



**Ricardo Aler Mur**

## **CLASIFICADORES KNN-I**

En esta clase se habla del aprendizaje de modelos de clasificación y regresión basados en instancias o ejemplares. En concreto:

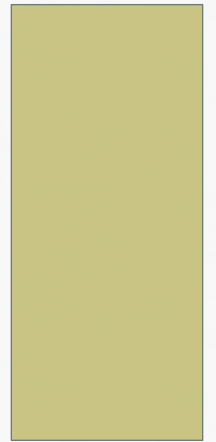
- Se define la clasificación y regresión basada en los vecinos más cercanos (KNN = k-nearest neighbours)
- Se explican algunas de las ventajas (no es necesario construir un modelo explícito) e inconvenientes de KNN, entre ellas, la lentitud en la clasificación cuando hay muchos datos, la dependencia de la función de distancia, la influencia nociva de atributos irrelevantes, la maldición de la dimensionalidad y la gran influencia del ruido en los datos.
- A continuación se muestra como se pueden solventar algunas de esas deficiencias:
  - Cómo seleccionar el parámetro K para combatir el ruido en los datos
  - Cómo seleccionar las instancias o datos más representativos para evitar la lentitud en la clasificación. Se explican las técnicas de edición (para eliminar datos engañosos/ruido) y las de condensación (para eliminar datos supérfluos), así como técnicas híbridas

- Por último, se explican los algoritmos RT1 / RT2 y RT3 que realizan edición y condensación de manera avanzada
- Se termina la clase introduciendo otra clase de algoritmos – los basados en prototipos - que clasifican basándose en la idea de vecindad, pero que eliminan el problema de la lentitud en la clasificación al utilizar como modelo un conjunto de prototipos en lugar del conjunto de datos completo. En concreto, se explica un algoritmo de posicionamiento de prototipos denominado Learning Vector Quantization o LVQ.

# Fases del análisis de datos

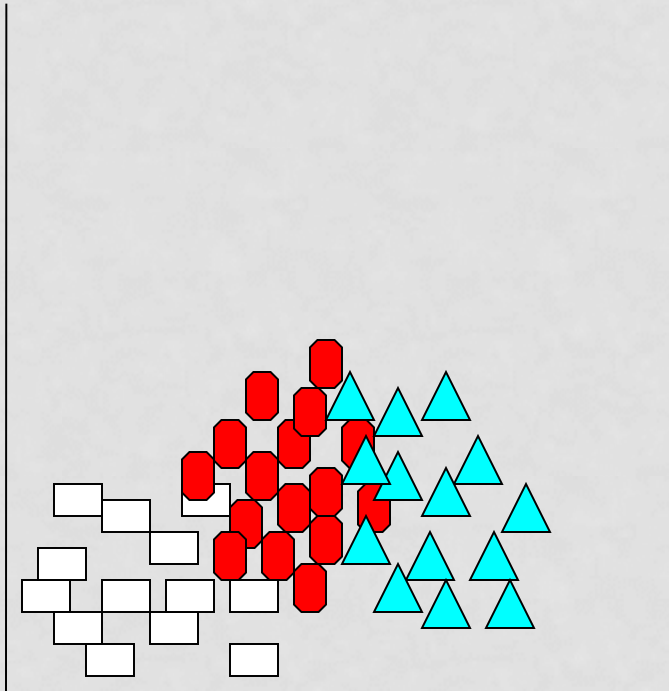
- Recopilación de los datos (tabla datos x atributos)
- Preproceso:
  - De los datos
  - **De los atributos:**
    - **Selección de atributos:**
      - **Ranking**
      - **Subset selection: CFS y WRAPPER**
    - Transformación / Generación de atributos:
      - No supervisada: **PCA, random projections, autoencoders**
      - Supervisada: **mediante redes de neuronas**
- Generación del modelo:
  - **Clasificación: árboles de decisión**
  - **Regresión:**
    - **Modelos lineales**
    - **Árboles de modelos**
- Evaluación: **validación cruzada, matriz de confusión**
- Despliegue y uso del modelo

KNN: VECINO(S) MAS  
CERCANO(S)



