

# Metaheurísticas

## Seminario 2.a. Problemas de optimización con técnicas basadas en trayectorias simples

---

### 1. Problema de Asignación Cuadrática (QAP)

- Definición del Problema
- Representación y Ejemplo: Diseño de un Hospital. Solución Greedy
- Búsquedas por Trayectorias Simples
- La Biblioteca QAPLIB

### 2. Problema de la Selección de Características

# Definición del Problema

---

- El Problema de Asignación Cuadrática (QAP) está considerado como uno de los problemas de optimización combinatoria más complejos
- Es NP-completo. Incluso la resolución de problemas pequeños de tamaño  $n > 25$  se considera una tarea computacionalmente muy costosa
- El problema general consiste en encontrar la asignación óptima de  $n$  unidades a  $n$  localizaciones, conociendo la distancia entre las primeras y el flujo existente entre las segundas

# Definición del Problema

---

- Sean  $n$  unidades ( $u_i, i=1,\dots,n$ ) y  $n$  localizaciones ( $l_j, j=1,\dots,n$ ). Se dispone de dos matrices  $F=(f_{ij})$  y  $D=(d_{ij})$ , de dimensión  $(n \times n)$  con la siguiente interpretación:
  - $F$  es la matriz de flujo, es decir,  $f_{ij}$  es el flujo que circula entre la unidad  $i$  y la  $j$
  - $D$  es la matriz de distancias, es decir,  $d_{kl}$  es la distancia existente entre la localización  $k$  y la  $l$
- El costo de asignar simultáneamente  $u_i$  a  $l_k$  y  $u_j$  a  $l_l$  es:

$$f_{ij} \cdot d_{kl}$$

# Definición del Problema

---

- La definición matemática del problema consiste en minimizar el costo de las asignaciones:

$$\min \quad \sum_{i,j=1}^n \sum_{k,p=1}^n f_{ij} d_{kp} x_{ij} x_{kp}$$

$$\text{s.a.} \quad \sum_{i=1}^n x_{ij} = 1 \quad 1 \leq j \leq n ,$$

$$\sum_{j=1}^n x_{ij} = 1 \quad 1 \leq i \leq n ,$$

$$x_{ij} \in \{0,1\} \quad 1 \leq i, j \leq n .$$

- Como se puede ver, la función de coste es cuadrática, lo que da nombre al problema y lo complica sustancialmente

# Definición del Problema

---

## ■ Problema de la asignación cuadrática, *QAP*:

*Dadas  $n$  unidades y  $n$  localizaciones posibles, el problema consiste en determinar la asignación óptima de las unidades en las localizaciones conociendo el flujo existente entre las primeras y la distancia entre las segundas*

- Si se considera una permutación para representar las asignaciones, se verifican directamente las restricciones del problema:

$$QAP = \min_{S \in \Pi_N} \left( \sum_{i=1}^n \sum_{j=1}^n f_{ij} \cdot d_{S(i)S(j)} \right)$$

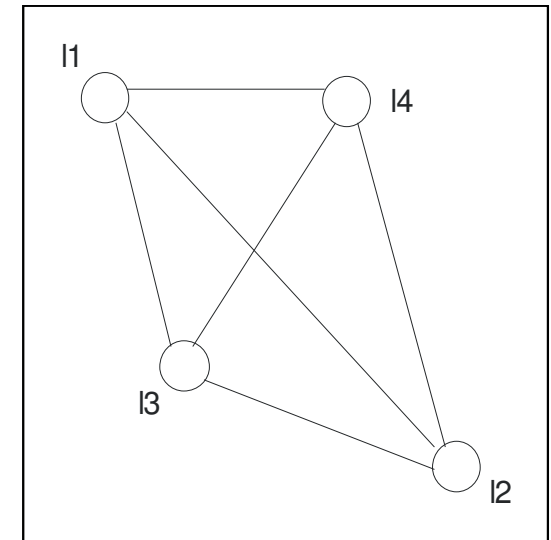
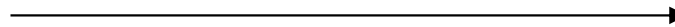
# Representación y Ejemplo: Diseño de un Hospital

---

Elshafei, A.N. (1977). Hospital Layout as a Quadratic Assignment Problem. *Operations Research Quarterly* 28, 167-179.

- Supongamos que se ha de diseñar un hospital que comprende cuatro unidades distintas:
  - u1: Maternidad
  - u2: Urgencias
  - u3: Unidad de Cuidados Intensivos
  - u4: Cirugía

Que han de ser situadas en  
un edificio con la siguiente  
distribución:

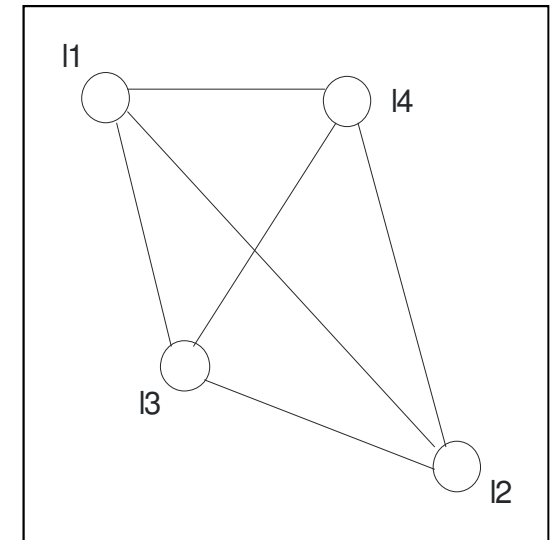


# Representación y Ejemplo: Diseño de un Hospital

- La matriz  $D$  contiene las distancias existentes entre las diferentes salas
- La matriz  $F$  recoge el número medio de pacientes que pasan de una unidad a otra cada hora (por ejemplo, podrían ser las medias mensuales, medias anuales, totales anuales, ...)

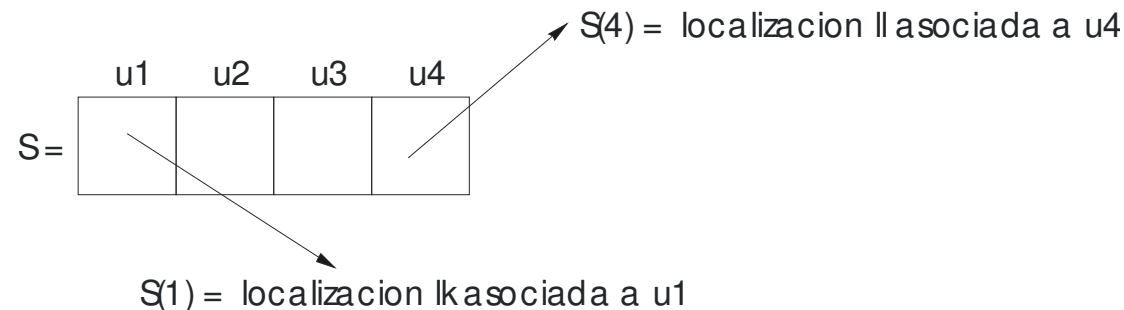
$$D = \begin{pmatrix} 0 & 12 & 6 & 4 \\ 12 & 0 & 6 & 8 \\ 6 & 6 & 0 & 7 \\ 4 & 8 & 7 & 0 \end{pmatrix}$$

$$F = \begin{pmatrix} 0 & 3 & 8 & 3 \\ 3 & 0 & 2 & 4 \\ 8 & 2 & 0 & 5 \\ 3 & 4 & 5 & 0 \end{pmatrix}$$



# Representación y Ejemplo: Diseño de un Hospital

- Las soluciones son permutaciones del conjunto  $N=\{1, 2, 3, 4\}$
- Para entenderlo mejor, podemos pensar en su representación en forma de vector permutación: *las posiciones se corresponden con las unidades y el contenido de las mismas con las localizaciones (salas del hospital) en la que se sitúan las unidades correspondientes:*



- En este caso, usamos una permutación para representar una asignación, al contrario que en el TSP en el que representa un orden
- Esto nos permite verificar de forma sencilla las restricciones del problema, lo que sería más complicado en otras representaciones como una matriz binaria



# Representación y Ejemplo: Diseño de un Hospital

---

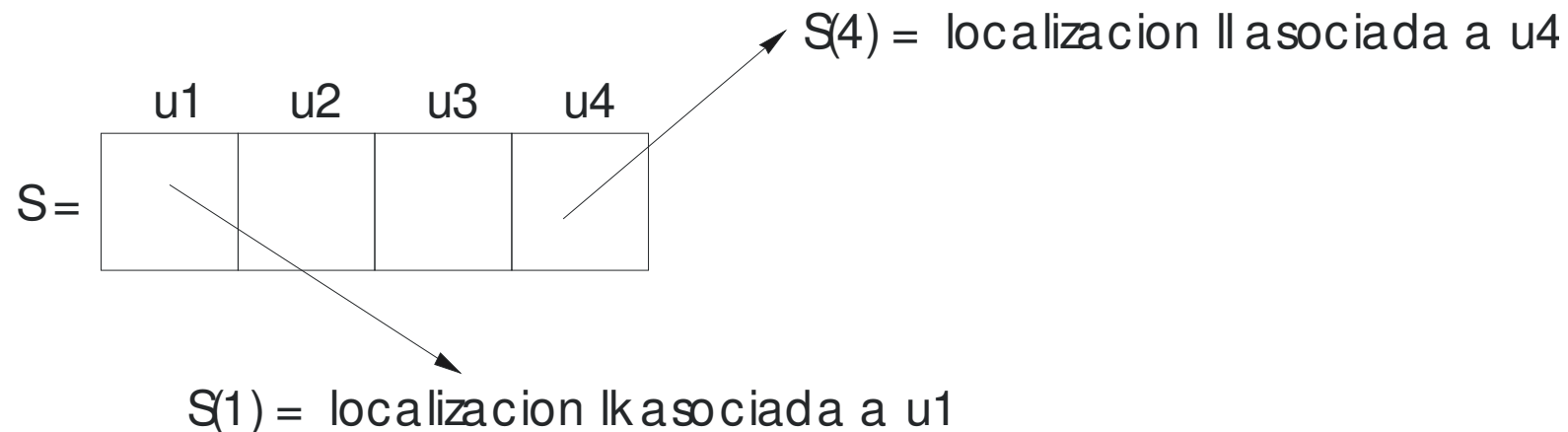
- Así, la solución  $S = \{3, 4, 1, 2\}$  representa la siguiente distribución de asignaciones:

$u1 \leftrightarrow I3$

$u2 \leftrightarrow I4$

$u3 \leftrightarrow I1$

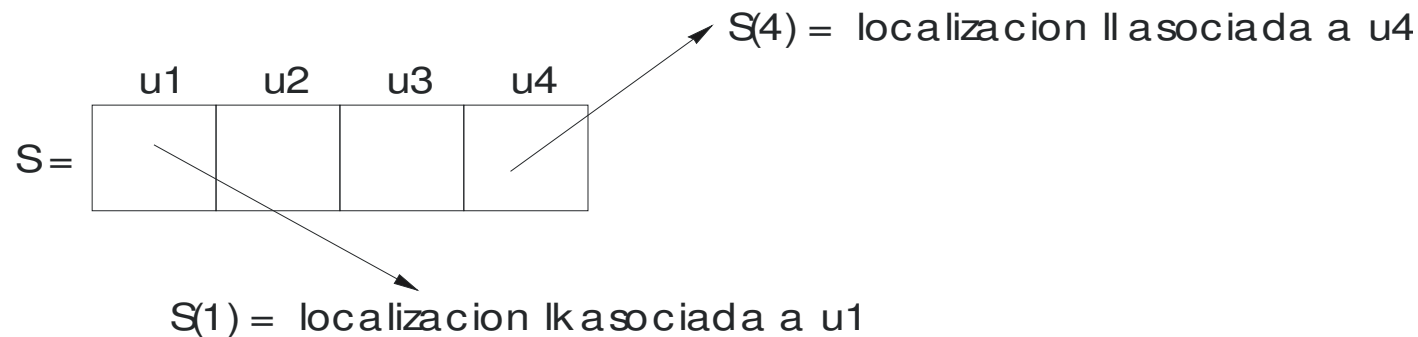
$u4 \leftrightarrow I2$



# Representación y Ejemplo: Diseño de un Hospital

- cuyo costo  $C(S)$  es:

$$\begin{aligned}
 & f_{12} \cdot d_{34} + f_{13} \cdot d_{31} + f_{14} \cdot d_{32} && 3 \cdot 7 + 8 \cdot 6 + 3 \cdot 6 + \\
 + & f_{21} \cdot d_{43} + f_{23} \cdot d_{41} + f_{24} \cdot d_{42} && 3 \cdot 7 + 2 \cdot 4 + 4 \cdot 8 + \\
 + & f_{31} \cdot d_{13} + f_{32} \cdot d_{14} + f_{34} \cdot d_{12} && 8 \cdot 6 + 2 \cdot 4 + 5 \cdot 12 + \\
 + & f_{41} \cdot d_{23} + f_{42} \cdot d_{24} + f_{43} \cdot d_{21} && 3 \cdot 6 + 4 \cdot 8 + 5 \cdot 12 && = 374
 \end{aligned}$$



- Cada asignación unidad-localización influye globalmente en la red, es decir, en todas las transferencias que se efectúan desde la unidad en cuestión

# Otras Aplicaciones

---

- **Diseño óptimo de teclados** para distintos idiomas en función de la frecuencia de pares de letras. Unid: letras; loc: teclas; flujo: frecuencia de pares de letras; dist: distancia entre las teclas en el teclado

Burkard, R.E. and J. Offerman (1977). Entwurf von Schreibmaschinentastaturen mittels quadratischer Zuordnungsprobleme. Z. Operations Research 21, B121-B123

- **Cableado óptimo de placas madre**: localización de componentes para reducir el cableado. Unid: componentes; loc: posición en la placa; flujo: nº de cables conectando componentes; dist: distancia entre locs.

