mathbf{c}_j) \right)^p$$

- Mide la separación interclúster con respecto a la homogeneidad intraclúster
- A mayor valor, mejor agrupamiento

Gracias