# K NEAREST NEIGHBORS (KNN)

Altura



□ Niño

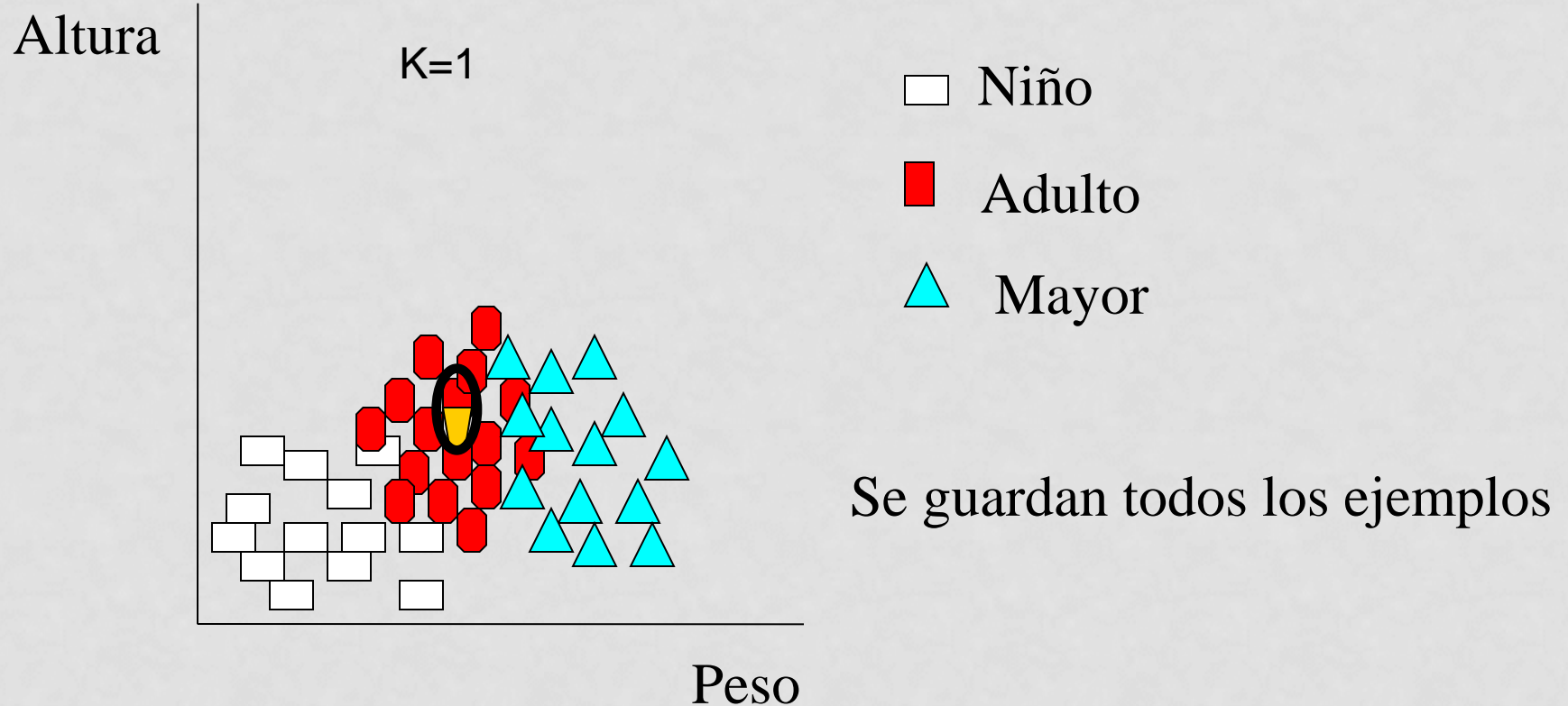
■ Adulto

▲ Mayor

Se guardan todos los ejemplos

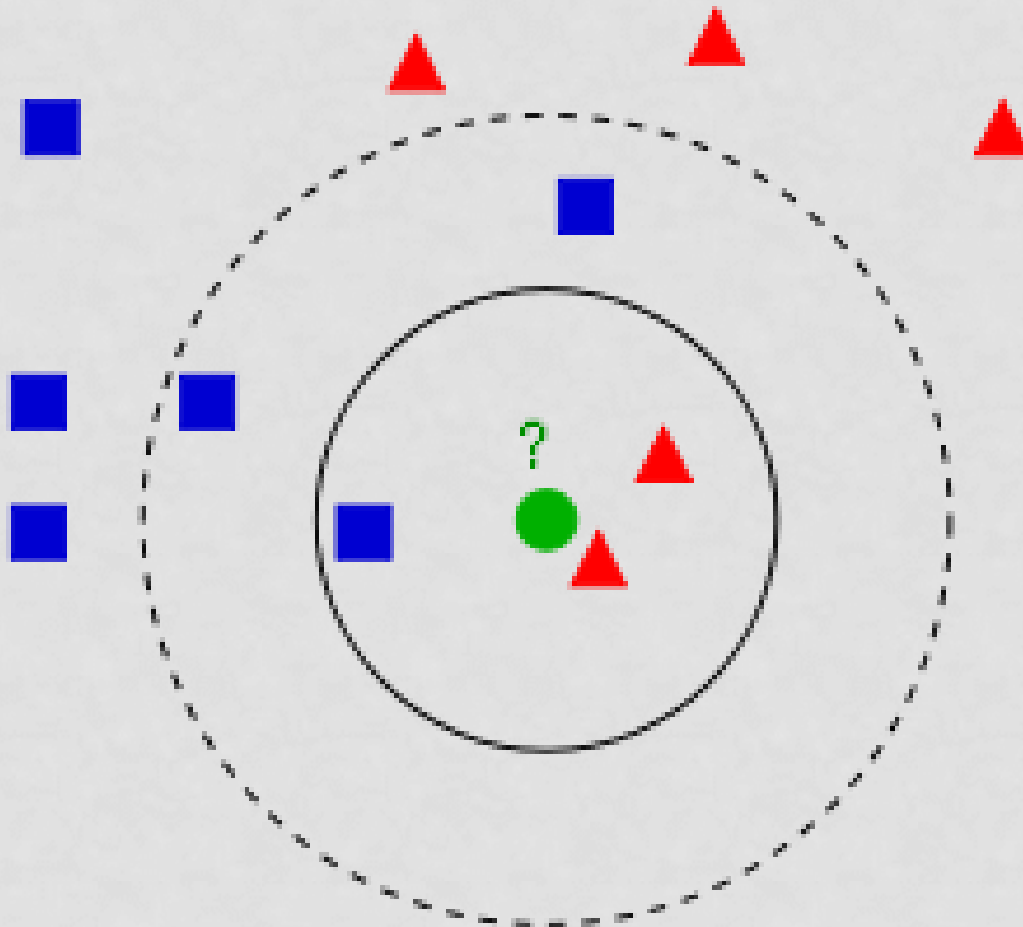
Peso

# K NEAREST NEIGHBORS (KNN)



# K NEAREST NEIGHBORS (KNN)

K=3



# CARACTERISTICAS KNN

- KNN: Clasifica nuevas instancias como la clase mayoritaria de entre los  $k$  vecinos mas cercanos de entre los datos de entrenamiento
- KNN es un algoritmo perezoso (lazy)
  - Durante el entrenamiento, sólo guarda las instancias, no construye ningún modelo (a diferencia de, por ejemplo, los árboles de decisión)
  - La clasificación se hace cuando llega la instancia de test
- Es no paramétrico (no hace suposiciones sobre la distribución que siguen los datos, a diferencia de por ejemplo, un modelo lineal)
  - El mejor modelo de los datos son los propios datos



# CARACTERISTICAS KNN

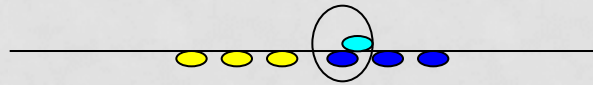
- KNN es local: asume que la clase de un dato depende sólo de los  $k$  vecinos mas cercanos (no se construye un modelo global)
- No necesita adaptación para mas de dos clases (a diferencia de clasificadores lineales)
- Si se quiere hace regresión, calcular la media de los  $k$  vecinos
- Es sencillo

# Problemas / limitaciones KNN

- Muy sensible a los atributos irrelevantes y la maldición de la dimensionalidad
- Muy sensible al ruido
- Lento, si hay muchos datos de entrenamiento ( $O(n*d)$  en almacenamiento y en tiempo)
- Depende de que la función de distancia sea la adecuada
  - Normalmente es la Euclidea:
    - En 2D:  $d(\mathbf{x}_i, \mathbf{x}_j)^2 = (x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2$ 
      - $d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2}$
    - En dD:  $d(\mathbf{x}_i, \mathbf{x}_j)^2 = (x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{id} - x_{jd})^2$

# Atributos irrelevantes

0 atributos irrelevantes

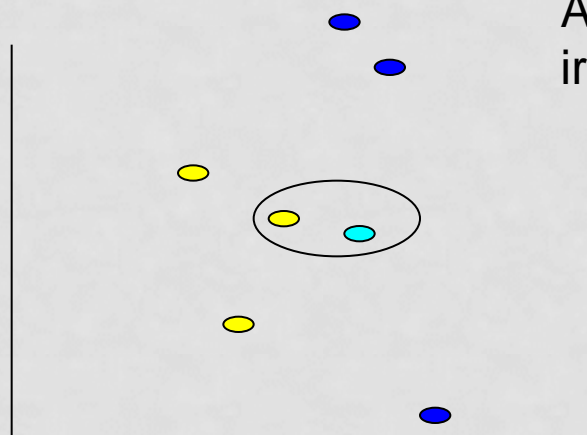


Con el atributo relevante, se clasifica bien

Atributo irrelevante

Atributo irrelevante

1 atributo irrelevante



Con el atributo irrelevante, se clasifica mal (las distancias cambian)

Atributo relevante

Vecino mas cercano  $k=1$

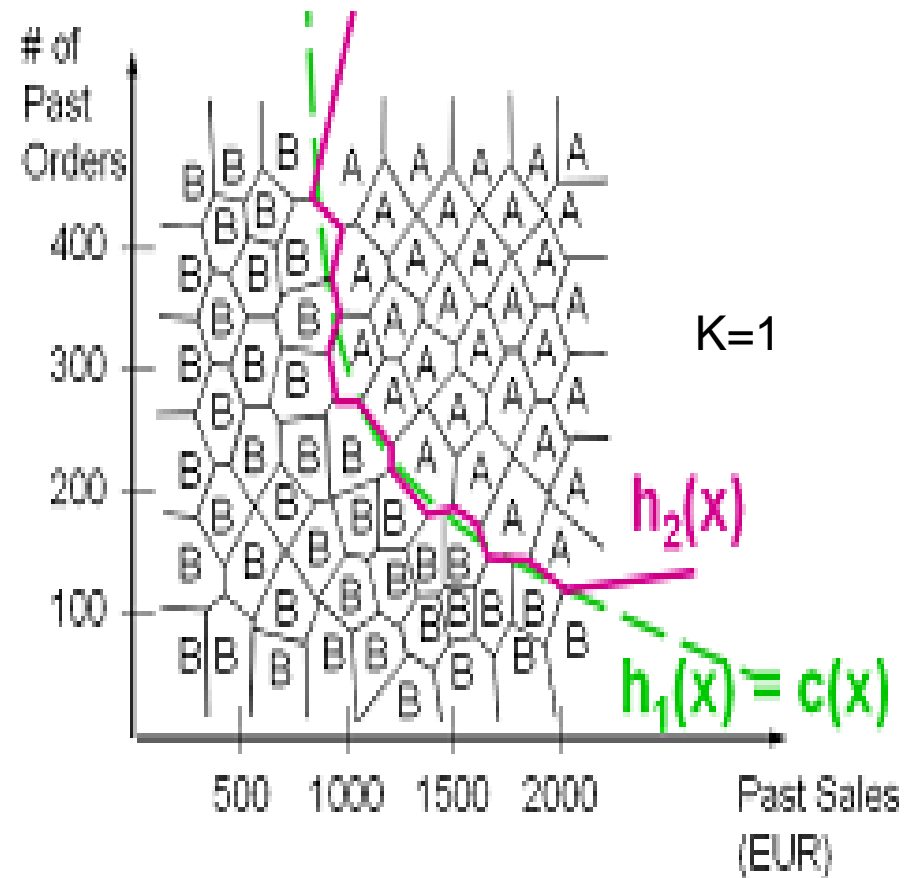
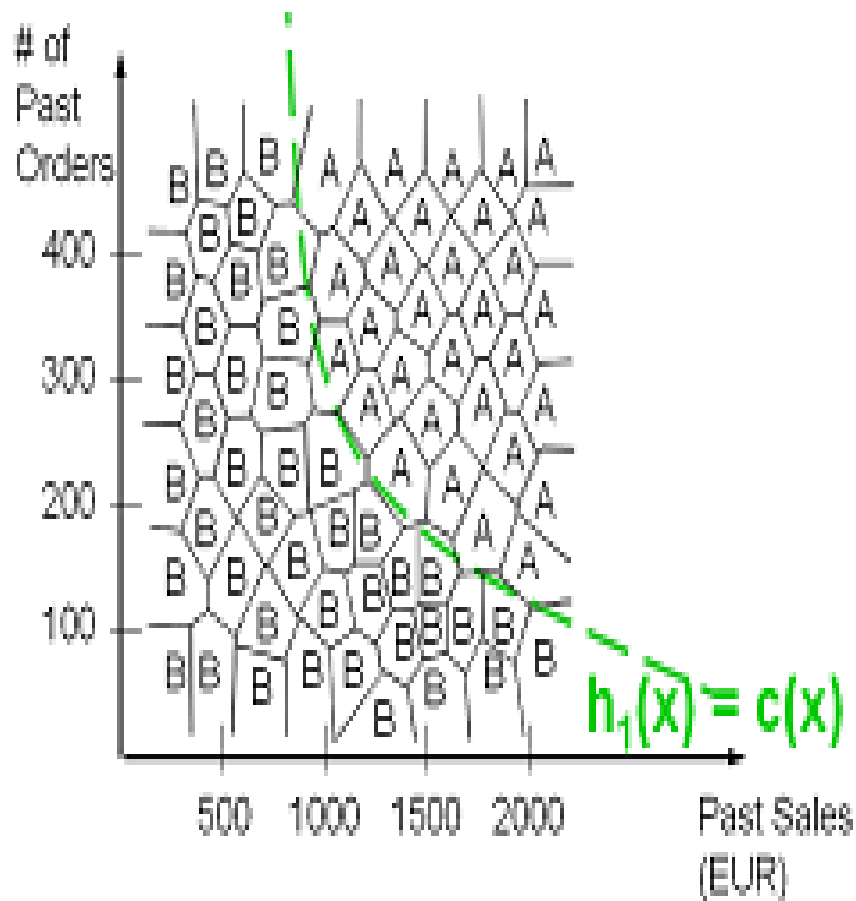
# La maldición de la dimensionalidad