Brixius, N.W. and K.M. Anstreicher (2001). The Steinberg Wiring Problem. In: The Sharpest Cut, The Impact of M. Padberg and His Work, M. Grötschel, ed., SIAM, 2004, 293-307

- **Asignación óptima de directores a oficinas**: Unid: directores; loc: oficinas; flujo: frecuencia de interacciones entre directores; dist: distancia entre oficinas

Hanan, M. and J.M. Kurtzberg (1972). A Review of the Placement and Quadratic Assignment Problems, SIAM Review 14, 324-342

# Otras Aplicaciones

---

- **Diseño de turbinas:** Localización de las aspas de la turbina (con masa ligeramente distinta por la fabricación) de modo que el centro de gravedad coincida con el eje del motor

J. Mosevich (1986). Balancing hydraulic turbine runners--A discrete combinatorial optimization problem, *European Journal of Operational Research* 26(2), 202-204

- **Diseño óptimo de campus**

J.W. Dickey, J.W. Hopkins (1972). Campus building arrangement using TOPAZ, *Transportation Research* 6, 59-68

- **Diseño de parques forestales**

Bos, J. (1993). A quadratic assignment problem solved by simulated annealing. *Journal of Environmental Management*, 37(2), 127-145

- **Diseño de líneas de producción**

Geoffrion, A.M., and G.W. Graves (1976). Scheduling Parallel Production Lines with Changeover Costs: Practical Applications of a Quadratic Assignment/LP Approach. *Operations Research* 24, 595-610

# Solución Greedy

---

- La complejidad del problema ha provocado que se hayan aplicado muchos algoritmos aproximados para su resolución
- Analizando la función objetivo podemos determinar que una buena fórmula heurística para resolver el problema es:

*Asociar unidades de gran flujo con localizaciones céntricas en la red y viceversa*

# Solución Greedy

---

- Podemos construir un algoritmo greedy usando esta heurística mediante dos vectores, el potencial de flujo y el de distancia:

$$\hat{f}_i = \sum_{j=1}^n f_{ij} ; \quad i = 1, \dots, n$$

$$\hat{d}_k = \sum_{l=1}^n d_{kl} ; \quad k = 1, \dots, n$$

- Cuanto mayor sea  $\hat{f}_i$ , más importante es la unidad en el intercambio de flujos y cuanto menor sea  $\hat{d}_k$ , más céntrica es la localización. Por tanto:

*el algoritmo irá seleccionando la unidad  $i$  libre con mayor  $\hat{f}_i$   
y le asignará la localización  $k$  libre con menor  $\hat{d}_k$*

# Solución Greedy

---

## Algoritmo Greedy-QAP

1. Calcular los potenciales  $\hat{f}_i$  y  $\hat{d}_k$ .
2.  $S \leftarrow \emptyset$ .
3. Repetir para  $x=1$  hasta  $n$  (asignaciones  $l$  a  $n$ ):
  - 3.1. Escoger la unidad  $u_i$  no asignada aún ( $u_i \notin S$ ) con mayor valor de  $\hat{f}_i$ .
  - 3.2. Escoger la localización  $l_k$  no asignada aún ( $l_k \notin S$ ) con menor valor de  $\hat{d}_k$ .
  - 3.3.  $a_x = (u_i, l_k)$ .  $S \leftarrow S \cup a_x$ .
4. Calcular el costo de  $S$ ,  $C(S)$ . Devolver  $S$  y  $C(S)$ .

# Solución Greedy

---

- Si aplicamos el algoritmo greedy propuesto al ejemplo del hospital obtenemos los siguientes resultados:

$$\hat{f} = \begin{bmatrix} 14 \\ 9 \\ 15 \\ 12 \end{bmatrix}$$

$$\hat{d} = \begin{bmatrix} 22 \\ 26 \\ 17 \\ 19 \end{bmatrix}$$

$$S = \{(u_3, l_3), (u_1, l_4), (u_4, l_1), (u_2, l_2)\}$$

$$\begin{aligned} C(S) = & f_{12} \cdot d_{42} + f_{13} \cdot d_{43} + f_{14} \cdot d_{41} & 3 \cdot 8 + 8 \cdot 7 + 3 \cdot 4 + \\ & + f_{21} \cdot d_{24} + f_{23} \cdot d_{23} + f_{24} \cdot d_{21} & 3 \cdot 8 + 2 \cdot 6 + 4 \cdot 12 + \\ & + f_{31} \cdot d_{34} + f_{32} \cdot d_{32} + f_{34} \cdot d_{31} & 8 \cdot 7 + 2 \cdot 6 + 5 \cdot 6 + \\ & + f_{41} \cdot d_{14} + f_{42} \cdot d_{12} + f_{43} \cdot d_{13} = & 3 \cdot 4 + 4 \cdot 12 + 5 \cdot 6 = 364 \end{aligned}$$



# Búsquedas por Trayectorias Simples

---

- **Representación:** Problema de asignación: una permutación  $\pi = [\pi(1), \dots, \pi(n)]$  en el que las posiciones del vector  $i=1, \dots, n$  representan las unidades y los valores  $\pi(1), \dots, \pi(n)$  contenidos en ellas las localizaciones. Permite verificar las restricciones
- **Operador de vecino de intercambio y su entorno:** El entorno de una solución  $\pi$  está formado por las soluciones accesibles desde ella a través de un movimiento de intercambio

Dada una solución (asignación de unidades a localizaciones) se escogen dos unidades distintas y se intercambia la localización asignada a cada una de ellas ( $Int(\pi, i, j)$ ):

$$\pi = [\pi(1), \dots, \pi(i), \dots, \pi(j), \dots, \pi(n)]$$

$$\pi' = [\pi(1), \dots, \pi(j), \dots, \pi(i), \dots, \pi(n)]$$

# Búsquedas por Trayectorias Simples

---

- $Int(\pi, i, j)$  verifica las restricciones, si la solución original  $\pi$  es factible siempre genera una solución vecina  $\pi'$  factible
- Su aplicación provoca que el tamaño del entorno sea:

$$|E(\pi)| = \frac{n \cdot (n-1)}{2}$$

- Las instancias del QAP no suelen ser demasiado grandes y el **cálculo factorizado del coste** de una solución se realiza de forma eficiente ( $O(n)$ ), permitiendo explorar el entorno completo

Aún así, dicha exploración requería  $O(n^3)$  por lo que es recomendable utilizar una estrategia avanzada, considerando una modalidad de lista de candidatos y seleccionado primero los movimientos más prometedores

# Búsqueda Local para el QAP

---

Stützle, Iterated local search for the quadratic assignment problem, European Journal of Operational Research 174 (2006) 1519–1539

- Algoritmo de **búsqueda local del primer mejor**: en cuanto se genera una solución vecina que mejora a la actual, se aplica el movimiento y se pasa a la siguiente iteración
  - Se detiene la búsqueda cuando se ha explorado el vecindario completo sin obtener mejora
- Se considera una **factorización** para calcular el coste de  $\pi'$  a partir del de  $\pi$  considerando sólo los cambios realizados por el movimiento de intercambio
- Se usa la técnica ***don't look bits*** para la lista de candidatos
  - Se emplea un vector binario de tamaño  $n$  que asocia un bit a cada unidad
  - Si dicho bit está activado en la iteración actual, no se considera ningún movimiento “que arranque” de la unidad en cuestión

# BL-QAP: Factorización del Movimiento de Intercambio

- Sea  $C(\pi)$  el coste de la solución original  $\pi$ . Para generar  $\pi'$ , el operador de vecino  $Int(\pi, r, s)$  escoge dos unidades  $r$  y  $s$  e intercambia sus localizaciones  $\pi(r)$  y  $\pi(s)$ :  
 $\pi = [\pi(1), \dots, \pi(r), \dots, \pi(s), \dots, \pi(n)]$

$$t \leftarrow \pi(r) ; \pi(r) \leftarrow \pi(s) ; \pi(s) \leftarrow t \quad \Rightarrow \quad \pi' = [\pi(1), \dots, \pi(s), \dots, \pi(r), \dots, \pi(n)]$$

- Por lo tanto, quedan afectados  $2 \cdot n$  sumandos, los relacionados con las dos **viejas** y las dos **nuevas** localizaciones de las dos unidades alteradas
- El coste del movimiento (la **diferencia de costes entre las dos soluciones**)  $\Delta C(\pi, r, s) = C(\pi') - C(\pi)$  se puede factorizar como:

$$\sum_{k=1, k \neq r, s}^n \left[ f_{rk} \cdot (d_{\pi(s)\pi(k)} - d_{\pi(r)\pi(k)}) + f_{sk} \cdot (d_{\pi(r)\pi(k)} - d_{\pi(s)\pi(k)}) + \right. \\ \left. f_{kr} \cdot (d_{\pi(k)\pi(s)} - d_{\pi(k)\pi(r)}) + f_{ks} \cdot (d_{\pi(k)\pi(r)} - d_{\pi(k)\pi(s)}) \right]$$

**nuevas**

**viejas**

# BL-QAP: Factorización del Movimiento de Intercambio

---

- Si  $\Delta C(\pi, r, s)$  es negativo ( $\Delta C(\pi, r, s) < 0$ ), la solución vecina  $\pi'$  es mejor que la actual  $\pi$  (el QAP es un problema de minimización) y se acepta. Si no, se descarta y se genera otro vecino
- El pseudocódigo de la BL del Primer Mejor del Tema 2 de Teoría quedaría:

## Repetir

$\pi' \leftarrow \text{GENERA\_VECINO}(\pi_{\text{act}});$

**Hasta** ( $\Delta C(\pi, r, s) < 0$ ) **O**  
(se ha generado  $E(\pi_{\text{act}})$  al completo)

- El coste  $C(\pi')$  de la nueva solución vecina es:  $C(\pi') = C(\pi) + \Delta C(\pi)$ .  
Sólo es necesario calcularlo para la solución vecina aceptada

# BL-QAP: Definición de la Lista de Candidatos: *Don't Look Bits*

---

- Técnica que permite focalizar la BL en una zona del espacio de búsqueda en la que potencialmente puede ocurrir algo
- Reduce significativamente el tiempo de ejecución con una reducción muy pequeña de la eficacia de la BL
- Sólo es aplicable con la BL del primer mejor
- En la primera iteración, todos los bits están a 0, es decir, todas las unidades están activadas en un bucle externo y todos sus movimientos pueden ser considerados para explorar el entorno:
  - Si tras probar todos los movimientos asociados a esa unidad, ninguno provoca una mejora, se pone su bit a 1 para desactivarla en el futuro
  - Si una unidad está implicada en un movimiento que genera una solución vecina con mejor coste, se pone su bit a 0 para reactivarla

# BL-QAP: Definición de la Lista de Candidatos: *Don't Look Bits*

procedure iterative improvement

  for  $i = 1$  to  $n$  do

    if  $dlb[i] = 0$  then

$improve\_flag \leftarrow false$

      for  $j = 1$  to  $n$  do

        CheckMove( $i, j$ )

        if move improves then

          ApplyMove( $i, j$ );  $dlb[i] \leftarrow 0, dlb[j] \leftarrow 0$

$improve\_flag \leftarrow true$

      endfor

      if  $improve\_flag = false$  then  $dlb[i] \leftarrow 1$

    end

end iterative improvement

**IMPORTANTE:** Este es el bucle interno de la BL, el de exploración del vecindario de la solución actual. Sea cual sea la BL usada, siempre habrá un bucle externo que repetirá el proceso mientras se produzca mejora

# Enfriamiento Simulado para el QAP

---

- **Representación**: permutación de asignación, igual que en la BL
- **Operador de generación de vecinos**: intercambio, igual que en la BL
- **Exploración del vecindario**: En cada iteración del bucle interno se genera una única solución vecina, de forma aleatoria, y se compara con la actual. **No se usa don't look bits**
- **Esquema de enfriamiento**: esquema de Cauchy modificado
- **Condición de enfriamiento  $L(T)$** : cuando se genere un número máximo de soluciones vecinas, *máx\_vecinos*, o cuando se acepte un número máximo de los vecinos generados, *máx\_éxitos*
- **Condición de parada**: cuando se alcance un número máximo de iteraciones o cuando el número de éxitos en el enfriamiento actual sea 0



# Búsqueda Tabú para el QAP

---

Skorin-Kapov, Tabu Search Applied to the Quadratic Assignment Problem,  
ORSA Journal on Computing 2:1 (1990) 33-40

- **Representación:** permutación de asignación, igual que en la BL
- **Operador de generación de vecinos:** intercambio, igual que en la BL
- **Lista de valores de atributos y posiciones tabú:** Cuando se acepta una nueva solución, se almacena información sobre ella en la lista tabú

