Dr. Jesús Ariel Carrasco Ochoa ariel@inaoep.mx
Oficina 8311

Contenido

- Introducción
- Estrategias de selección
- Técnicas filter
 - Condensación
 - Edición
- Técnicas wrapper

- Selección de instancias (Instance selection), Selección de prototipos (Prototype selection)
- Preprocesamiento
 - Clasificación Supervisada
 - Regresión

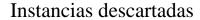
- Selección de instancias (Instance selection), Selección de prototipos (Prototype selection)
- Preprocesamiento
 - Clasificación Supervisada
 - Regresión

Por qué hacer selección de objetos en clasificación supervisada

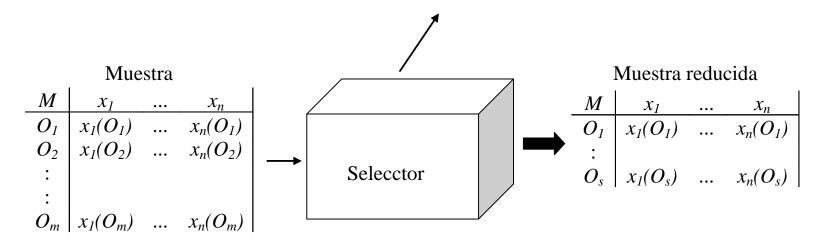
- Mejorar los resultados de la clasificación
- Reducir el costo de la clasificación
- Acelerar el proceso de clasificación

El objetivo es seleccionar un subconjunto reducido de objetos útiles para la clasificación eliminando

- Objetos irrelevantes
- Objetos redundantes
- Objetos ruidosos



$$egin{array}{c|cccc} M & x_1 & \dots & x_n \\ \hline O_i & x_I(O_i) & \dots & x_n(O_i) \\ dots & & & & & \\ O_k & x_I(O_k) & \dots & x_n(O_k) \\ \hline \end{array}$$



Estrategias de Selección

o *filter*.- La selección se hace con un criterio independiente del clasificador.

 wrapper.- La selección se hace usando información del mecanismo de clasificación.

Estrategias filter

- Se busca eliminar diferentes tipos de objetos.
 - Condensación: Tratan de Eliminar objetos irrelevantes o redundantes
 - Edición: Tratan de eliminar ouliers (ruido).
 - Híbridos

Métodos de Condensación

- Tratan de eliminar objetos que no son útiles para clasificar correctamente al conjunto de entrenamiento
- o CNN (Hart, 1968)
- o RNN (Gates, 1972)
- CNN+RNN (Gates, 1972)

CNN (Hart, 1968)

- Inicializar en S=Ø
- 2. Para cada objeto O∈D
 - 3. Si el vecino más cercano a O en S es de clase diferente
 - 4. $S = S \cup \{0\}$
 - 5. $D = D \{O\}$
- 6. Si se agregaron objetos a S ir a 2
- 7. Tomar D=S

RNN (Gates, 1972)

- Inicializar en S=D
- 2. Para cada objeto O∈S
 - Si ∀O'∈D el vecino más cercano a O' en S-{O} es de su misma clase
 - 4. $S = S \{0\}$
- 5. Tomar D=S

CNN+RNN (Gates, 1972)

- Inicializar en S=Ø
- 2. Para cada objeto O∈D
 - 3. Si el vecino más cercano a O en S es de clase diferente
 - 4. $S = S \cup \{0\}$
 - 5. $D = D \{0\}$
- 6. Si se agregaron objetos a S ir a 2
- 7. Para cada objeto O∈S
 - 8. Si ∀O'∈D el vecino más cercano a O' en S-{O} es de su misma clase
 - 9. $S = S \{0\}$
- 10. Tomar D=S

Métodos de Edición

 Tratan de eliminar objetos que no son bien clasificados o causan la mala clasificación de otros

- ENN (Wilson, 1972)
- o RENN (Tomek, 1976)
- MultiEdit (Kittler, 1980)

ENN (Wilson, 1972)

- Inicializar en S=D
- 2. Para cada objeto O∈S
 - 3. Si la mayoría de sus k vecinos más cercanos en D son de diferente clase
 - 4. $S = S \{0\}$
- 5. D=S

RENN (Tomek, 1976)

- Inicializar en S=D
- 2. Para cada objeto O∈S
 - 3. Si la mayoría de sus k vecinos más cercanos en S son de diferente clase
 - 4. $S = S \{O\}$
- 5. Tomar D=S
- 6. Si se eliminó algún objeto de S ir a 1