- Le afecta mucho, dado que KNN no crea un modelo global sino que utiliza los datos para representar las fronteras de separación
- Recordemos que el número de datos necesario para representar la superficie de una esfera en  $d$  dimensiones crece exponencialmente con  $d$ :
  - $O(r^d)$

# Normalización

- Si los atributos son nominales, usar distancia de Hamming:
  - Si el atributo  $e$  es nominal, en lugar del componente  $(x_{ie}-x_{je})^2$  se usa  $\delta(x_{ie}, x_{je})$ : 0 si  $x_{ie}=x_{je}$  y 1 en caso contrario
- Es necesario normalizar los atributos para que atributos con mucho rango no tengan mas peso que los demás.  
Normalizaciones:
  - Estandarización:  $x'_{1j} = (x_{1j}-\mu_j)/\sigma_j$
  - Normalización:  $x'_{1j} = (x_{1j}-\min_j)/(\max_j-\min_j)$
  - Es como usar una distancia ponderada:
    - $d(\mathbf{x}_i, \mathbf{x}_j) = m_1*(x_{i1}-x_{j1})^2 + m_2*(x_{i2}-x_{j2})^2 + m_3*(x_{i3}-x_{j3})^2 + \dots + m_d*(x_{id}-x_{jd})^2$

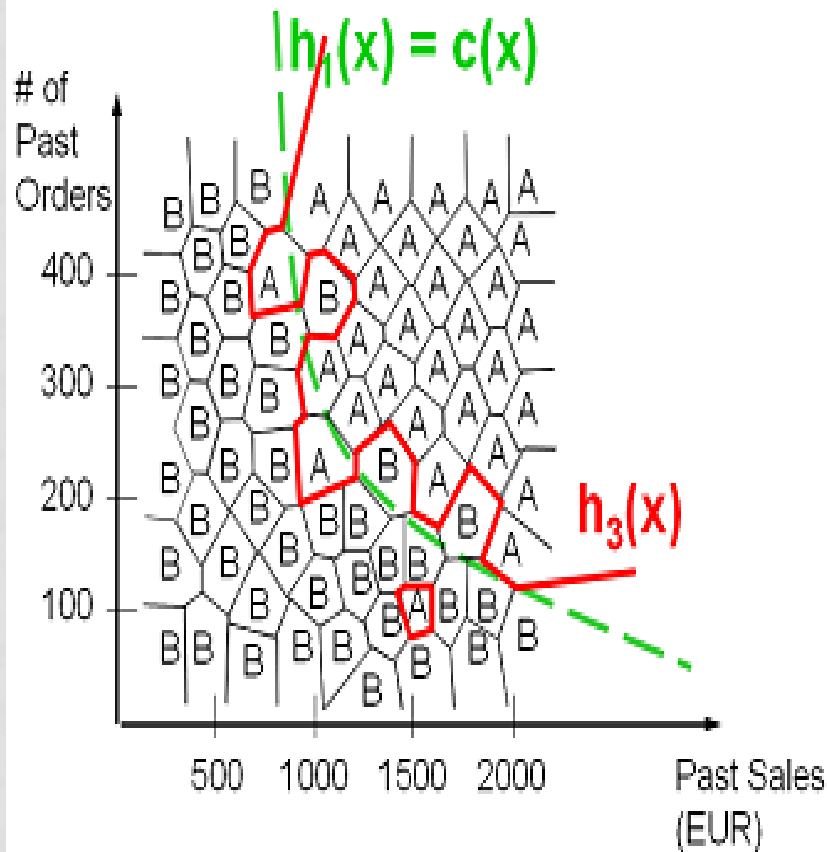
# TESELACIÓN DE VORONOI



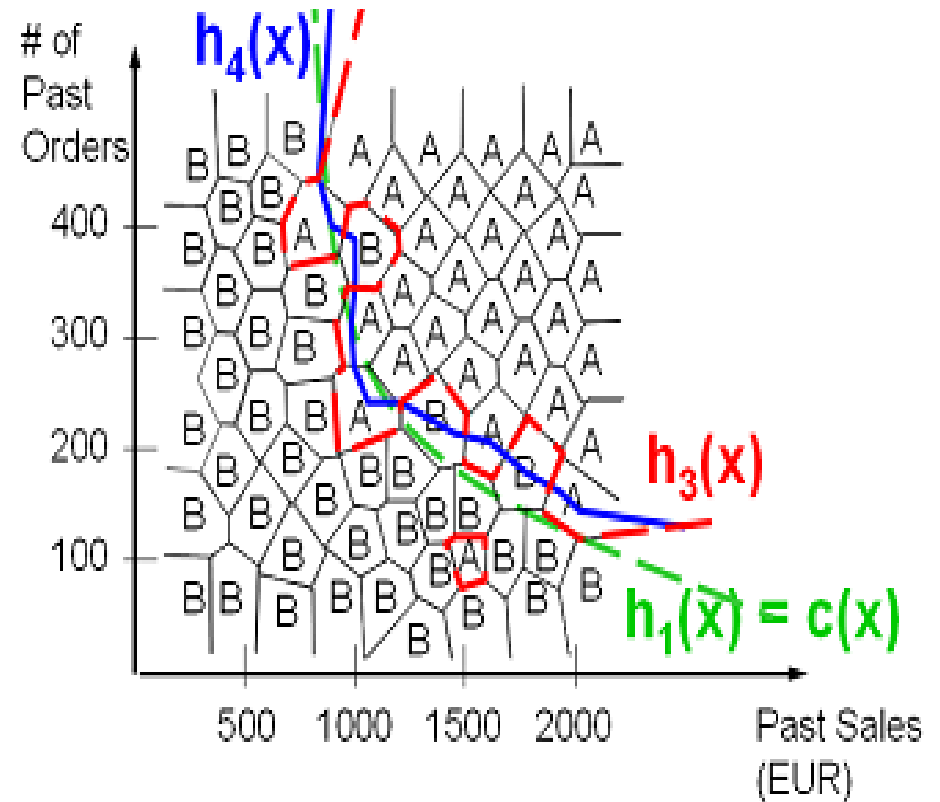
# ¿PORQUÉ USAR K VECINOS?

- Con  $k=1$ , las instancias que son ruido (o solape entre clases) tienen mucha influencia
- Con  $k>1$ , se consideran mas vecinos y el ruido pierde influencia (es como hacer una promediado)
- Si  $k$  es muy alto, se pierde la idea de localidad
  - ¿En que se convierte KNN si  $k$ =número de datos?
- Para evitar que los vecinos lejanos tengan mucha influencia, se puede hacer que cada vecino vote de manera inversamente proporcional a la distancia  $1/d$
- Si hay dos clases, usar  $k$  impar para deshacer empates
- Se puede elegir  $k$  mediante validación cruzada

# ¿PORQUÉ USAR K VECINOS?



(a) 1-NN on noisy data



(b) 3-NN and noisy data



# METODOLOGÍA ANÁLISIS DE DATOS

- Recopilación de los datos (tabla datos x atributos)
- Preproceso:
  - **De los datos: Normalización**
  - **De los atributos:**
    - **Selección de atributos:**
      - **Ranking**
      - **Subset selection: CFS y WRAPPER**
    - Transformación / Generación de atributos:
      - No supervisada: **PCA, random projections, autoencoders**
      - Supervisada: **mediante redes de neuronas**
- **GENERACIÓN DE MODELOS / AJUSTE DE PARÁMETROS / SELECCIÓN DE MODELO**
  - **Clasificación: árboles de decisión, reglas, KNN**
  - **Regresión: modelos lineales, árboles de modelos, KNN**
- Evaluación: **validación cruzada, matriz de confusión**
- Despliegue y uso del modelo

# SELECCIÓN DE K

- A la elección de un modelo para un problema concreto se le denomina selección del modelo (model selection)
- Se puede seleccionar la familia o tipo del modelo:
  - Separador lineal
  - Árboles de decisión
  - KNN
  - ...
- Si se ha seleccionado una familia de modelos, hay que dar valor a ciertos parámetros
  - Árboles de decisión: confidence factor
  - KNN: K
  - ...