Se guardan tanto las **localizaciones intercambiadas como las posiciones que ocupan** (es decir, las unidades a las pasan a estar asignadas):  $(i, j, \text{Pos}(i), \text{Pos}(j))$

No se permiten futuros intercambios que den lugar a que la localización  $i$  vuelva a ocupar la posición (estar asignada a la unidad)  $\text{Pos}(i)$  Y/O la localización  $j$  la posición  $\text{Pos}(j)$

# Búsqueda Tabú para el QAP

---

La lista tabú es un vector de registros con esos cuatro valores

Ejemplo:

$$\begin{array}{cccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ \pi_{\text{act}} & = & (\mathbf{1} & \mathbf{2} & \underline{\mathbf{4}} & \mathbf{3} & \mathbf{8} & \underline{\mathbf{5}} & \mathbf{7} & \mathbf{6}) \\ \pi' & = & (\mathbf{1} & \mathbf{2} & \underline{\mathbf{5}} & \mathbf{3} & \mathbf{8} & \underline{\mathbf{4}} & \mathbf{7} & \mathbf{6}) \end{array}$$

$LT \leftarrow LT + \{(4,5,6,3)\}$  (si está llena, eliminar el registro más antiguo)

Para decidir si un movimiento es tabú activo, hay que tener en cuenta que  $(i, j, \text{Pos}(i), \text{Pos}(j)) = (j, i, \text{Pos}(j), \text{Pos}(i))$

- **Estrategia de selección de vecino:** Examinar un número de soluciones vecinas generadas aleatoriamente y escoger la mejor que verifique los criterios tabú o que no los verifique y supere el criterio de aspiración

# Búsqueda Tabú para el QAP

---

- **Criterio de aspiración:** Tener menor coste que la mejor solución obtenida hasta el momento
- **Estrategias de reinicialización:**
  - Generar una solución aleatoria y continuar el algoritmo a partir de ella (**diversificación**)
  - Volver a la mejor solución obtenida hasta el momento y continuar el algoritmo a partir de ella (**intensificación**)
  - Generar una nueva solución poco frecuente a partir de la memoria a largo plazo y continuar el algoritmo a partir de ella (**diversificación controlada**)

# Búsqueda Tabú para el QAP

---

- **Memoria a largo plazo:** Se usa una matriz *frec* de tamaño  $n \times n$  que almacena el número de veces que se ha producido cada asignación unidad-localización en las soluciones aceptadas durante la búsqueda
- Para reinicializar con diversidad, se van realizando una a una las asignaciones que se dieron con menor frecuencia en el pasado:
  - Hacer  $U \leftarrow \{1, \dots, n\}$ ,  $L \leftarrow \{1, \dots, n\}$
  - Repetir  $n$  veces
    - Sea  $(u, l) \in U \times L$  tal que  $\text{frec}[u, l] \leq \text{frec}[i, j]$ ,  $\forall i \in U$ ,  $\forall j \in L$
    - Hacer  $S[u] \leftarrow l$
    - Hacer  $U \leftarrow U - \{u\}$  y  $L \leftarrow L - \{l\}$

# Búsqueda Tabú para el QAP

---

- Se realiza una reinicialización cuando transcurran un cierto número de iteraciones sin mejorar la mejor solución obtenida hasta el momento en la ejecución del algoritmo (mejor global)
- Para realizar una reinicialización, se escoge una de las tres opciones según la siguiente distribución de probabilidades:
  - 0,25: Nueva solución actual aleatoria
  - 0,25: Volver a la mejor solución generada
  - 0,50: Solución poco frecuente a partir de la memoria a largo plazo
- Independientemente del tipo de reinicialización realizada, se modifica el tamaño de la lista tabú para provocar un cambio más efectivo en el comportamiento del algoritmo. El tamaño de cada lista se incrementa o decrementa en un 50% (modificación sobre el tamaño actual)
- Tras la reinicialización, la lista de corto plazo se inicializa a cero, mientras que la de largo plazo no se modifica

# La Biblioteca QAPLIB

---

- La QAPLIB es una biblioteca que contiene distintas instancias del QAP llevando un registro de las mejores soluciones obtenidas hasta el momento para las mismas y de las cotas teóricas de la calidad de la mejor solución que se puede obtener
- Es accesible en la Web en las direcciones siguientes:  
  
<http://www.opt.math.tu-graz.ac.at/qaplib> (hasta Feb. 2002)  
  
<http://www.seas.upenn.edu/qaplib/> (hasta la actualidad)
- En dicha dirección pueden encontrarse tanto los datos como las soluciones de distintas instancias del problema, relacionadas con diferentes aplicaciones

# La Biblioteca QAPLIB

---

- El formato de los ficheros de datos es:

$$\begin{array}{c} n \\ \textcircled{A} \\ \textcircled{B} \end{array} \min_{S \in \Pi_N} \left( \sum_{i=1}^n \sum_{j=1}^n a_{ij} \cdot b_{S(i)S(j)} \right)$$

donde  $n$  es el tamaño del problema y  $A$  y  $B$  son, respectivamente las matrices de flujo y distancia, de acuerdo a nuestra codificación de las soluciones del problema

- El formato de los ficheros de soluciones (óptimas o mejores conocidas) es:

$$\begin{array}{c} n \quad sol \\ p \end{array}$$

donde  $n$  es el tamaño del problema,  $sol$  es el coste de la solución y  $p$  es la permutación correspondiente

# La Biblioteca QAPLIB

- La siguiente tabla es un ejemplo de la información que proporciona la QAPLIB:

name	n	feas.sol.	permutation/bound	gap
-----				
<a href="#">Tai12a</a>	12	<a href="#">224416</a> (OPT)	(8,1,6,2,11,10,3,5,9,7,12,4)	
<a href="#">Tai12b</a>	12	<a href="#">39464925</a> (OPT)	(9,4,6,3,11,7,12,2,8,10,1,5)	
<a href="#">Tai15a</a>	15	<a href="#">388214</a> (OPT)	(5,10,4,13,2,9,1,11,12,14,7,15,3,8,6)	
<a href="#">Tai15b</a>	15	<a href="#">51765268</a> (OPT)	(1,9,4,6,8,15,7,11,3,5,2,14,13,12,10)	
<a href="#">Tai17a</a>	17	<a href="#">491812</a> (OPT)	(12,2,6,7,4,8,14,5,11,3,16,13,17,9,1,10,15)	
<a href="#">Tai20a</a>	20	<a href="#">703482</a> (OPT)	(10,9,12,20,19,3,14,6,17,11,5,7,15,16,18,2,4,8,13,1)	
* <a href="#">Tai20b</a>	20	<a href="#">122455319</a> (OPT)	(8,16,14,17,4,11,3,19,7,9,1,15,6,13,10,2,5,20,18,12)	
* <a href="#">Tai25a</a>	25	<a href="#">1167256</a> (OPT)	(9,4,6,11,5,1,15,10,14,3,17,12,19,18,23,8,21,2,22,7,16,20,24,25,13)	
* <a href="#">Tai25b</a>	25	<a href="#">344355646</a> (OPT)	(4,15,10,9,13,5,25,19,7,3,17,6,18,20,16,2,22,23,8,11,21,24,14,12,1)	
<a href="#">Tai30a</a>	30	<a href="#">1818146</a> (Ro-TS)	1706855 (L&P)	6.12 %
* <a href="#">Tai30b</a>	30	<a href="#">637117113</a> (OPT)	(4 8 11 15 17 20 21 5 14 30 2 13 6 29 10 26 27 24 28 22 12 9 7 23 19)	
<a href="#">Tai35a</a>	35	<a href="#">2422002</a> (Ro-TS)	2216627 (L&P)	8,48 %



# Metaheurísticas

## Seminario 2.a. Problemas de optimización con técnicas basadas en trayectorias simples

---

1. Problema de Asignación Cuadrática (QAP)
2. Problema de la Selección de Características
  - Definición del Problema de Clasificación
  - Ejemplo de Clasificador: k-NN
  - Definición del Problema de Selección de Características
  - Representaciones
  - Solución Greedy
  - Búsquedas por Trayectorias Simples
  - Bases de Datos a Utilizar

# Definición del Problema de Clasificación

---

Disponemos de una muestra de objetos ya clasificados  $w_1, \dots, w_n$ , representados en función de sus valores en una serie de atributos:

$w_i$  tiene asociado el vector  $(x_1(w_i), \dots, x_n(w_i))$

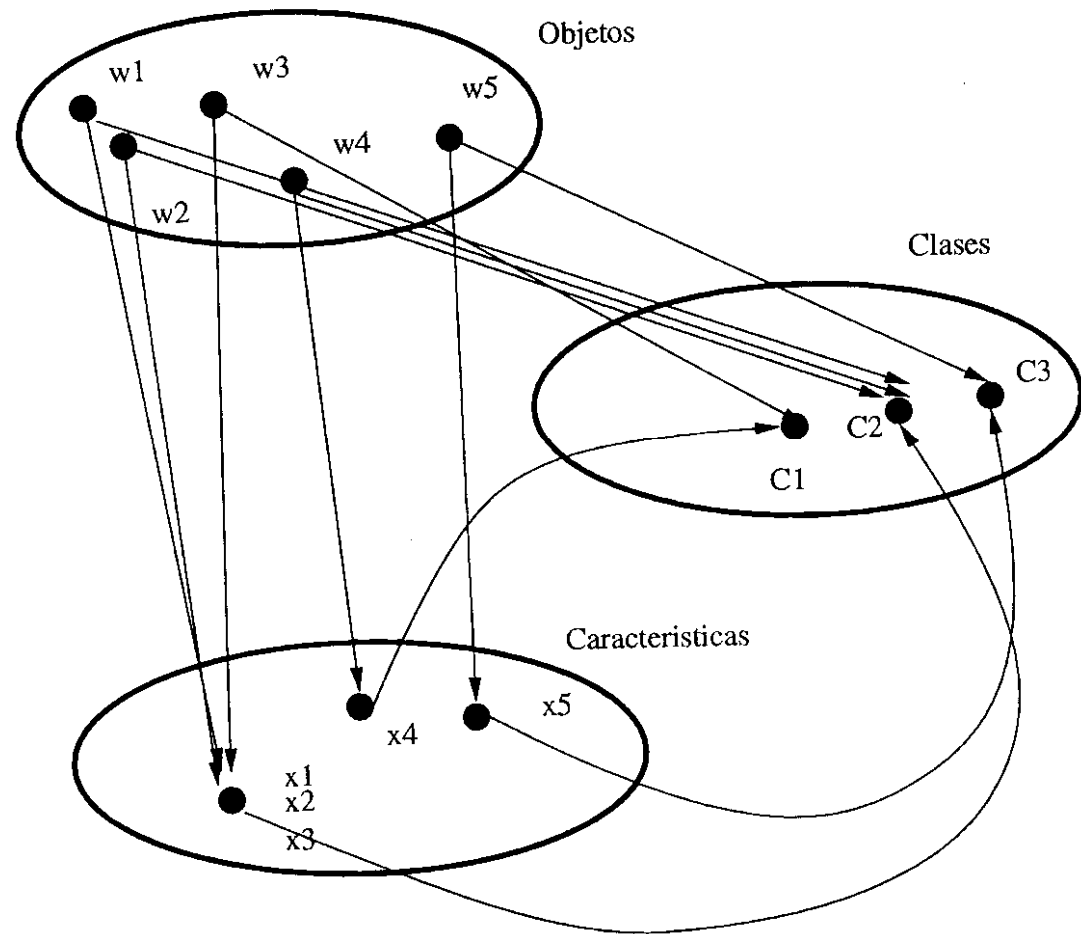
Cada objeto pertenece a una de las clases existentes  $\{C_1, \dots, C_M\}$

**OBJETIVO: Obtener un sistema que permita clasificar dichos objetos de modo automático**

$$\begin{array}{lll} w_1 = (x_1(w_1), \dots, x_n(w_1)) & \rightarrow & C_{i_1} \\ \vdots & & \vdots \\ w_k = (x_1(w_k), \dots, x_n(w_k)) & \rightarrow & C_{i_k} \end{array} \quad i_j \in \{1, \dots, M\}, \quad j \in \{1, \dots, k\}$$

# Definición del Problema de Clasificación

El problema fundamental de la clasificación está directamente relacionado con la separabilidad de las clases



# Definición del Problema de Clasificación

---

## Concepto de aprendizaje supervisado en clasificación

- Se conocen las clases existentes en el problema
- Se conoce la clase concreta a la que pertenece cada objeto del conjunto de datos

Existen una gran cantidad de técnicas para el aprendizaje supervisado de Sistemas de Clasificación:

- Técnicas estadísticas:  $k$  vecinos más cercanos, discriminadores bayesianos, etc...
- Árboles de clasificación, Sistemas basados en reglas, Redes Neuronales, Máquinas de Soporte Vectorial, ...