MultiEdit (Kittler, 1980)

- 1. Hacer una partición aleatoria (D_0, D_1, D_r) de D tal que todos los D_i sean de aproximadamente el mismo tamaño.
- 2. Para cada D_i aplicar ENN pero buscando los vecinos en $D_{(i+1) \mod r}$
- 3. Tomar $D=S_0 \cup S_1 \cup ... \cup S_r$
- 4. Si se eliminó algún objeto ir a 1

Métodos de Híbridos

- Combinan ideas de Condensación y Edición.
 - Tratan de eliminar objetos que no son útiles para clasificar correctamente al conjunto de entrenamiento
 - Tratan de eliminar objetos que no son bien clasificados o causan la mala clasificación de otros
- DROP 1-5 (Wilson & Martínez, 2000)
- ICF (Brighton & Mellish, 1999, 2002)

DROP 1-5 (Wilson & Martínez, 2000)

Decremental Reduction Optimization Procedure.

- Inicializar en S=D
- 2. Para cada objeto O∈S
 - 3. Si la mayoría de sus k vecinos más cercanos en D son de diferente clase
 - 4. $S = S \{0\}$
- 5. D=S

DROP 1 (Wilson & Martínez, 2000)

- Inicializar en S=D
- 2. Para cada objeto O∈S
 - 3. Si ∀O'∈S el vecino más cercano a O' en S-{O} es de su misma clase
 - 4. $S = S \{0\}$
- 5. Tomar D=S

DROP 2 (Wilson & Martínez, 2000)

- Inicializar en S=D
- Ordenar S descendentemente de acuerdo a la distancia con su vecino más cercano de diferente clase
- 3. Para cada objeto O∈S
 - 4. Si ∀O'∈D el vecino más cercano a O' en S-{O} es de su misma clase
 - 5. $S = S \{0\}$
- 6. Tomar D=S

DROP 3 (Wilson & Martínez, 2000)

- Aplicar ENN
- Inicializar en S=D
- Ordenar S descendentemente de acuerdo a la distancia con su vecino más cercano de diferente clase
- 4. Para cada objeto O∈S
 - 5. Si ∀O'∈D el vecino más cercano a O' en S-{O} es de su misma clase
 - 6. $S = S \{0\}$
- 7. Tomar D=S

DROP 4 (Wilson & Martínez, 2000)

- Inicializar en S=D
- 2. Para cada objeto O∈S
 - 3. Si la mayoría de sus k vecinos más cercanos en D son de diferente clase y ∀O'∈D el vecino más cercano a O' en S-{O} es de su misma clase
 - 4. $S = S \{0\}$

DROP 4 (Wilson & Martínez, 2000)

- Ordenar S descendentemente de acuerdo a la distancia con su vecino más cercano de diferente clase
- 6. Para cada objeto O∈S
 - 7. Si ∀O'∈D el vecino más cercano a O' en S-{O} es de su misma clase
 - 8. $S = S \{0\}$
- 9. Tomar D=S

DROP 5 (Wilson & Martínez, 2000)

- Inicializar en S=D
- Ordenar S Ascendentemente de acuerdo a la distancia con su vecino más cercano de diferente clase
- 3. Para cada objeto O∈S
 - 4. Si ∀O'∈D el vecino más cercano a O' en S-{O} es de su misma clase
 - 5. $S = S \{0\}$
- 6. Tomar D=S

ICF (Brighton & Mellish, 1999, 2002)

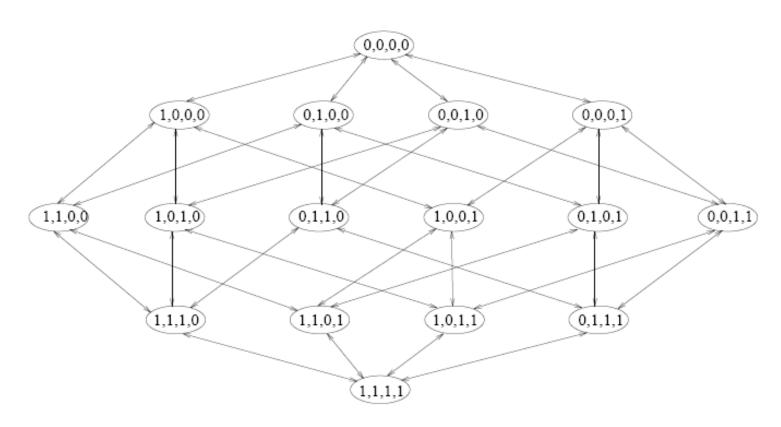
- Aplicar ENN
- 2. Inicializar en S=D
- 3. Para cada objeto O∈S
 - 4. Reachable(O)=número de objetos de su clase que son más cercanos que el objeto más cercano de otra clase.
 - 5. Coverage(O)= número de objetos O' de su clase para los cuales O es más cercano a O' que el más cercano a O' de diferente clase
- 6. Para cada objeto O∈S
 - Si Reachable(O)>Coverage(O)
 - 8. $S = S \{0\}$
- Tomar D=S