# SELECCIÓN DE K

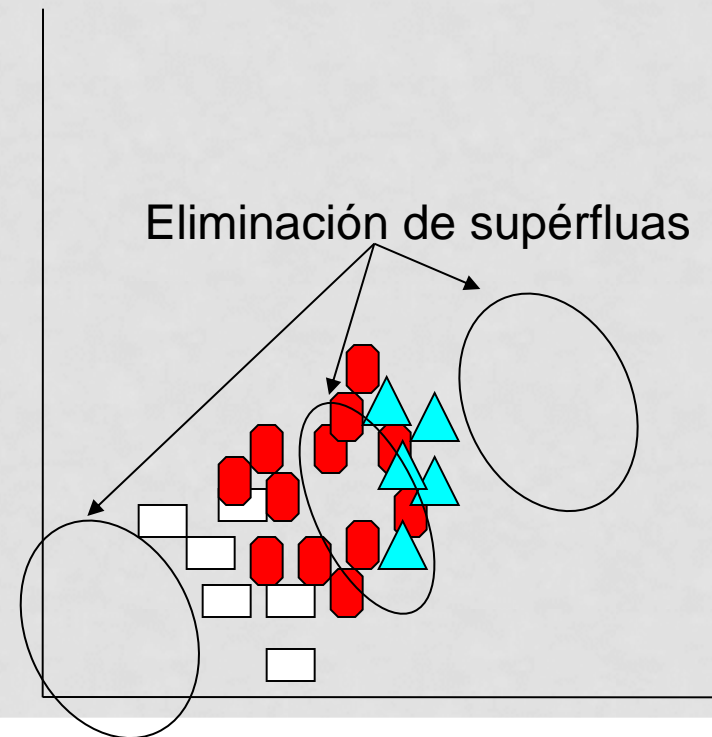
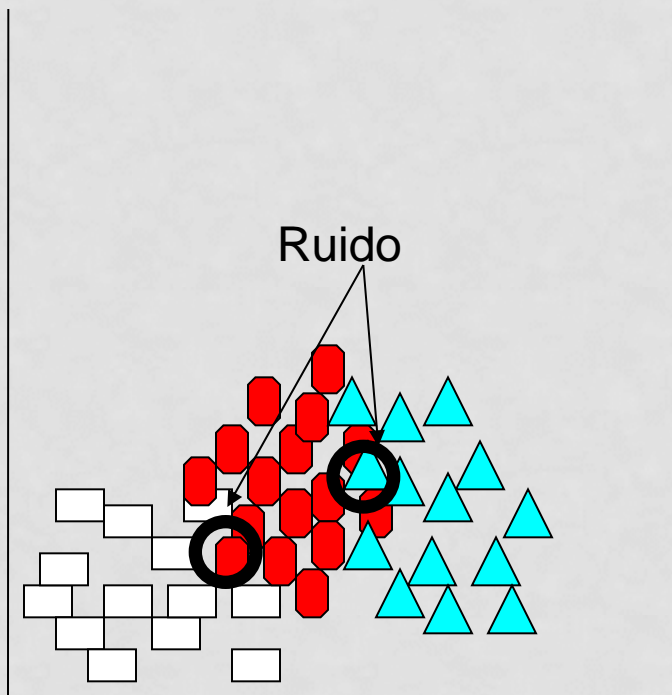
- Se le puede dar valor a  $k$  mediante validación cruzada de los datos de entrenamiento (si es demasiado costoso en tiempo, se puede utilizar un solo conjunto de validación)
- No confundir con la validación cruzada que se hace para evaluar el modelo

# Problemas / limitaciones KNN

- Lento, porque tiene que guardar todos los datos de entrenamiento
  - Solución: eliminación de instancias supérfluas
- Muy sensible al ruido
  - Solución: eliminación de instancias con ruido
- Depende de que la función de distancia sea la adecuada
  - Solución (parcial): aprendizaje de la función distancia
- Muy sensible a los atributos irrelevantes y la maldición de la dimensionalidad:
  - Solución: selección de atributos

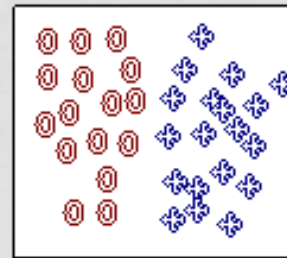
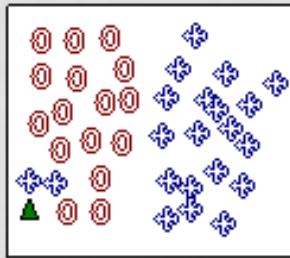
# SELECCIÓN DE INSTANCIAS

- Hay instancias supérfluas: no son necesarias para clasificar. Si las borramos, se decrementará el tiempo de clasificación
- Hay instancias que son ruido (o solape entre clases): confunden al clasificador. Si las borramos, mejorará el porcentaje de aciertos esperado

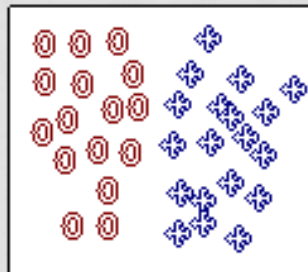


# SELECCIÓN DE INSTANCIAS

Engañosas



Supérfluas

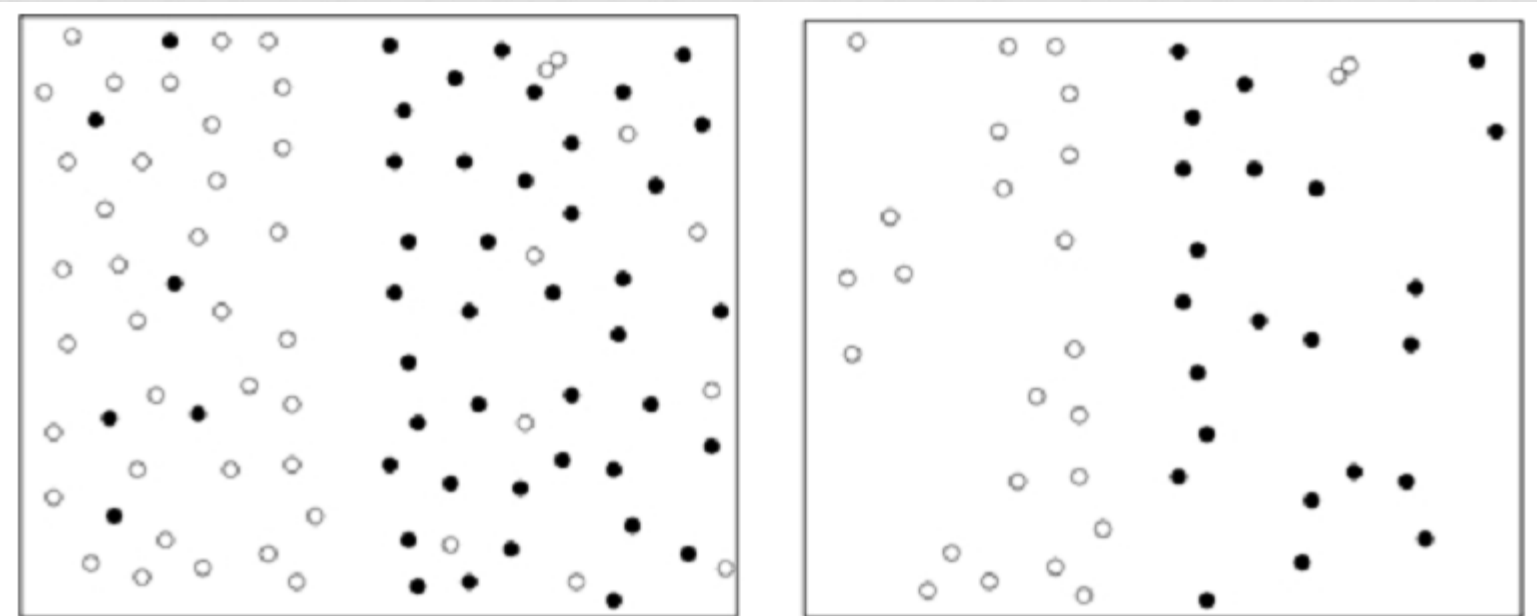


# DOS TIPOS DE TÉCNICAS

- **Editing:** eliminar instancias engañosas (ruido)
  - Típicamente se eliminarán sólo unas pocas
- **Condensación:** eliminar instancias supérfluas:
  - Se pueden llegar a eliminar muchas, las del interior, manteniendo las de la frontera