# Definición del Problema de Clasificación

---

## **Ejemplo:** *Diseño de un Clasificador para la flor del Iris*

- *Problema simple muy conocido: clasificación de lirios*
- *Tres clases de lirios: setosa, versicolor y virgínica*
- *Cuatro atributos: longitud y anchura de pétalo y sépalo, respectivamente*
- *150 ejemplos, 50 de cada clase*
- *Disponible en <http://www.ics.uci.edu/~mlearn/MLRepository.html>*



***setosa***



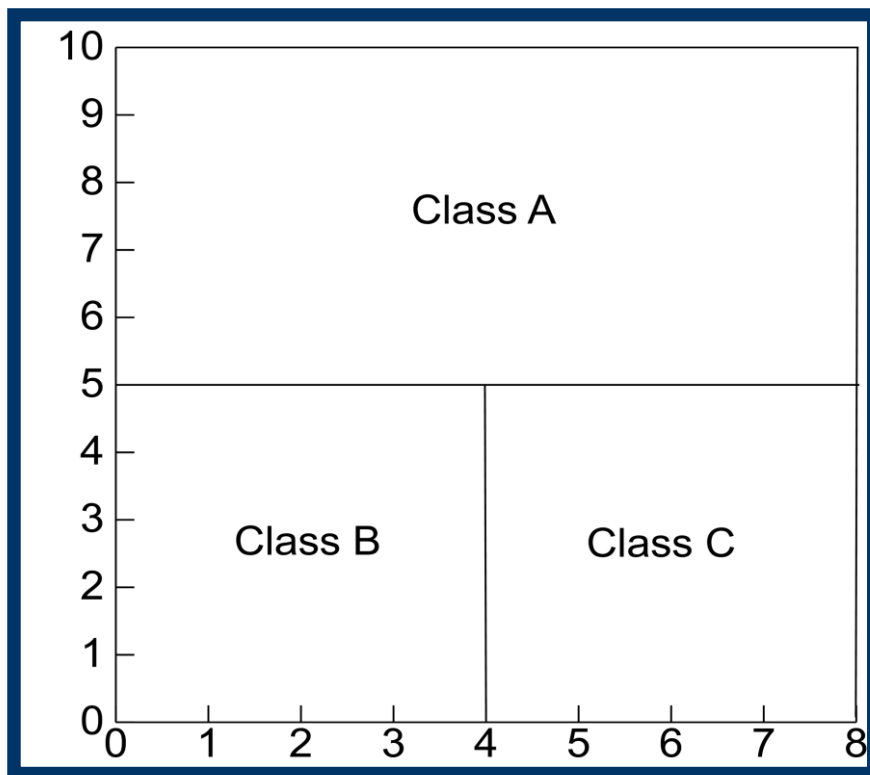
***versicolor***



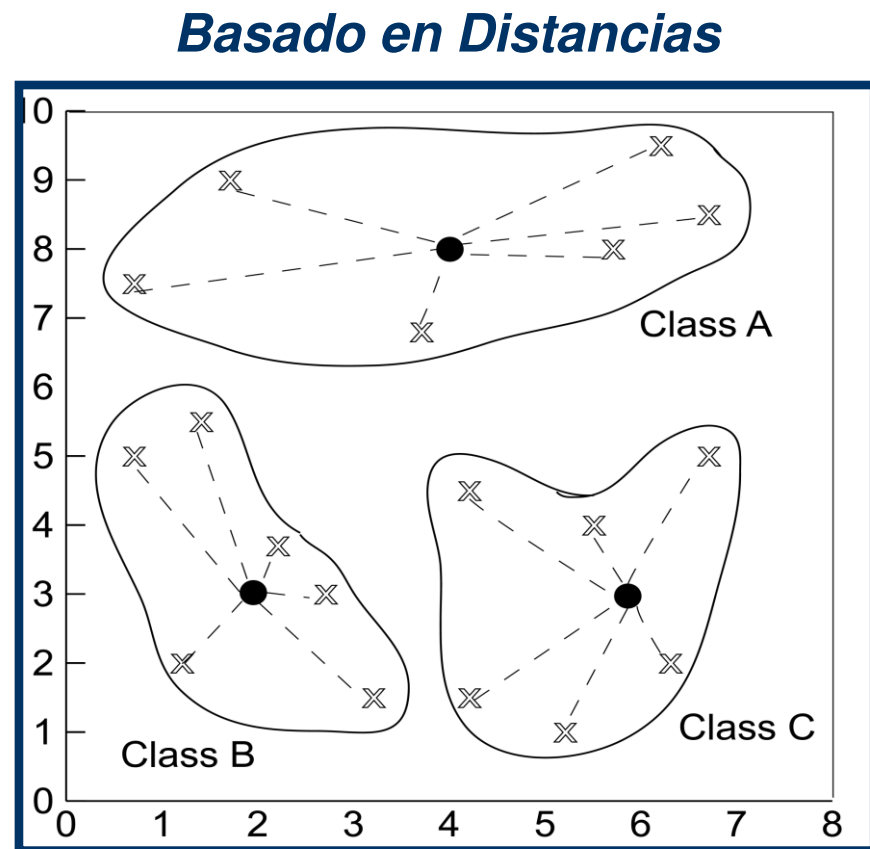
***virgínica***

# Definición del Problema de Clasificación

## Ejemplos de clasificación sobre clases definidas: Basada en particiones y en distancias



*Basado en Particiones*

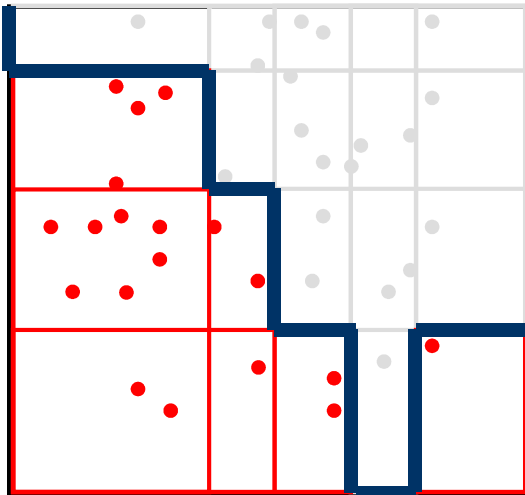


# Definición del Problema de Clasificación

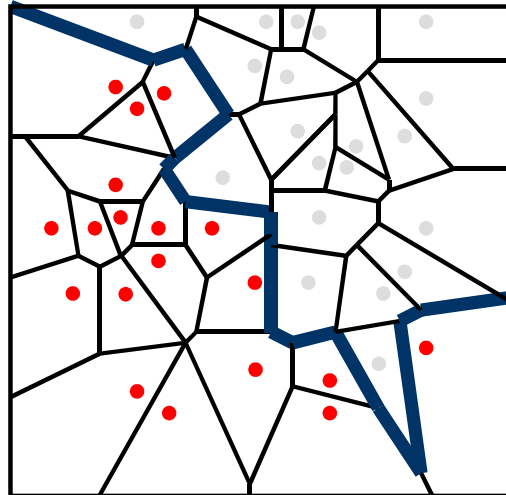
---

## Ejemplos de clasificación sobre clases definidas

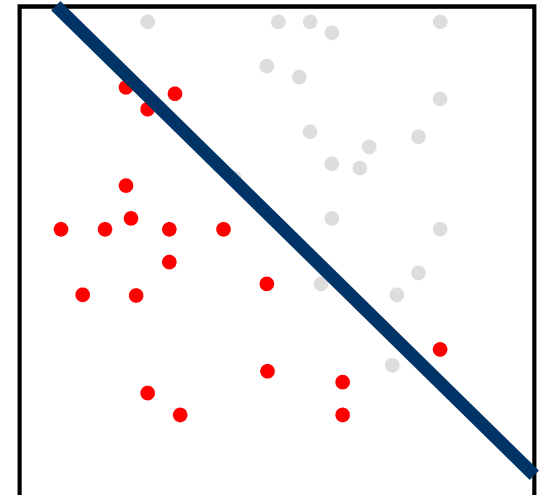
- *Reglas intervalares*



- *Basado en distancias*



- *Clasificador lineal*



# Definición del Problema de Clasificación

---

Para diseñar un clasificador, son necesarias dos tareas:

## **Aprendizaje y Validación**

El conjunto de ejemplos se divide en dos subconjuntos:

- **Entrenamiento:** Utilizado para aprender el clasificador
- **Prueba:** Se usa para validarlo. Se calcula el porcentaje de clasificación sobre los ejemplos de este conjunto (desconocidos en la tarea de aprendizaje) para conocer su poder de generalización

Para mayor seguridad, se suele hacer varias particiones entrenamiento-prueba

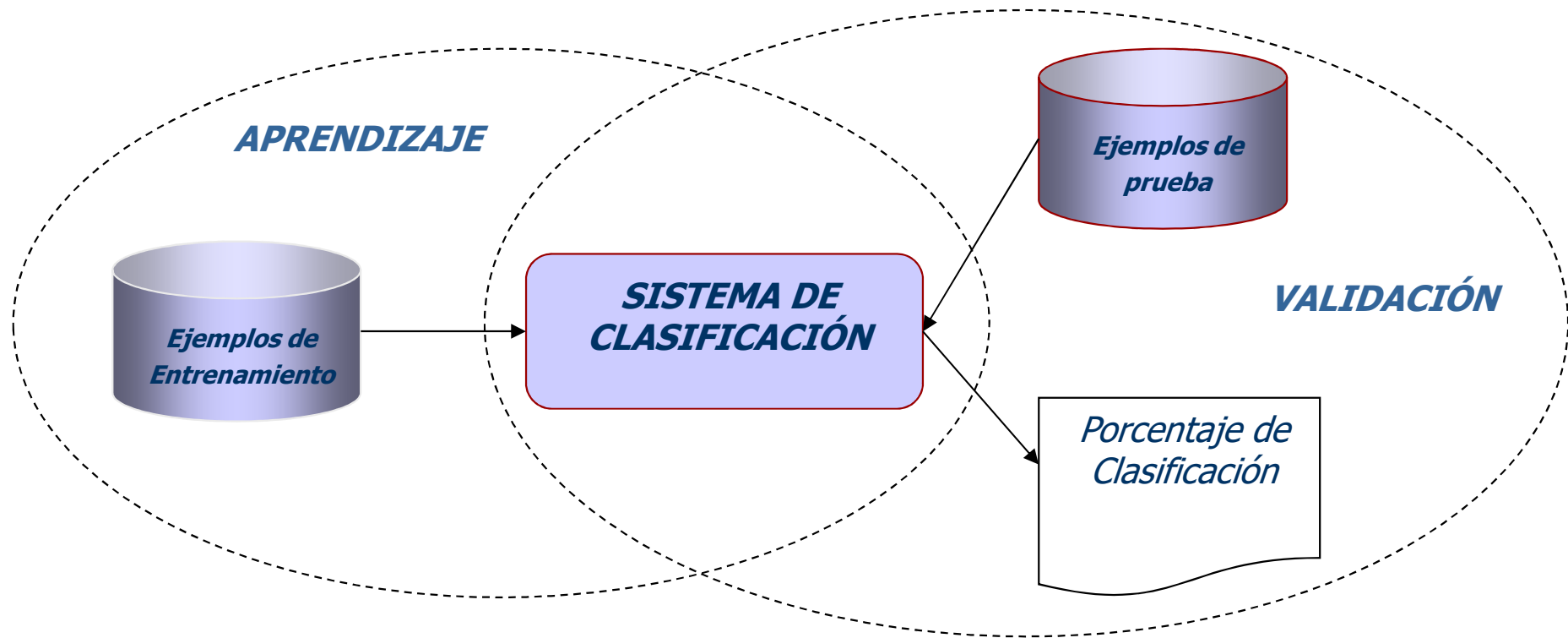
Para cada una, se diseña un clasificador distinto usando los ejemplos de entrenamiento y se valida con los de prueba