Estrategias wrapper

- Evalúan subconjuntos de objetos utilizando un clasificador.
- Para evitar la búsqueda exhaustiva siguen alguna estrategia de búsqueda.
 - Comúnmente estrategias ávidas o aleatorias

Estrategias wrapper

 Para n objetos, el espacio de búsqueda es de tamaño 2ⁿ



Estrategias wrapper

- Búsqueda exhaustiva
- Búsqueda secuencial
 - Hacia atrás (backward)
 - Hacia adelante (forward)
 - Flotante (floating)
- Búsqueda aleatoria
 - Algoritmos genéticos
 - Búsqueda tabú
 - ...

Búsqueda exhaustiva

- El tamaño del espacio de búsqueda es 2ⁿ
- Si se busca un número predefinido de objetos el espacio de búsqueda es de tamaño

$$\binom{n}{k}$$

 Para seleccionar 200 objetos de 500 el espacio de búsqueda es de tamaño

$$\begin{pmatrix} 500 \\ 200 \end{pmatrix} \approx 5 \times 10^{-144} \approx 2^{480}$$

Sequential Forward Selection

Sea D el conjunto de todos los objetos

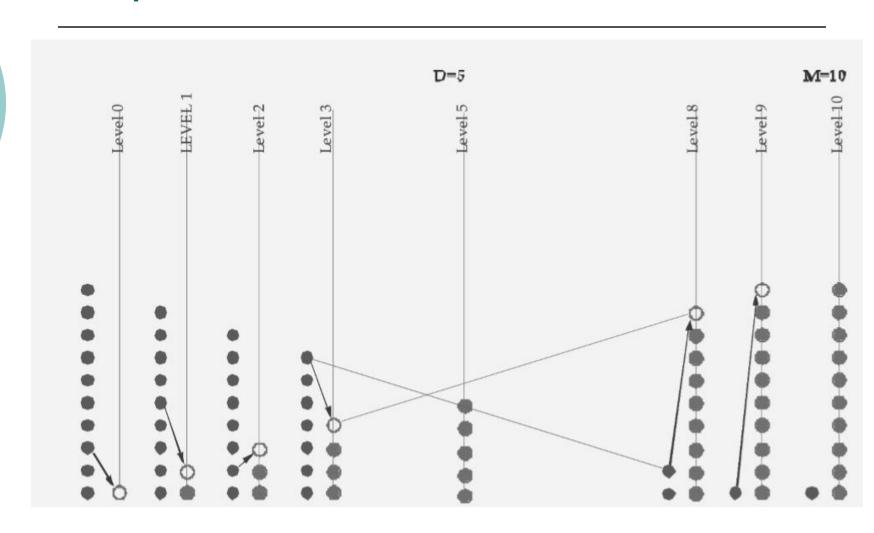
- o S=Ø
- Repetir

$$O = \max_{O \in D \setminus S} (q(S \cup \{O\}))$$

$$S = S \cup \{O\}$$

o Hasta $|S|=k/q(S)<q(S\setminus\{O\})/q(S)>t$

Sequential Forward Selection



Sequential Bacward Selection

Sea D en conjunto de todos los objetos

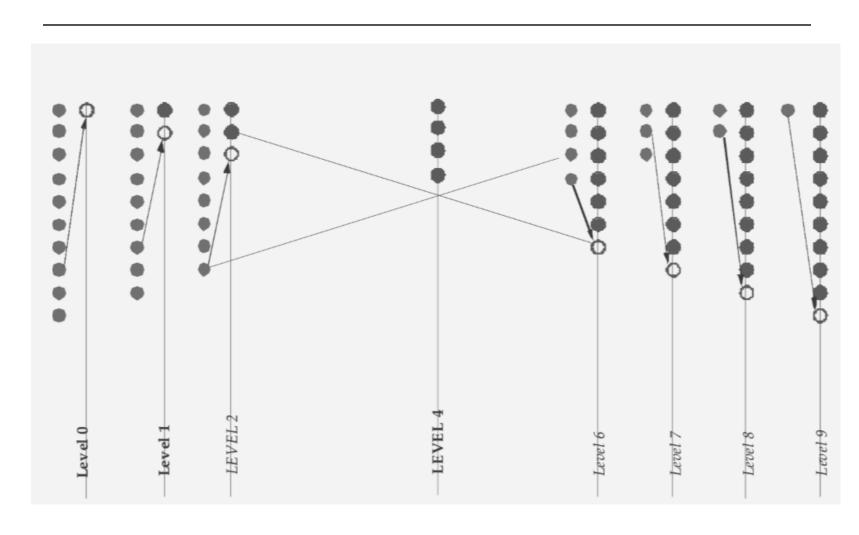
- S=D
- Repetir

$$O = \max_{O \in S} (q(S \setminus \{O\}))$$

$$S = S \setminus \{O\}$$

Hasta |S|=k / q(S)<q(S∪{O}) / q(S)>t

Sequential Backward Selection



Sequential Floating Forward Selection