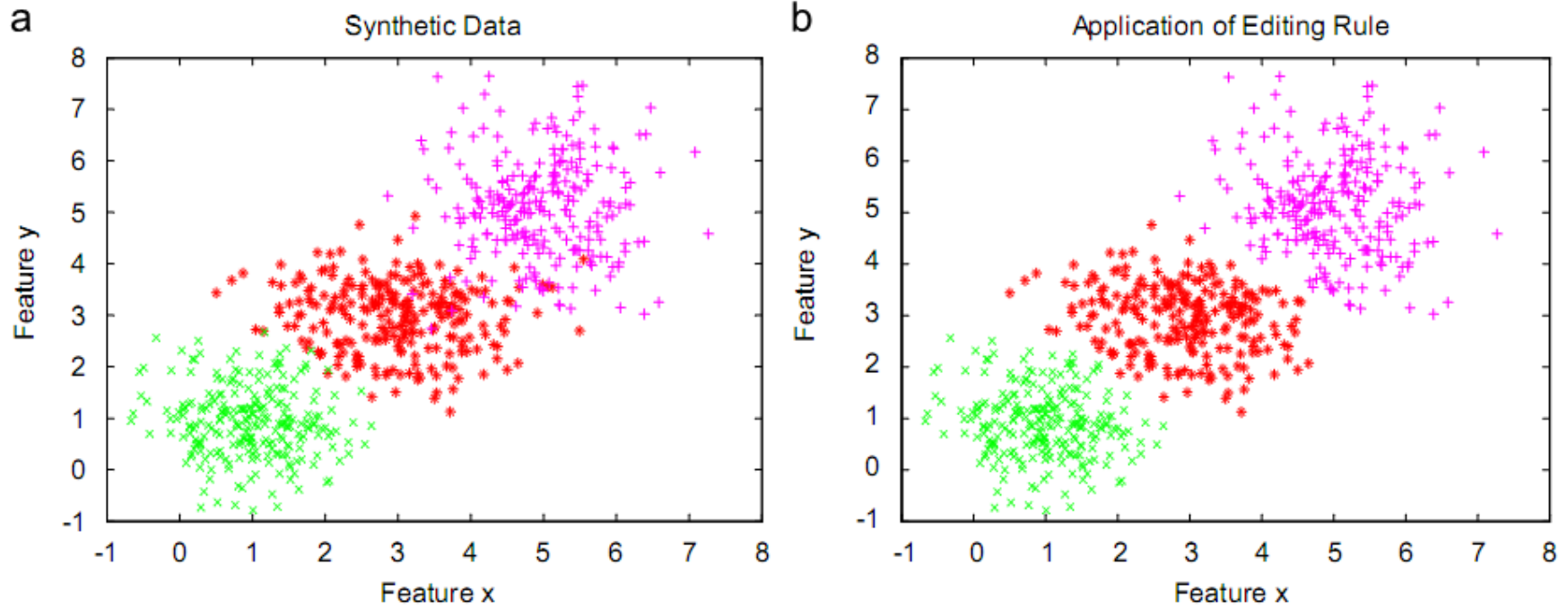
# Wilson editing

- Wilson editing: elimina instancia  $x_i$  si es clasificada incorrectamente por sus  $k$  vecinos:
  - Excepciones en el interior de una clase
  - Algunos puntos en la frontera (suaviza las fronteras)
- Repeated Wilson editing: repite Wilson editing hasta que no se pueda aplicar mas



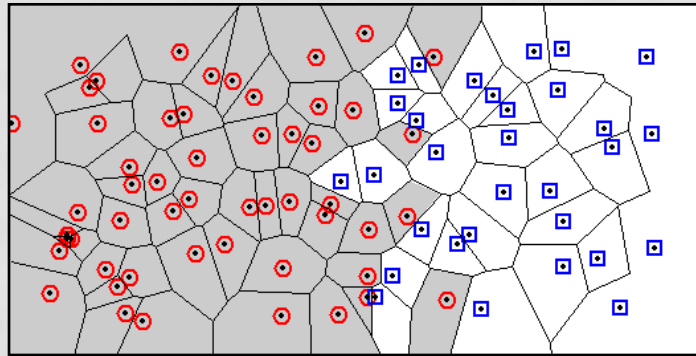


# Wilson editing: ejemplo 2

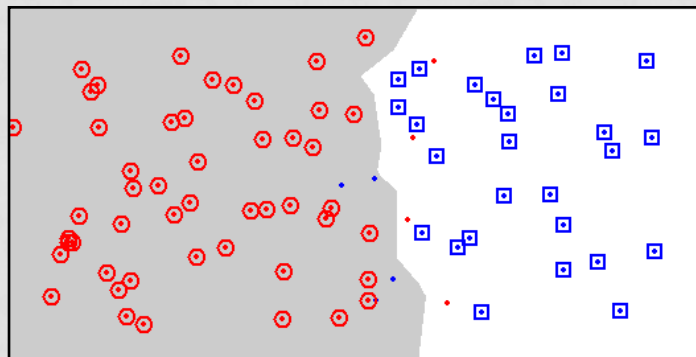


# Wilson editing: ejemplo 3

Overlapping classes



Original data



Wilson editing with  $k=7$

# Wilson editing

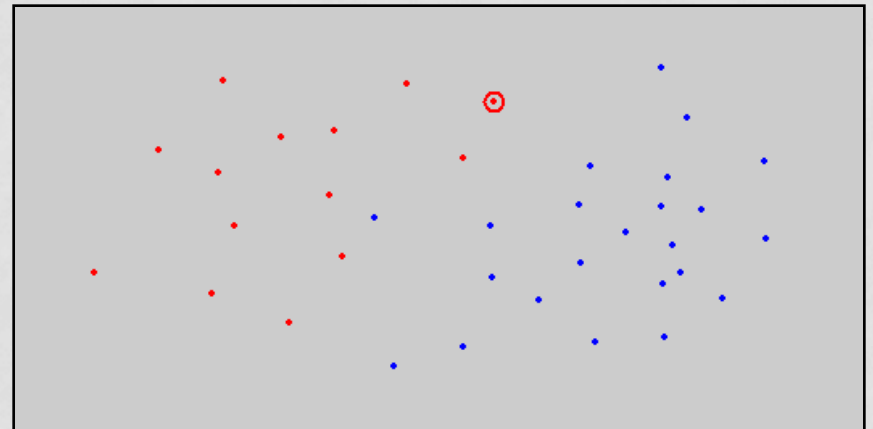
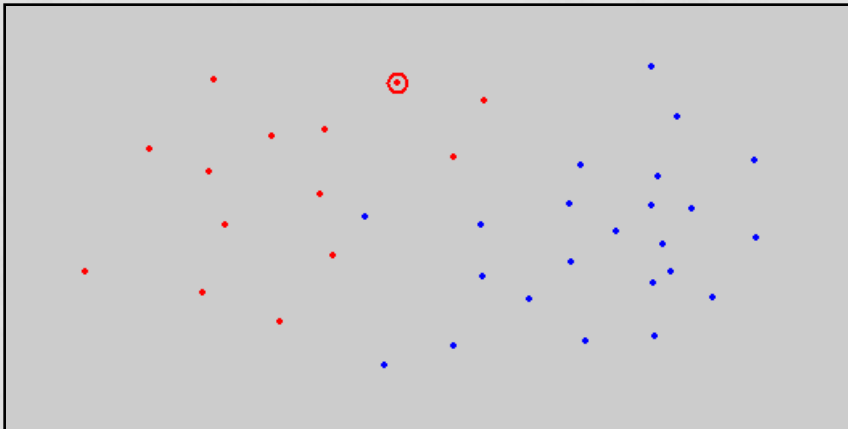
- No quita demasiados datos (es decir, no hay gran mejora en eficiencia)
- Funciona bien si no hay demasiado ruido.
  - Si hay mucho ruido, las instancias con ruido clasificarán bien a otras instancias con ruido

# Condensed Nearest Neighbor (CNN)

- Intenta reducir el número de instancias, eliminando las supérfluas
- Va recorriendo las instancias, y si esa instancia ya está bien clasificada con las que ya hay guardadas, no la guarda. Sólo guarda aquellas que no se clasifican bien con las ya existentes (y por tanto, es crítica, es necesaria)

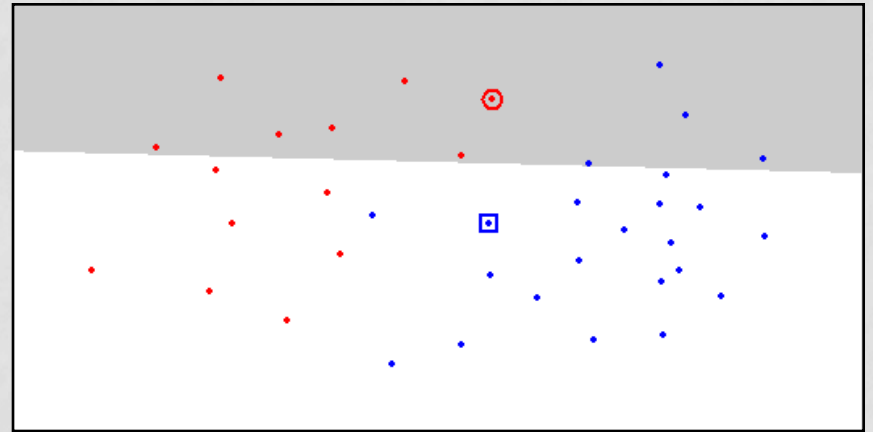
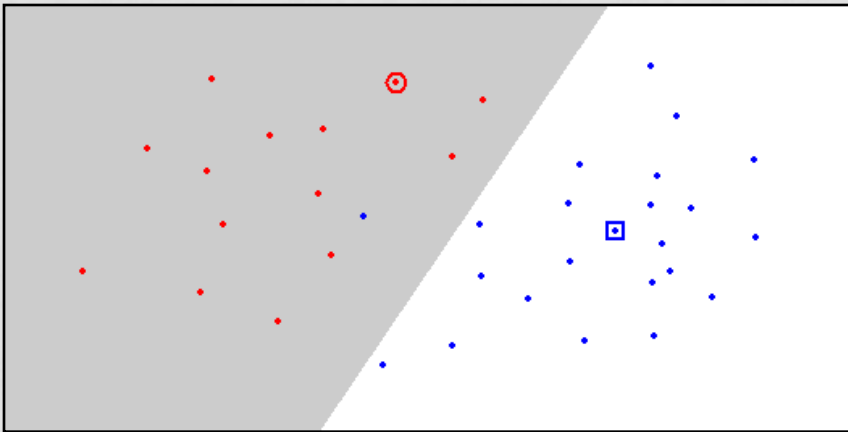
# CNN

1. Inicializa *store* con  $\mathbf{x}_1$
2. Elige un  $\mathbf{x}_i$  fuera de *store* mal clasificado según *store*. Muévelo a *store*
3. Repite 2 hasta que no se muevan mas instancias a *store*
4. Para clasificar con KNN, usar *store* en lugar de los datos originales