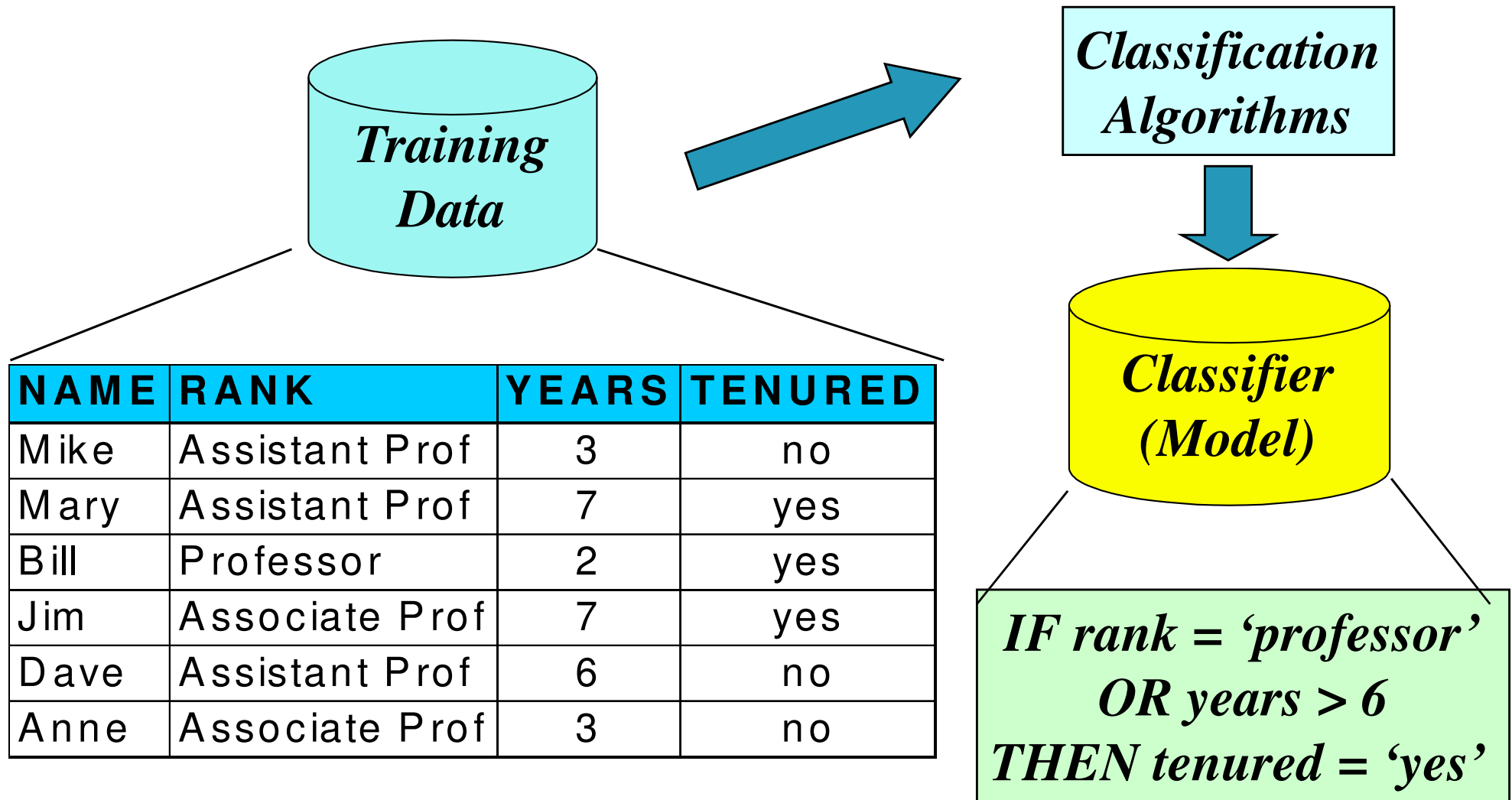
# Definición del Problema de Clasificación

---

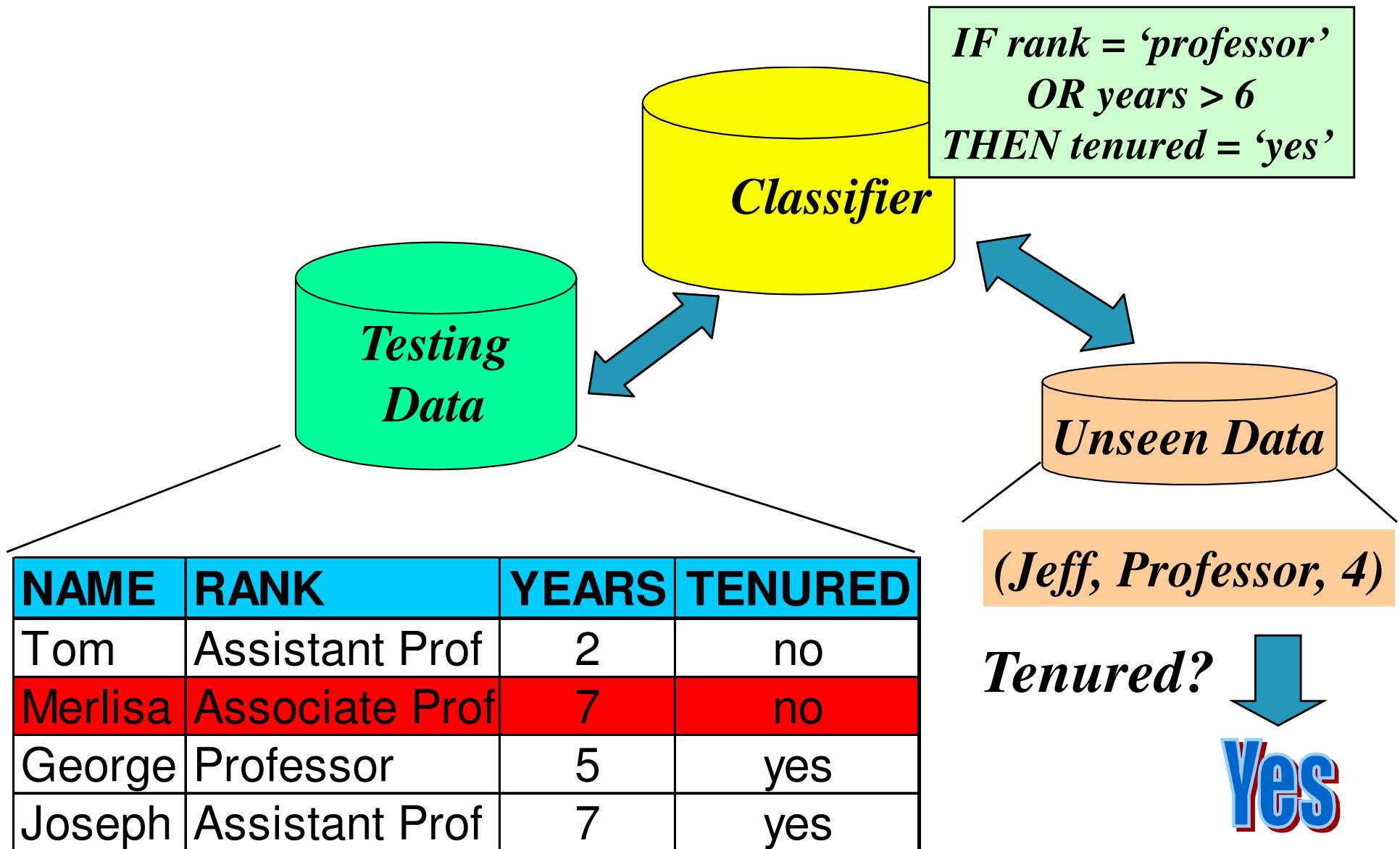
## Esquema de aprendizaje en clasificación



# Definición del Problema de Clasificación



# Definición del Problema de Clasificación



# Definición del Problema de Clasificación

---

Nosotros aplicaremos la técnica de validación llamada **5x2-cross validation** (validación cruzada)

- Usaremos 5 particiones distintas de los datos al 50%
- Aprenderemos un clasificador con la mitad de los datos y lo validaremos con la otra mitad. Luego repetiremos el proceso a la inversa
- Por tanto, de cada partición obtendremos dos valores de porcentaje de clasificación en el conjunto de prueba
- La calidad del método de clasificación se medirá con un único valor, correspondiente a la media de los 10 porcentajes de clasificación (2 de cada partición)

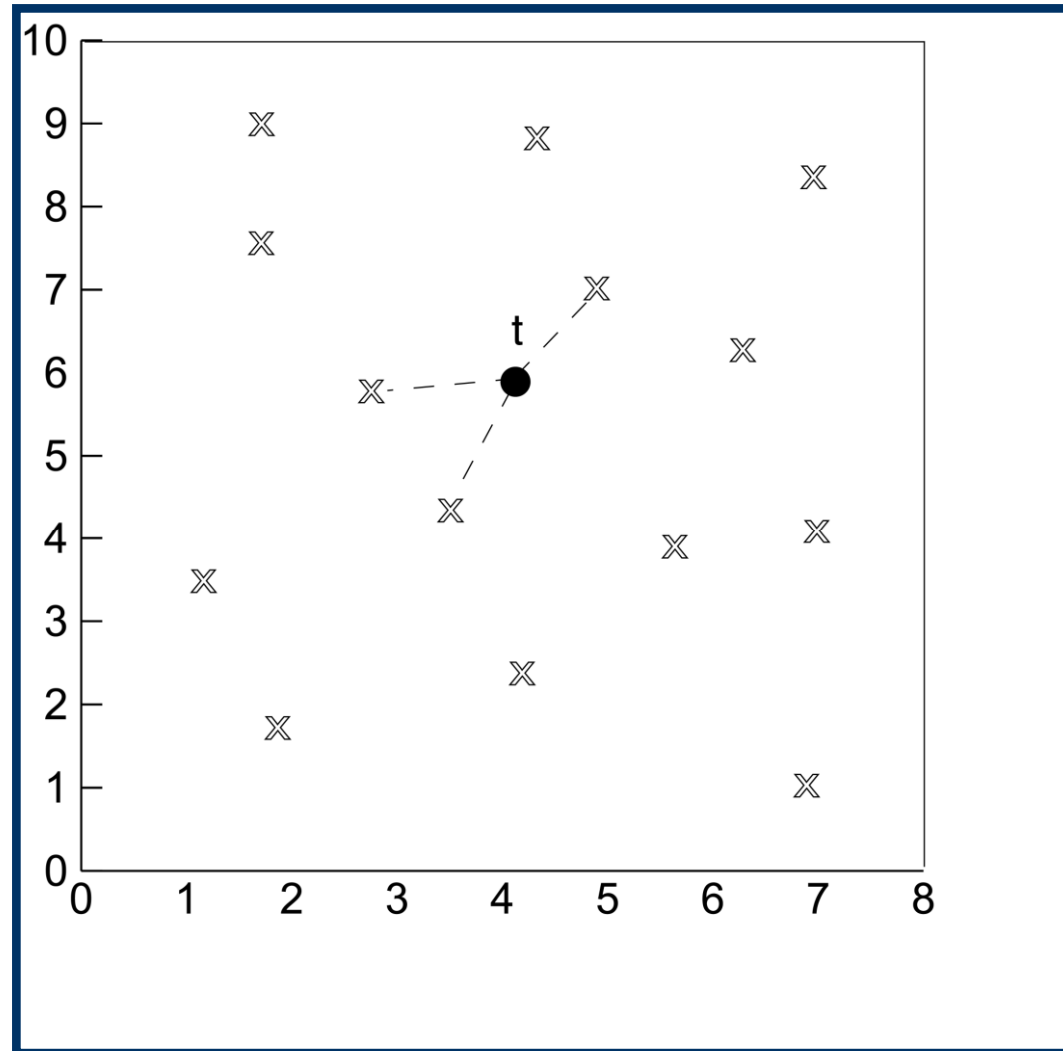
# Ejemplo de Clasificador: k-NN

---

*El k-NN (k vecinos más cercanos) es uno de los clasificadores más utilizados por su simplicidad*

- i. El proceso de aprendizaje de este clasificador consiste en almacenar una tabla con los ejemplos disponibles, junto a la clase asociada a cada uno de ellos
- ii. Dado un nuevo ejemplo a clasificar, se calcula su distancia (usaremos la Euclídea) a los  $n$  ejemplos existentes en la tabla y se escogen los  $k$  más cercanos
- iii. El nuevo ejemplo se clasifica según la clase mayoritaria de esos  $k$  ejemplos más cercanos
- iv. El caso más simple es cuando  $k = 1$  (1-NN)

# Ejemplo de Clasificador: k-NN



**$k = 3$**

# Ejemplo de Clasificador: k-NN

Dado el siguiente conjunto con 4 instancias, 3 atributos y 2 clases:

$x_1$ :	0.4	0.8	0.2	positiva
$x_2$ :	0.2	0.7	0.9	positiva
$x_3$ :	0.9	0.8	0.9	negativa
$x_4$ :	0.8	0.1	0.0	negativa

Calculamos la distancia del ejemplo con todos los de la tabla:

*Queremos clasificar con 1-NN el ejemplo:*  
 $x_q$ : 0.7 0.2 0.1

$$d(x_1, x_q) = \sqrt{(0.4 - 0.7)^2 + (0.8 - 0.2)^2 + (0.2 - 0.1)^2} = 0.678$$

$$d(x_2, x_q) = \sqrt{(0.2 - 0.7)^2 + (0.7 - 0.2)^2 + (0.9 - 0.1)^2} = 1.068$$

$$d(x_3, x_q) = \sqrt{(0.9 - 0.7)^2 + (0.8 - 0.2)^2 + (0.9 - 0.1)^2} = 1.020$$

$$d(x_4, x_q) = \sqrt{(0.8 - 0.7)^2 + (0.1 - 0.2)^2 + (0.0 - 0.1)^2} = 0.173$$

Por tanto, el ejemplo se clasificará con respecto a la clase negativa

**IMPORTANTE: Los atributos deben estar normalizados en  $[0,1]$  para no priorizar unos sobre otros**

# Ejemplo de Clasificador: k-NN

---

Para normalizar los datos, hay que saber el intervalo de dominio de cada uno de los atributos

Dado un valor  $x_j$  perteneciente al atributo  $j$  del ejemplo  $x$  y sabiendo que el dominio del atributo  $j$  es  $[Min_j, Max_j]$ , el valor normalizado de  $x_j$  es:

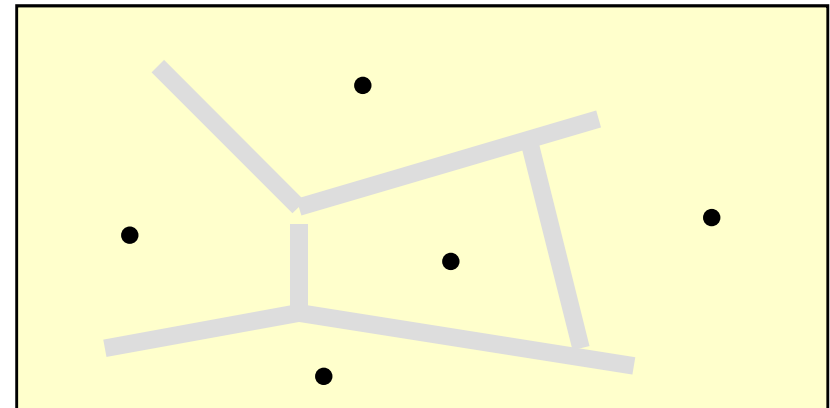
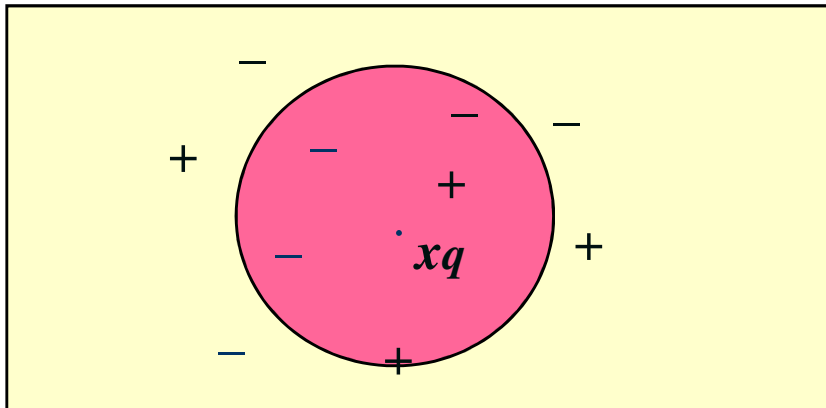
$$x_j^N = \frac{x_j - Min_j}{Max_j - Min_j}$$



# Ejemplo de Clasificador: k-NN

---

- El  $k$ -NN devuelve la clase más repetida de entre los  $k$  ejemplos de entrenamiento más cercanos a  $x_q$
- Diagrama de Voronoi: superficie de decisión inducida por 1-NN para un conjunto dado de ejemplos de entrenamiento



# Ejemplo de Clasificador: k-NN

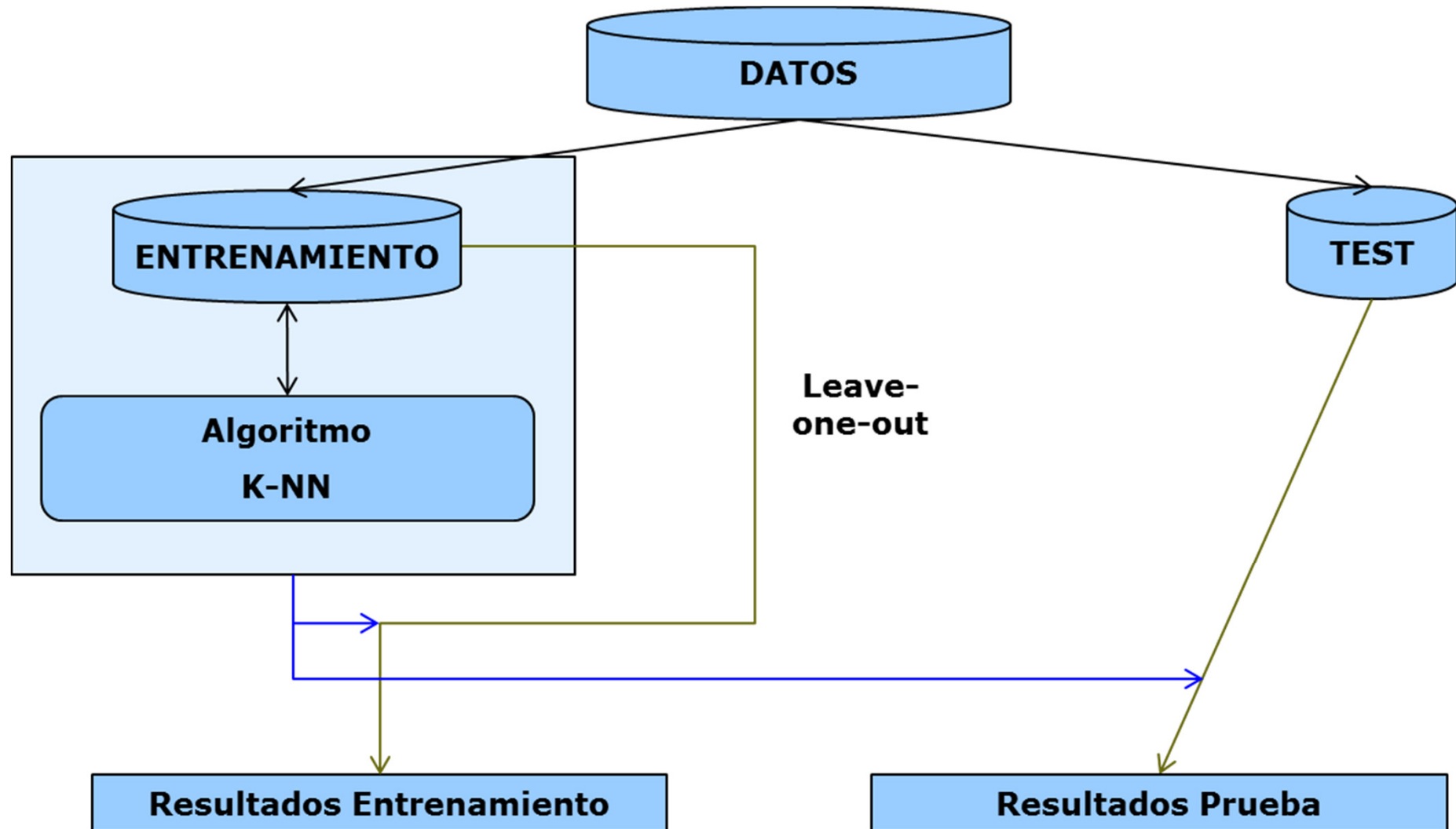
---

## *Cálculo del Porcentaje de Entrenamiento en 1-NN:*

- *En el algoritmo 1-NN, no es posible calcular el porcentaje de acierto sobre el conjunto de entrenamiento de un modo directo.*
- *Si intentásemos clasificar un ejemplo del conjunto de entrenamiento directamente con el clasificador 1-NN, el ejemplo más cercano sería siempre él mismo, con lo que se obtendría un 100% de acierto.*
- *Para remediar esto, se debe seguir el procedimiento denominado dejar uno fuera (“leave one out”).*
- *Para clasificar cada ejemplo del conjunto de entrenamiento, se busca el ejemplo más cercano **sin considerar a él mismo**.*
- *Por lo demás, se opera igual: cada vez que la clase devuelta coincida con la clase real del ejemplo, se contabiliza un acierto.*
- *El porcentaje final de acierto es el número de aciertos entre el número total de ejemplos.*

# Ejemplo de Clasificador: k-NN

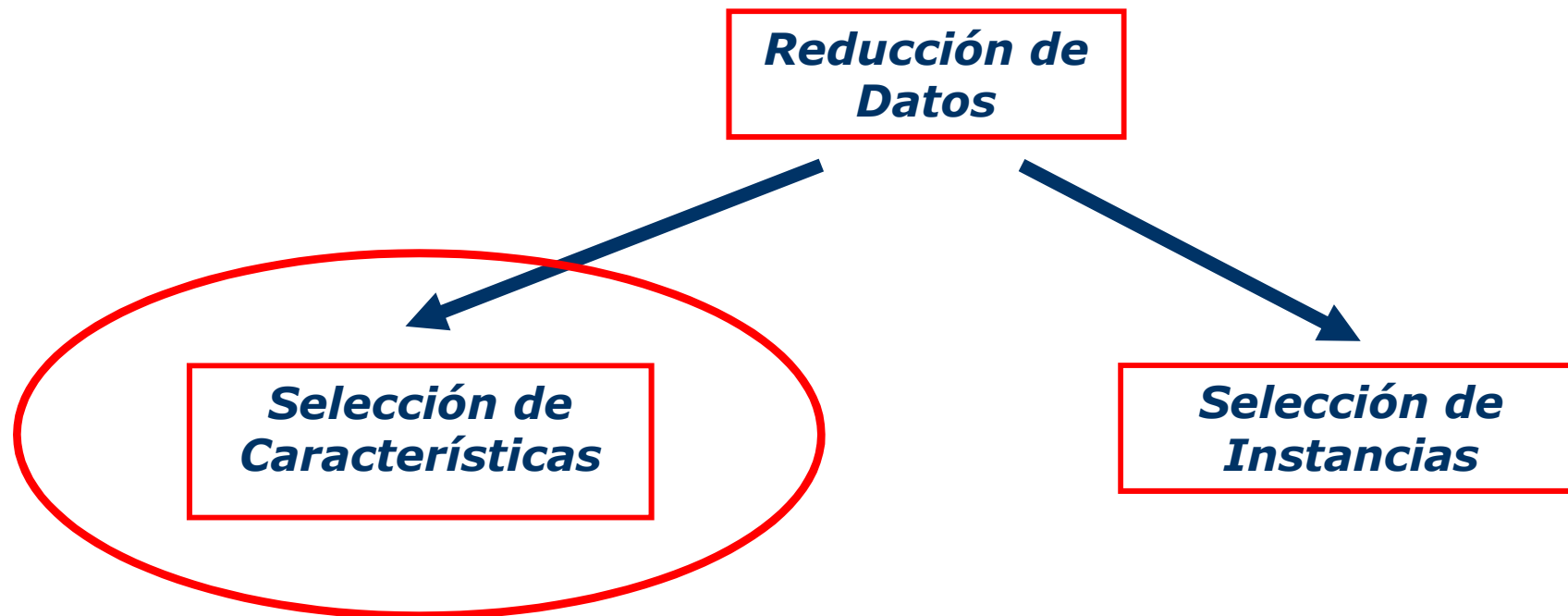
---



# Definición del Problema de Selección de Características

---

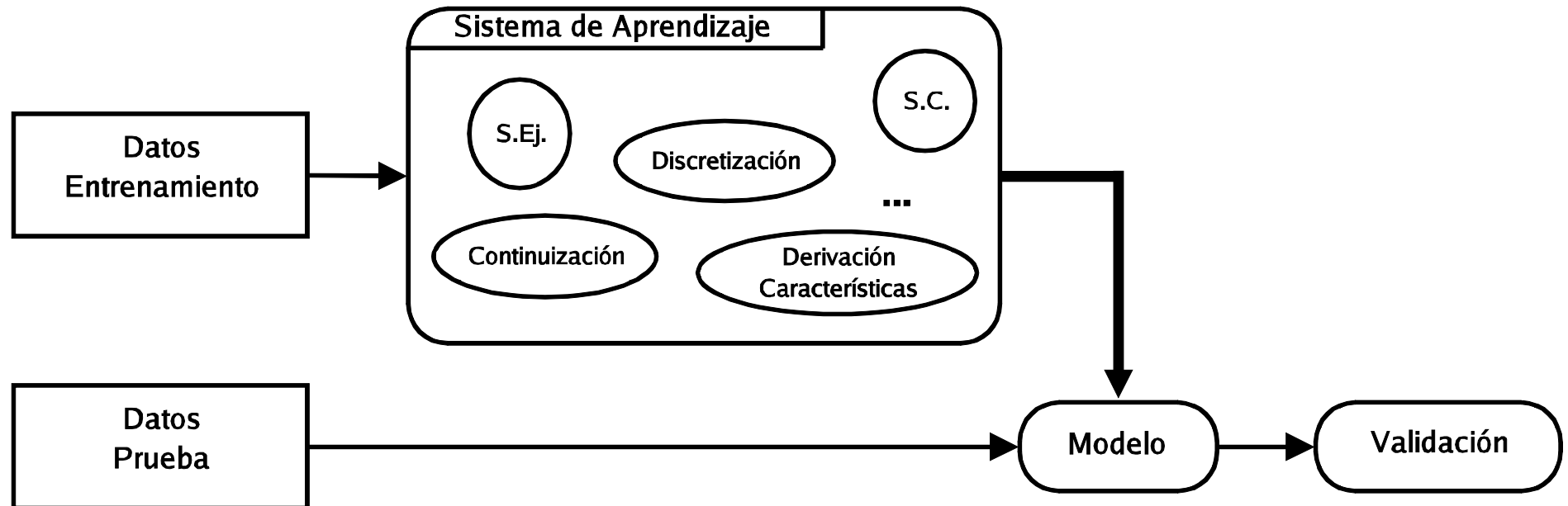
*A veces el problema de clasificación es demasiado complejo y necesita ser procesado previamente para ser tratable → **Reducción de datos***



# Definición del Problema de Selección de Características

---

## Proceso de aprendizaje



# Definición del Problema de Selección de Características

---

*Situación **ideal** en construcción de clasificadores:  
disponer del máximo de información → el mayor  
número de características*

- Inconvenientes:
  - Características redundantes
  - Características irrelevantes
  - Costes en captación, almacenamiento, eficiencia

# Definición del Problema de Selección de Características

*Var. 1.*

*Var. 5*

*Var. 13*

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
A	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1
B	0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1
C	0	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1
D	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
E	0	1	0	0	0	1	1	0	1	1	0	0	0	0	1	0
F	1	1	1	0	1	1	0	0	1	0	1	0	0	1	0	0