• S=
$$\emptyset$$
• Repetir $O = \max_{O \in D \setminus S} (q(S \cup \{O\}))$

$$S = S \cup \{O\}$$
Repetir $O = \max_{O \in S} (q(S \setminus \{O\}))$

$$S = S \setminus \{O\}$$
Hasta $q(S) < q(S \cup \{O\})$
• Hasta $|S| = k / q(S) > t$

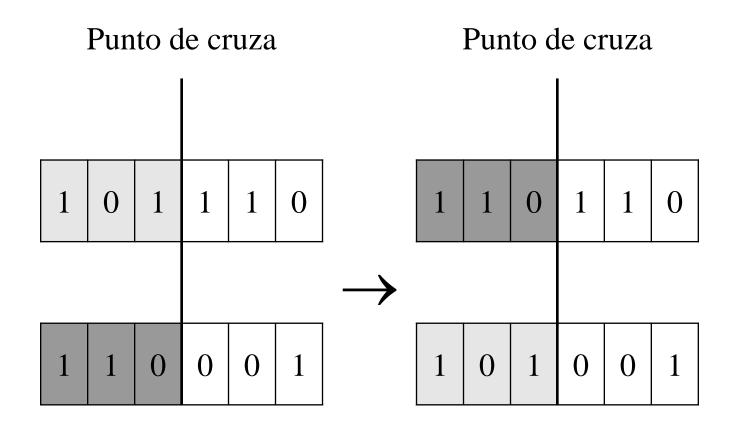
Sequential Floating Backward Selection

```
o S=Ø
o Repetir O = \max (q(S \setminus \{O\}))
                       O \in S
                S = S \setminus \{O\}
                Repetir O = \max (q(S \cup \{O\}))
                                 O \in D \setminus S
                           S = S \cup \{O\}
                Hasta q(S) \leq q(S \setminus \{O\})
                S\{O}
\circ Hasta |S|=k/q(S)>t
```

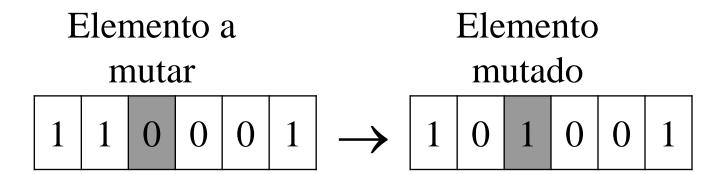
Individuos (suponiendo n objetos)

x_1	x_2	• • •	\mathcal{X}_n
0	1	0/1 0/1	1

Cruza (combinación de individuos)



Mutación (alteración de individuos)



- Generar población inicial P
- Para i=1, ..., numGeneraciones evalúa(P) P2 = cruza(P)evalúa(P2) $P3 = mutación (P \cup P2)$ evalúa(P3) $P = selecciona(P \cup P2 \cup P3)$ o salida = mejorElemento(P)

Evaluación de selectores de objetos

- Utilizando un clasificador
 - Seleccionar un conjunto de bases de datos
 - Utilizar algún método de validación aplicando selección+clasificación
 - Utilizar alguna medida de evaluación de calidad de clasificación

Dr. Jesús Ariel Carrasco Ochoa ariel@inaoep.mx
Oficina 8311