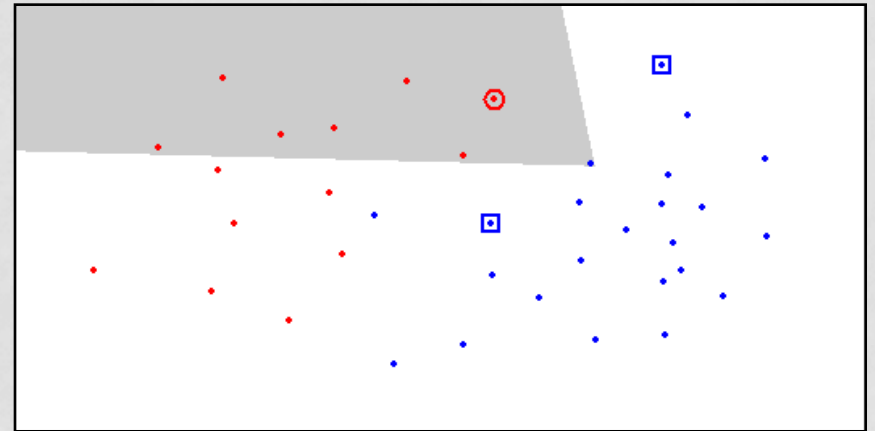
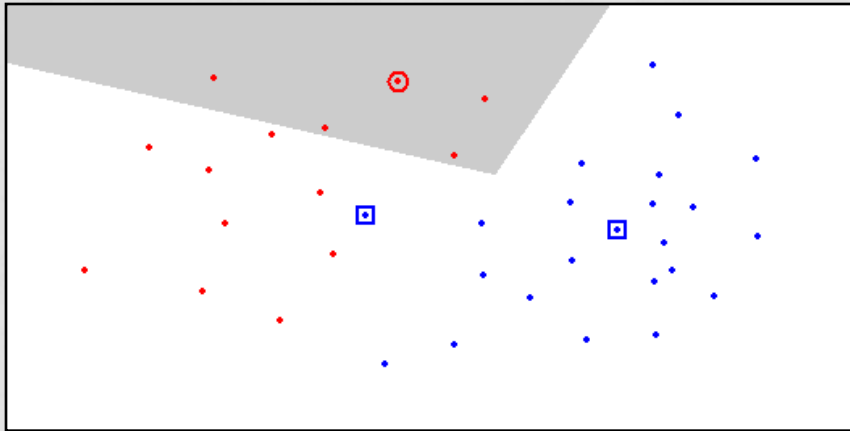


# CNN

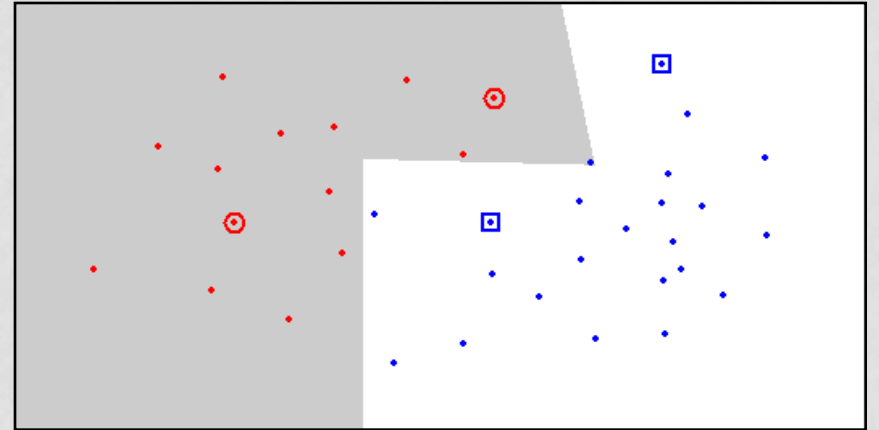
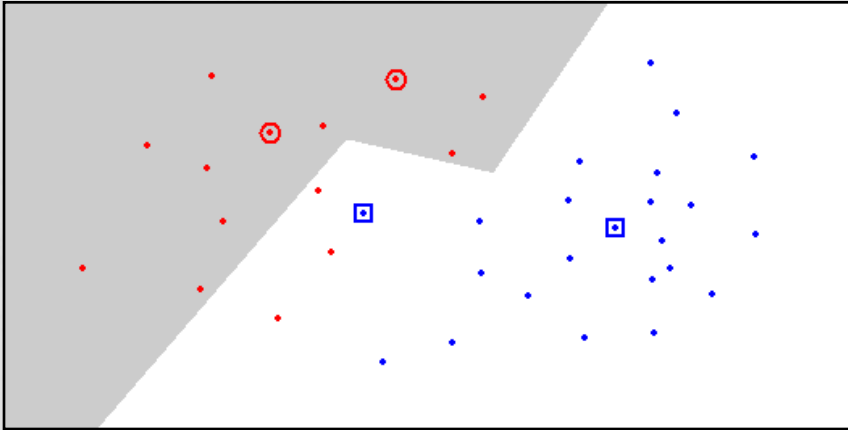
CNN depende mucho del orden en el que se toman las instancias



# CNN

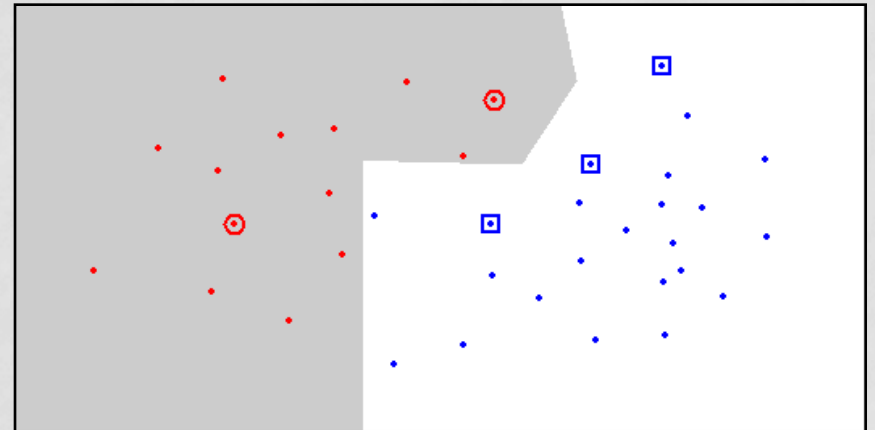
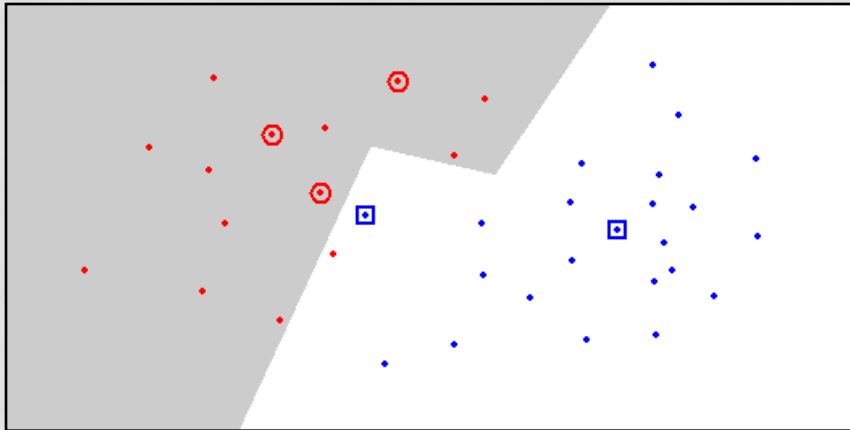


# CNN

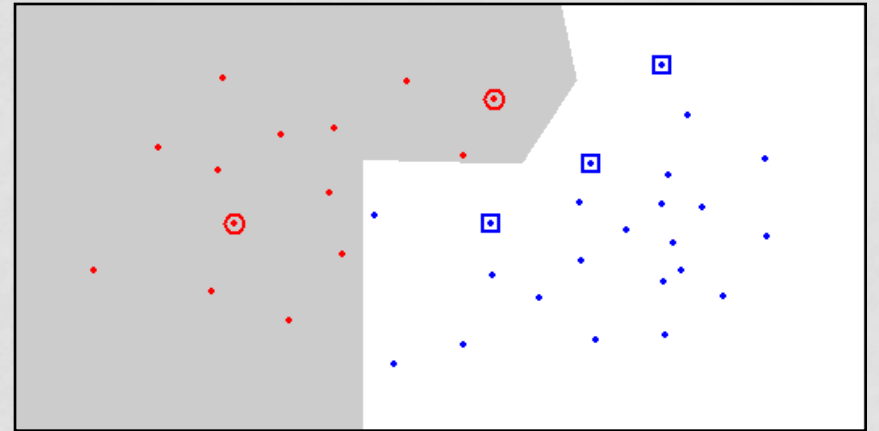
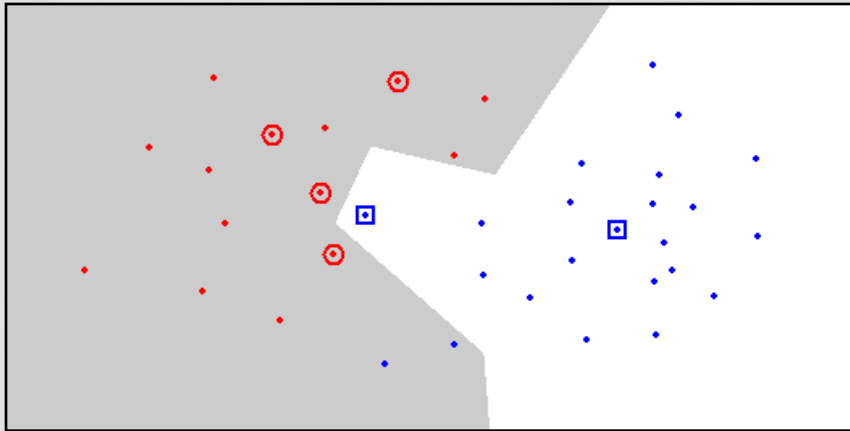