**Ejemplo**

# Definición del Problema de Selección de Características

## Ejemplo

$F = \{f_1 \quad f_2 \quad f_3 \quad \dots \quad f_n\}$						$Cl$
$e_1$	8.1	Sí	ESO	...	Casad.	A
$e_2$	3.4	No	Diplom.	...	Divor.	B
$e_3$	-2.5	No	Doctor	...	Solt.	A
$\vdots$						
$e_n$	-1.0	No	Ingeni.	...	Viud.	C

$\Pi_S \rightarrow$

$S = \{f_1 \quad f_3\}$				$Cl$
$e_1$	8.1	ESO		A
$e_2$	3.4	Diplom.		B
$e_3$	-2.5	Doctor		A
$\vdots$				
$e_n$	-1.0	Ingeni.		C



# Definición del Problema de Selección de Características

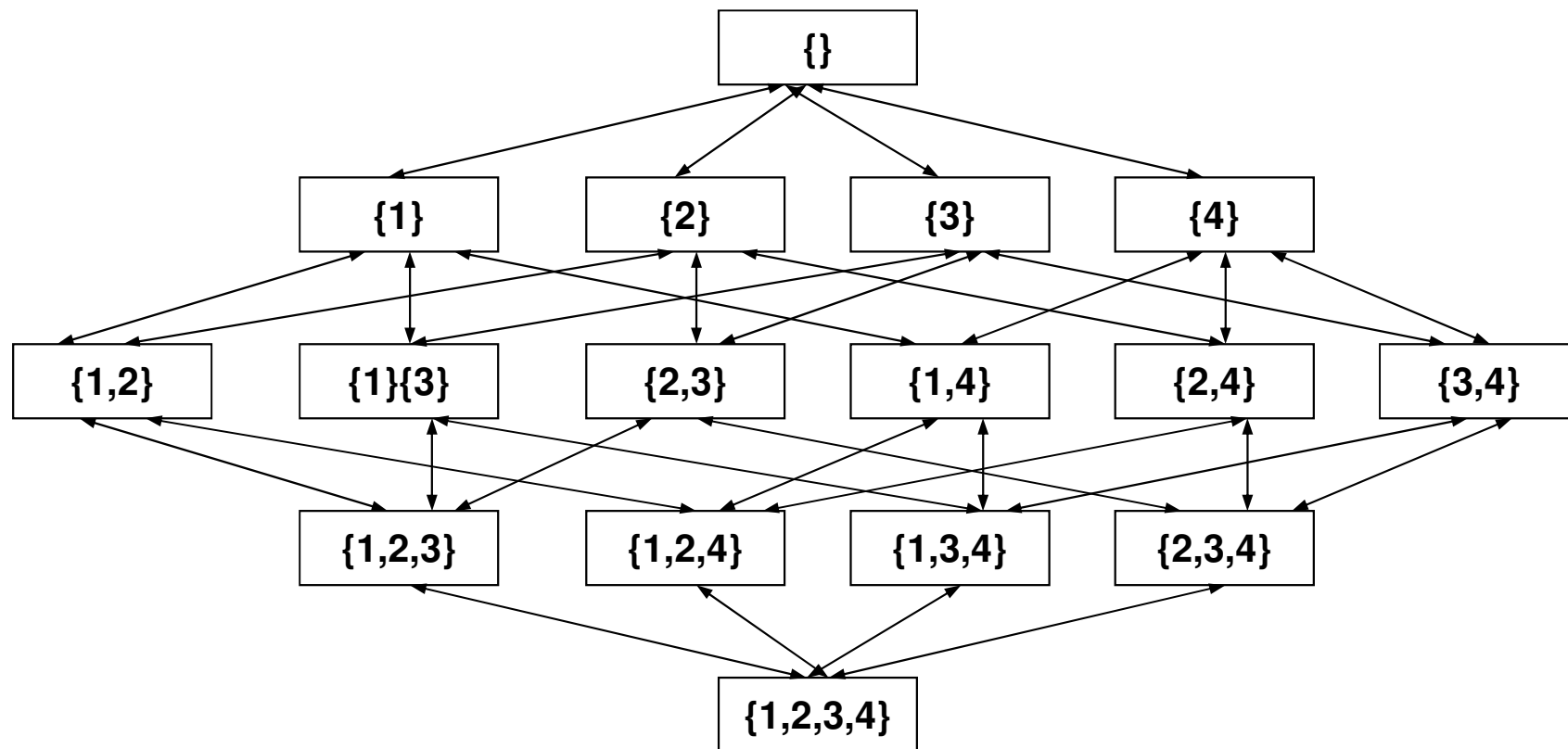
---

## Objetivo:

- Seleccionar un subconjunto del conjunto total de características, de tal forma que los clasificadores que se construyan a partir de él sean *mejores*
- Existen distintos criterios para determinar cuándo el clasificador generado es mejor
- Es un problema de **búsqueda** en el espacio de todos los posibles subconjuntos de características
- Hay  $2^n$  posibles subconjuntos para  $n$  características. Es un problema **NP-duro**

# Definición del Problema de Selección de Características

**La selección de características (SC) se puede considerar como un problema de búsqueda:**



**Complejidad**

## Ventajas de la SC:

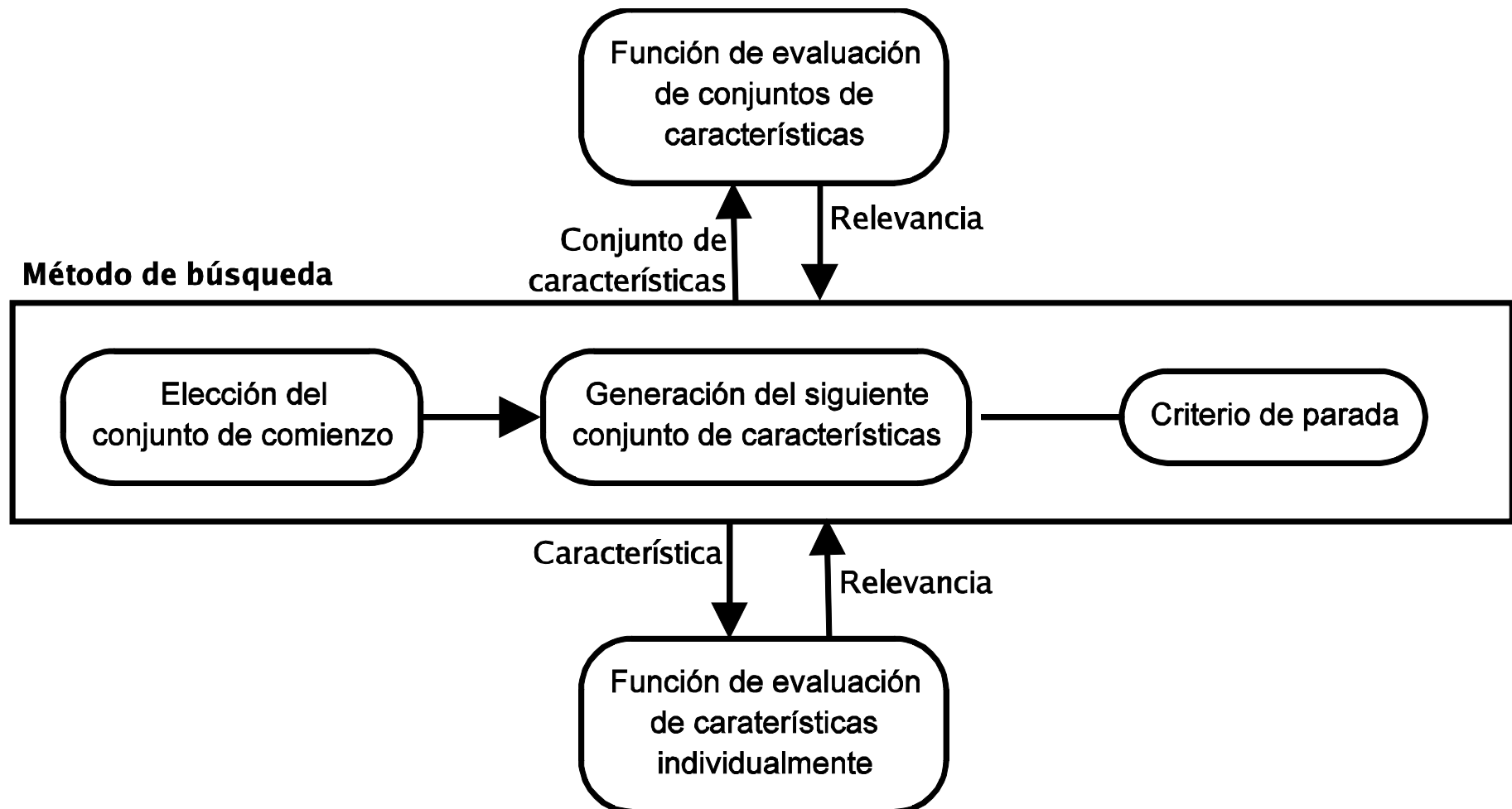
- **Eficiencia.** El tiempo que se requiere para construir un clasificador depende del número de características
- **Mejora en los resultados.** Principio de la “navaja de Occam”
- **Reducción de costes** de adquisición y almacenamiento
- Permite aplicar técnicas de diseño de clasificadores que sufran de la “**maldición de la dimensionalidad**”
- **Mejora la interpretabilidad**

## Tipos de métodos de SC:

- **Tipo filtro (*filter*)**. La SC es previa a la construcción del clasificador  
Ej: métodos basados en información mutua
- **Tipo envoltante (*wrapper*)**. La SC y la construcción del clasificador son etapas conjuntas
- **Tipo incrustado (*embedded*)**. El método de diseño del clasificador realiza de forma implícita una SC  
Ej: árboles de clasificación

# Definición del Problema de Selección de Características

## Descomposición modular del problema de SC:



### **Funciones de evaluación:**

- Medidas basadas en **consistencia**
- Medidas basadas en **Teoría de la Información**
- Medidas basadas en **distancia**
- Medidas basadas en **rendimiento de clasificadores**

## Nuestra función de evaluación:

- Rendimiento promedio de un clasificador 3-NN (considerando  $k=3$  vecinos) aplicando validación sobre el conjunto  $T$  de datos: *tasa\_clas*
- *tasa\_clas* mide el porcentaje de instancias correctamente clasificadas pertenecientes a  $T$ :

$$tasa\_clas = 100 \cdot \frac{\text{n}^\circ \text{ instancias bien clasificadas de } T}{\text{n}^\circ \text{ instancias en } T}$$

- El objetivo es obtener el subconjunto de características que maximiza esta función

# Representaciones

---

- Hay varias representaciones para el problema de SC: binaria, entera, permutación
- **Su uso puede depender de si se especifica el número de características deseado a priori o si lo aprende el método**
- Cuando no se especifica, en un problema con  $n$  características existen  $2^n$  posibles subconjuntos de características
- En caso de indicar un número de características  $m$  predefinido, el tamaño del espacio de búsqueda varía



# Representaciones

---

## 1. Representación binaria: Vector binario de tamaño $n$ :

$$S = (f_1, f_2, \dots, f_n), \quad \text{donde } f_i \in \{0, 1\}$$

$f_1$	$f_2$	.....	$f_{n-1}$	$f_n$
-------	-------	-------	-----------	-------

Un 1 en la posición  $f_i$  indica la selección de la característica en cuestión, un 0 su no selección

La representación binaria es adecuada para el **aprendizaje del número de características**. El tamaño del espacio es  $2^n$

Si se desea extraer un número fijo, no es una buena opción porque requiere el manejo de restricciones. Una solución sólo es factible (válida) si tiene exactamente  $m$  1's

# Representaciones

**2. Representación entera:** Vector de tamaño  $m \leq n$  con valores  $\{1, \dots, n\}$  que indican el no. de la característica seleccionada:

$$S = (\text{sel}_1, \text{sel}_2, \dots, \text{sel}_m), \quad \text{donde } \text{sel}_i \in \{1, \dots, n\}$$

$\text{sel}_1$	$\text{sel}_2$	.....	$\text{sel}_{m-1}$	$\text{sel}_m$
----------------	----------------	-------	--------------------	----------------

Esta representación está más pensada para el caso en que se desee extraer **un número de características concreto**. El tamaño del espacio de búsqueda es el número combinatorio  $\binom{n}{m}$

Aún así, hay que manejar restricciones. Para que una solución sea factible, no puede tener características repetidas

Además, ¡es redundante! Hay  $m!$  vectores distintos que representan el mismo subconjunto de características → **El espacio de búsqueda es más grande que el espacio de soluciones real**

# Representaciones

---

- 3. Representación de orden:** También es posible codificar una solución al problema de la SC con una permutación

Se tiene una permutación  $\pi$  de tamaño  $n$ . Las características seleccionadas son las primeras del orden definido:

$S = \pi = (of_1, of_2, \dots, of_n)$ , donde  $of_i \in \{1, \dots, n\}$ , sin repetidos

$of_1$	$of_2$	.....	$of_{n-1}$	$of_n$
--------	--------	-------	------------	--------

El tamaño del espacio de búsqueda es  $n!$

# Representaciones

---

Esta representación es muy flexible:

- Se puede usar con un **número de características prefijado**, escogiendo las  $m$  primeras
- Se puede usar para **aprender el número de características** considerando que cada permutación  $\pi$  codifica varias soluciones (una con la primera característica, otra con las dos primeras, ...)