# CNN

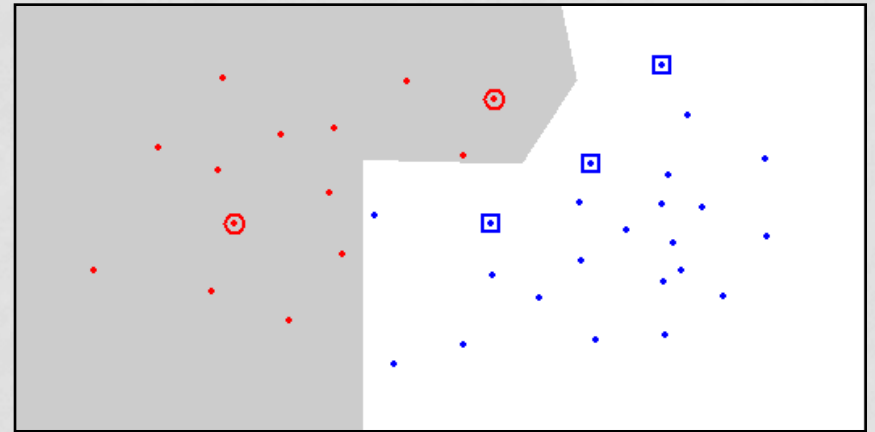
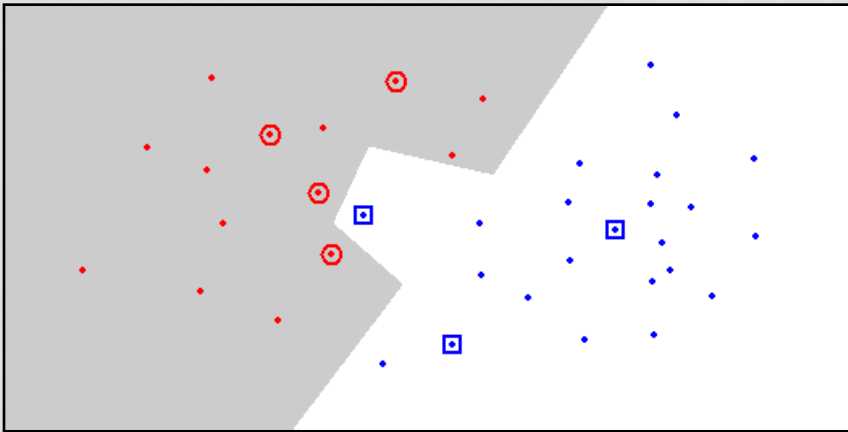


# CNN



# CNN

CNN depende mucho del orden en el que se toman las instancias



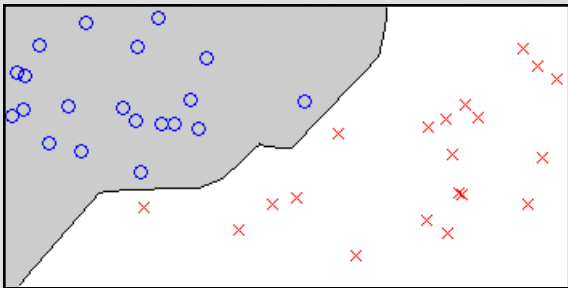
# Características de CNN

- +: positivo, -: negativo
- +: elimina todas aquellas instancias no críticas para la clasificación (reduce mucho la necesidad de almacenamiento)
- -: pero tiende a conservar aquellas instancias con ruido (puesto que son mal clasificadas por las instancias en *Store*)
- -: CNN depende mucho del orden en el que se toman las instancias

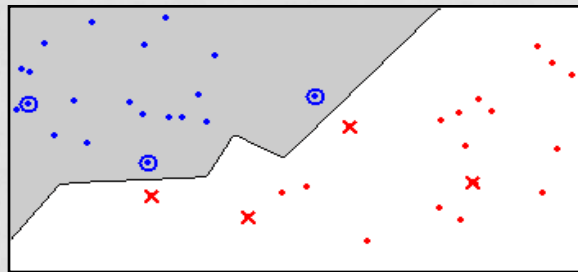
# Reduced Nearest Neighbor rule (RNN)

- Es como CNN, pero comienza con todos los datos y va quitando aquellos que, al quitarlos, no hagan que alguna otra instancia pase a estar mal clasificada.
  - Guarda las instancias críticas / necesarias para una clasificación correcta
  - A diferencia de CNN, permite eliminar instancias ruido (puesto que no contribuyen a clasificar correctamente otras instancias, todo lo contrario)

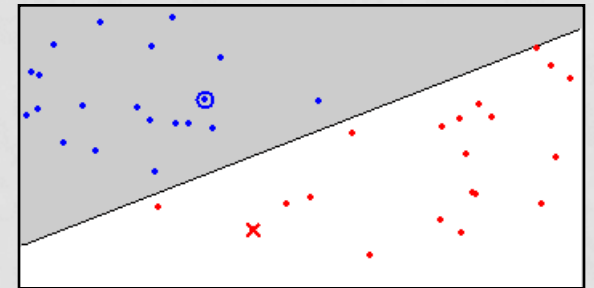
- No hay garantía de que CNN o RNN encuentren el conjunto consistente mínimo



**Datos originales**



**Datos condensados**



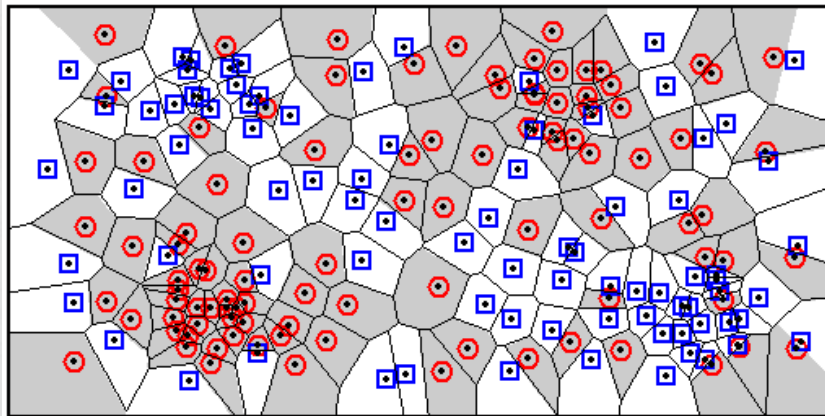
**Conjunto consistente mínimo**

# Edición y condensación

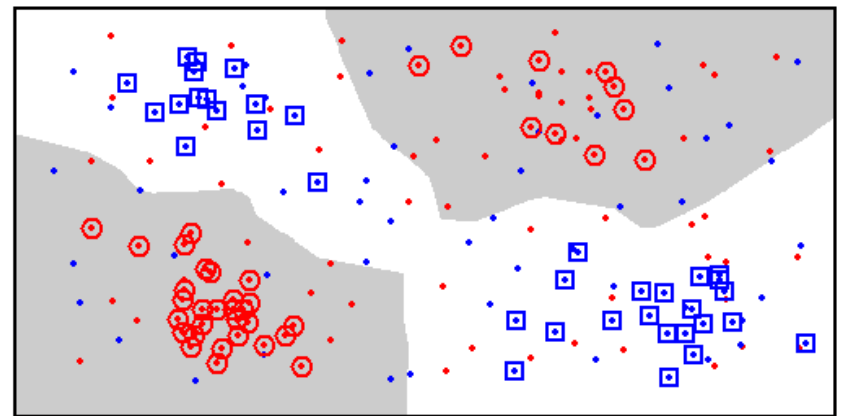
- La edición elimina el ruido y suaviza las fronteras, pero mantiene la mayor parte de los datos (mejora la capacidad de generalización pero no mejora la eficiencia)
- La condensación elimina gran cantidad de datos superfluos, pero mantiene los datos con ruido, puesto que CNN y RSS mantienen aquellos datos clasificados mal por los demás datos (y el ruido tiene siempre esta propiedad)

# Híbridos: 1-Editar 2-Condensar

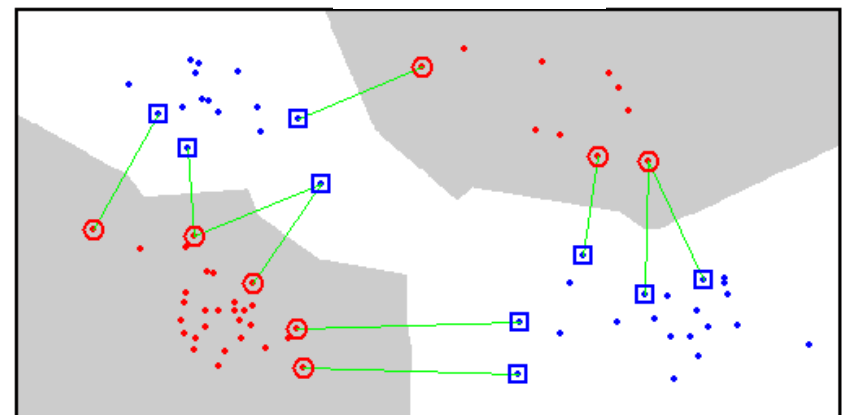
- Primero editar, luego condensar



Original



Edición



Condensación



# Algoritmos avanzados

- Editar / Condensar:
  - RT3:
    - D. Randall Wilson, Tony R. Martinez: Instance Pruning Techniques. ICML 1997: 403-411
  - Iterative case filtering (ICF)
    - Henry Brighton, Chris Mellish: Advances in Instance Selection for Instance-Based Learning Algorithms. Data Min. Knowl. Discov. 6(2): 153-172 (2002)