A la hora de evaluar la solución, se decodificarían todas, se evaluarían y se asignaría a esa permutación  $\pi$  el mejor valor de función objetivo de todas

En ambos casos, sigue siendo redundante. Distintas ordenaciones de las mismas características representan la misma selección

## Métodos de búsqueda:

- Enumeración explícita (búsqueda completa)
  - Garantiza el óptimo
- Búsqueda secuencial
  - Método voraz (*greedy*) que parte de un conjunto vacío de características y va añadiendo de una en una (*forward*), o parte del conjunto total y va eliminando de una en una (*backward*)
- Búsqueda probabilística
  - Métodos MonteCarlo y Las Vegas
- Búsqueda con metaheurísticas

# Solución *greedy*

---

## **Método *Sequential Forward Selection* (SFS):**

- Parte del conjunto vacío de características
- En cada paso añade la característica más prometedora
- La más prometedora es la que produce una mayor ganancia en la función de evaluación con respecto al conjunto de características ya seleccionado
- El algoritmo para cuando ninguna de las características que quedan por seleccionar produce una mejora en la función de evaluación
- Usaremos la función de evaluación indicada, basada en el acierto de clasificación del k-NN con  $k=3$

# Solución *greedy*

---

## Algoritmo SFS:

$F \leftarrow \{f_1, \dots, f_n\}$

$S \leftarrow \emptyset$

$fin \leftarrow falso$

*Mientras* ( $F \neq \emptyset$  y  $!fin$ )

$fp \leftarrow Característica\_más\_prometedora(F)$

*Si*  $S \cup \{fp\}$  mejora la solución actual  $S$  ( $f(S \cup \{fp\}) \geq f(S)$  )

$S \leftarrow S \cup \{fp\}$

*en otro caso*

$fin \leftarrow verdad$

$F \leftarrow F - \{fp\}$

*Devolver*  $S$

# Búsquedas por Trayectorias Simples

---

- **Representación binaria:** Problema de selección: un vector binario  $s=(s_1, \dots, s_n)$  en el que cada posición  $i$  representa una característica y su valor 0/1 indica si está o no seleccionada

No tiene restricciones al no exigir un número predeterminado de características

- **Operador de vecino de inversión y su entorno:** El entorno de una solución  $s$  está formado por las soluciones accesibles desde ella a través de un movimiento de inversión  $Flip(s,i)$ :

Dada una solución (subconjunto de características) se escoge una característica concreta  $i$  y se intercambia su selección (pasa a valer 1, si estaba a 0, o viceversa):

$$s=(s_1, \dots, s_i, \dots, s_n) \rightarrow s'=(s_1, \dots, s'_i, \dots, s_n)$$



# Búsquedas por Trayectorias Simples

---

- $Flip(s,i)$  verifica las restricciones, si la solución original  $s$  es factible siempre genera una solución vecina  $s'$  factible
- Su aplicación provoca que el tamaño del entorno sea:

$$|E(s)| = n$$

- El problema de la SC **no permite realizar un cálculo factorizado del coste** de forma sencilla

Tampoco es fácil definir una preferencia entre las características para explorar el entorno. Podría considerarse alguna medida de cantidad de información para esta tarea

# Búsqueda Local para la SC

---

- Algoritmo de **búsqueda local del primer mejor**: en cuanto se genera una solución vecina que mejora a la actual, se aplica el movimiento y se pasa a la siguiente iteración
- Se detiene la búsqueda cuando se ha explorado el vecindario completo sin obtener mejora
- No se considera ningún tipo de **factorización** ni ningún mecanismo específico de exploración del entorno, se estudian los vecinos por orden

# Enfriamiento Simulado para la SC

---

- **Representación**: vector binario, igual que en la BL
- **Operador de generación de vecinos**: inversión, igual que en la BL
- **Exploración del vecindario**: En cada iteración del bucle interno se genera una única solución vecina, de forma aleatoria, y se compara con la actual
- **Esquema de enfriamiento**: esquema de Cauchy modificado
- **Condición de enfriamiento  $L(T)$** : cuando se genere un número máximo de soluciones vecinas, *máx\_vecinos*, o cuando se acepte un número máximo de los vecinos generados, *máx\_éxitos*
- **Condición de parada**: cuando se alcance un número máximo de iteraciones o cuando el número de éxitos en el enfriamiento actual sea 0

# Búsqueda Tabú para la SC

---

- **Representación**: vector binario, igual que en la BL
- **Operador de generación de vecinos**: inversión, igual que en la BL
- **Lista de valores de atributos y posiciones tabú**: Cuando se acepta una nueva solución, se almacena información sobre ella en la lista tabú

Se guarda el **movimiento que generó la solución** (es decir, el índice de la característica de la que se alteró la selección): (*i*)

No se permiten futuros intercambios que provoquen que **el estatus de seleccionada/no seleccionada de la característica *i* se vuelva a modificar**

Ejemplo:

1	2	3	4	5	6	7	8
(0	1	<u>1</u>	0	0	1	1	0)
(0	1	<u>0</u>	0	0	1	1	0)

LT = LT + {3} (si está llena, eliminar el registro más antiguo)

# Búsqueda Tabú para la SC

---

- **Estrategia de selección de vecino**: Examinar un número de soluciones vecinas generadas aleatoriamente y escoger la mejor que verifique los criterios tabú o que no los verifique y supere el criterio de aspiración
- **Criterio de aspiración**: Tener mayor coste que la mejor solución obtenida hasta el momento
- **Estrategias de reinicialización**:
  - Generar una solución aleatoria y continuar el algoritmo a partir de ella (**diversificación**)
  - Volver a la mejor solución obtenida hasta el momento y continuar el algoritmo a partir de ella (**intensificación**)
  - Generar una nueva solución poco frecuente a partir de la memoria a largo plazo y continuar a partir de ella (**diversificación controlada**)

# Búsqueda Tabú para la SC

---

- **Memoria a largo plazo:** Se usa un vector *frec* de tamaño  $n$  que almacena el número de veces que cada característica ha sido seleccionada en las soluciones aceptadas durante la búsqueda
- Para reinicializar con diversidad, se genera una nueva solución inicial aleatoria con una probabilidad inversa a la de la memoria de frecuencias:

$$s_i = \begin{cases} 1, & \text{si } u_{\text{aleatorio}} < 1 - \frac{frec_i}{num\_soluciones} \\ 0, & \text{si } u_{\text{aleatorio}} \geq 1 - \frac{frec_i}{num\_soluciones} \end{cases}$$

- La elección de qué mecanismo de reinicialización aplicar se hace igual que en el QAP

# Bases de Datos a Utilizar

---

- En la actualidad, hay muchas bases de datos que se utilizan como bancos de prueba (*benchmarks*) para comprobar el rendimiento de los algoritmos de clasificación
- El UCI es un repositorio de bases de datos para aprendizaje automático muy conocido
- Está accesible en la Web en:  
<http://www.ics.uci.edu/~mlearn/MLRepository.html>

# Bases de Datos a Utilizar

---

- A partir de este repositorio, se ha desarrollado un formato para definir todas las cualidades de una base de datos en un único fichero
- Se trata del formato ARFF, utilizado en WEKA (<http://www.cs.waikato.ac.nz/ml/weka/>)
- Un fichero ARFF contiene dos partes:
  - Datos de cabecera: Líneas que comienzan por @. Contienen información acerca de los atributos: significado y tipo
  - Datos del problema: Líneas de datos. Son los ejemplos en sí. Cada línea se corresponde con un ejemplo y los valores están separados por comas



# Bases de Datos a Utilizar

---

## **Ejemplo fichero ARFF:**

```
@relation iris  
@attribute sepalLength real  
@attribute sepalWidth real  
@attribute petalLength real  
@attribute petalWidth real  
@attribute class {Iris-setosa, Iris-versicolor, Iris-virginica}  
@data  
5.1, 3.5, 1.4, 0.2, Iris-setosa  
4.9, 3.0, 1.4, 0.2, Iris-setosa  
7.0, 3.2, 4.7, 1.4, Iris-versicolor  
6.0, 3.0, 4.8, 1.8, Iris-virginica
```

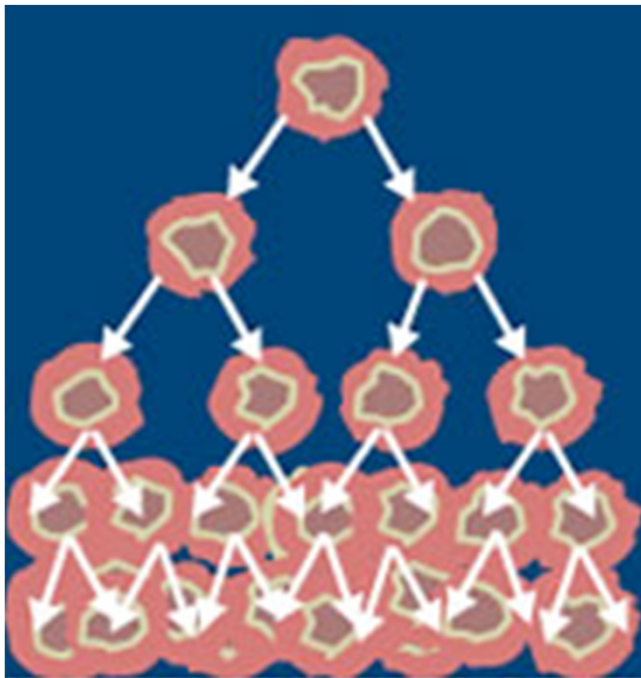
- 5 Atributos, los 4 primeros de tipo real y el último de tipo nominal
- La clase es el último atributo (atributo de salida) con los posibles valores definidos
- Los datos van a continuación de la directiva @data

# Bases de Datos a Utilizar

---

Trabajaremos con tres bases de datos: Wdbc, Movement\_Libras y Arrhythmia

**Wdbc** es una base de datos contiene atributos calculados a partir de una imagen digitalizada de una aspiración con aguja fina de una masa en la mama. Se describen las características de los núcleos de las células presentes en la imagen. La tarea consiste en determinar si un tumor encontrado es benigno o maligno..



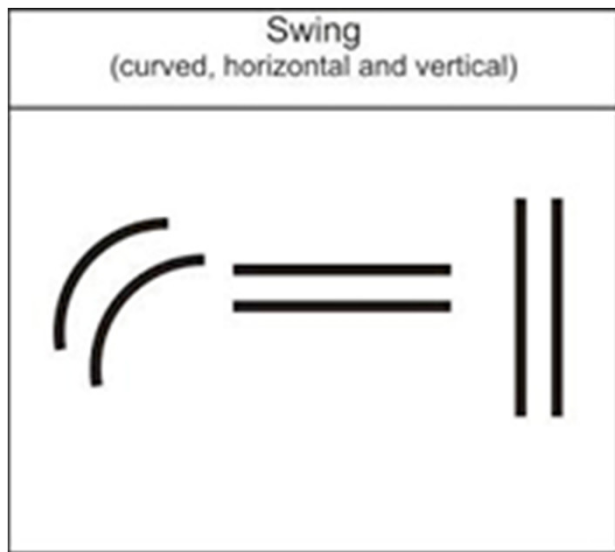
- Consta de 569 ejemplos
- Consta de 31 atributos (clase incluida)
- Consta de 2 clases (M = maligno, B = benigno)
- Atributos: Valores reales calculados de cada núcleo de la célula (10 valores \* 3 células): radio, textura, perímetro, área, concavidad, simetría, ...

# Bases de Datos a Utilizar

---

Trabajaremos con tres bases de datos: Wdbc, Movement\_Libras y Arrhythmia

**Movement\_Libras** es una base de datos en la que cada clase referencia a un tipo de movimiento de la mano en LIBRAS (nombre portugués 'Língua Brasileira de Sinais', la lengua oficial de señales brasileña)



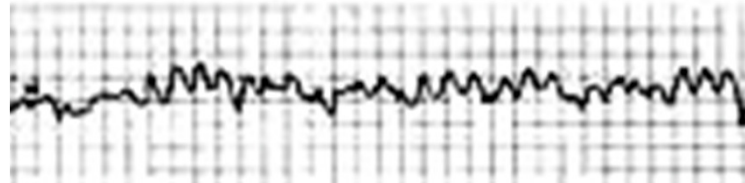
- Consta de 360 ejemplos
- Consta de 91 atributos (clase incluida)
- Consta de 15 clases (24 ejemplos por clase)
- Atributos:
  - Representación en 90 características de las coordenadas de los movimientos.

# Bases de Datos a Utilizar

---

Trabajaremos con tres bases de datos: Wdbc, Movement\_Libras y Arrhythmia

**Arrhythmia** contiene datos que se utilizan para distinguir entre la presencia y la ausencia de arritmia cardiaca y clasificarlo en uno de los 5 grupos mayoritarios.



- Consta de 386 ejemplos (seleccionados de los 452 originales, eliminando los ejemplos de clases minoritarias e imputando valores perdidos con la media/moda no supervisada)
- Consta de 279 atributos (clase incluida)
- Consta de 5 clases (de 16 originales, solo se mantienen las que superan un mínimo de 20 representantes)
- Atributos: Edad, sexo, altura, peso, y muchos parámetros de medición (duración QRS, intervalo P-R, Q-T, ondas, canales, ...)