# RT1 / RT2 / RT3

- Son algoritmos híbridos (condensación + edición)
- Objetivos de RT1 / RT2 / RT3 para mejorar a CNN:
  - Que no dependa del orden (como le ocurría a CNN)
  - Que elimine instancias con ruido
  - Que elimine instancia supérfluas

<http://cgm.cs.mcgill.ca/~athens/cs644/Projects/2004/SumedhaAhuja-EdithLaw/hybrid.html>

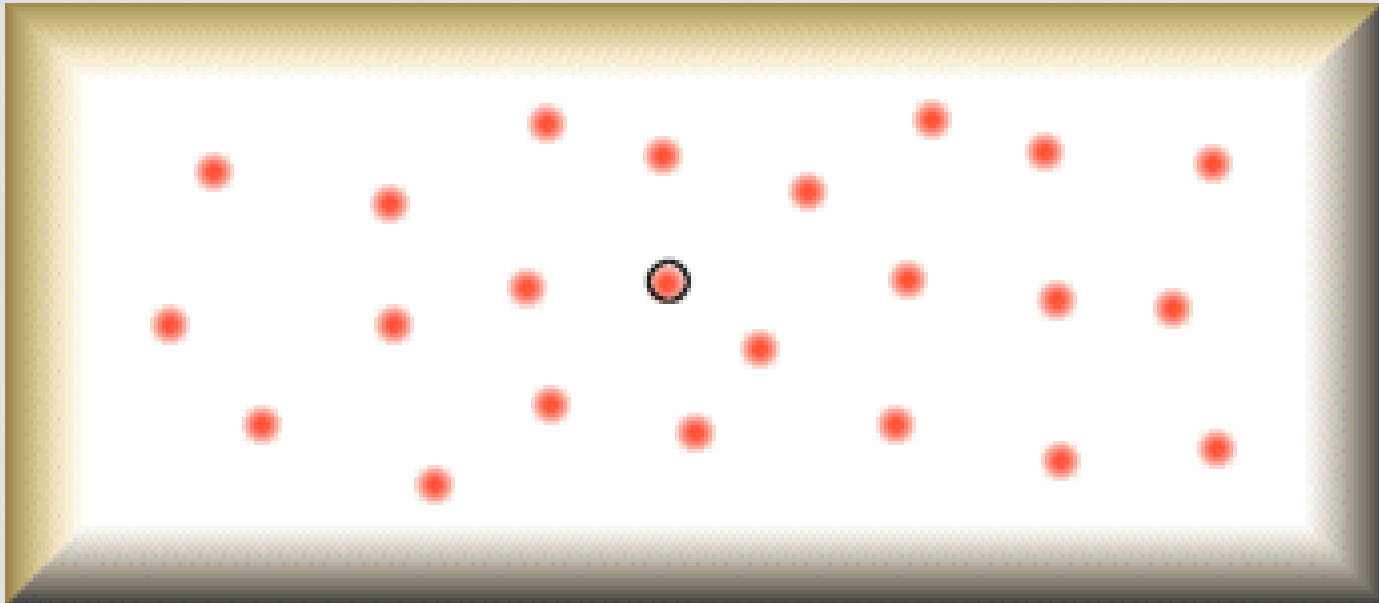
# RT1 / RT2 / RT3

- RT1 está inspirado en RNN (que quitaba una instancia si con eso no hacía que otras instancias pasaran a estar mal clasificadas)
- Para cada instancia P, RT1 calcula sus asociados, Un asociado es una instancia que tiene a P como uno de sus k-vecinos.
  - Es decir una instancia asociada a P es una instancia cuya clasificación puede verse afectada si quitamos P
- RT1 quita una instancia P si el número de asociados clasificados correctamente **sin** P es mayor o igual que los clasificados correctamente **con** P

# RT1

$k=3$

La instancia tiene 6 asociados



# Algoritmo de RT1

```
1  RT1(Training set  $T$ ): Instance set  $S$ .
2    Let  $S = T$ .
3    For each instance  $P$  in  $S$ :
4      Find  $P.N_{1..k+1}$ , the  $k+1$  nearest neighbors of  $P$  in  $S$ .
5      Add  $P$  to each of its neighbors' lists of associates.
6    For each instance  $P$  in  $S$ :
7      Let  $with = \#$  of associates of  $P$  classified correctly with  $P$  as a neighbor.
8      Let  $without = \#$  of associates of  $P$  classified correctly without  $P$ .
9      If  $(without - with) \geq 0$ 
10         Remove  $P$  from  $S$ .
11         Remove  $P$  from its associates' lists of nearest neighbors, and find
12           the next nearest neighbor for each of these associates.
13         Remove  $P$  from its neighbors' lists of associates.
14     Endif
15  Return  $S$ .
```

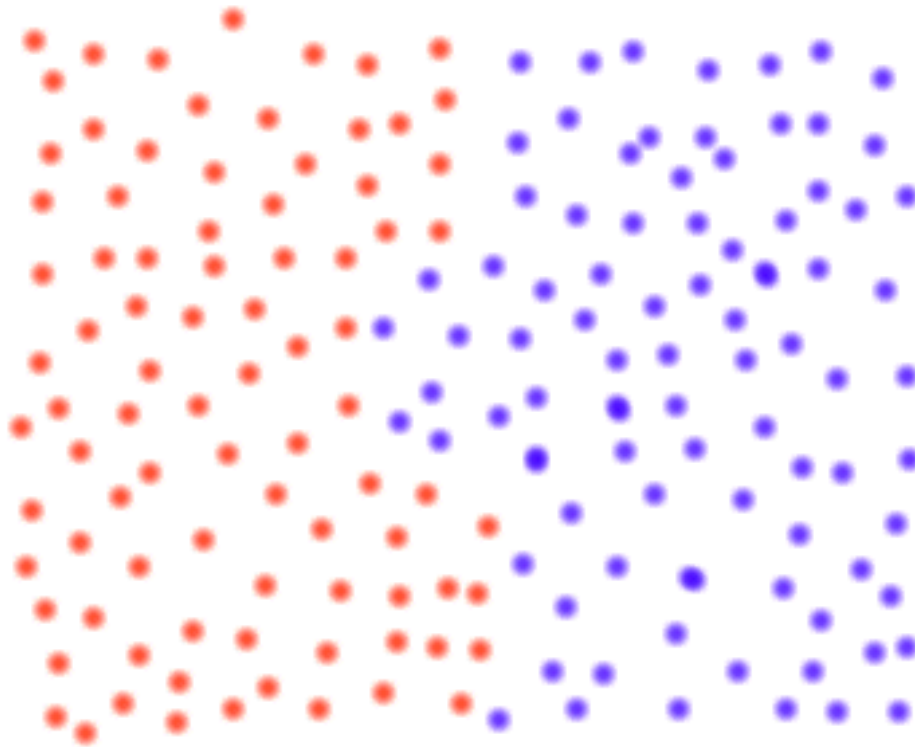
# RT1

- RT1 quita instancias con ruido, puesto que si se las quita eso no perjudica a la clasificación de los asociados, todo lo contrario
- RT1 quita instancias supérfluas en el centro de los clusters, porque en esos lugares todas las instancias son de la misma clase, por lo que quitar vecinos no es perjudicial
- RT1 tiende a guardar instancias no-ruido en la frontera, porque si la quitamos, otras instancias pueden pasar a estar mal clasificadas

## RT2 / RT3

- RT2 intenta quitar las instancias de los centros de los clusters primero
- Para ello, ordena las instancias por distancia a su vecino mas cercano que pertenezca a una clase distinta
- La idea es considerar antes aquellas instancias lejos de la frontera, pero para ello, hay que eliminar aquellas instancias con ruido en el interior de los clusters.
- Así, RT3 hace una primera pasada para quitar las instancias ruidosas (Wilson editing rule)
- RT3 por tanto tiende a conservar las instancias cerca de la frontera

# RT2 / RT3



**Data Set**



# Algoritmo de RT2

```
1  RT1(Training set  $T$ ): Instance set  $S$ .
2    Let  $S = T$ .
3    For each instance  $P$  in  $S$ :
4      Find  $P.N_{1..k+1}$ , the  $k+1$  nearest neighbors of  $P$  in  $S$ .
5      Add  $P$  to each of its neighbors' lists of associates.
6    For each instance  $P$  in  $S$ :
7      Let  $with = \#$  of associates of  $P$  classified correctly with  $P$  as a neighbor.
8      Let  $without = \#$  of associates of  $P$  classified correctly without  $P$ .
9      If  $(without - with) \geq 0$ 
10         Remove  $P$  from  $S$ .
11         Remove  $P$  from its associates' lists of nearest neighbors, and find
12           the next nearest neighbor for each of these associates.
13         Remove  $P$  from its neighbors' lists of associates.
14       Endif
15    Return  $S$ .
```

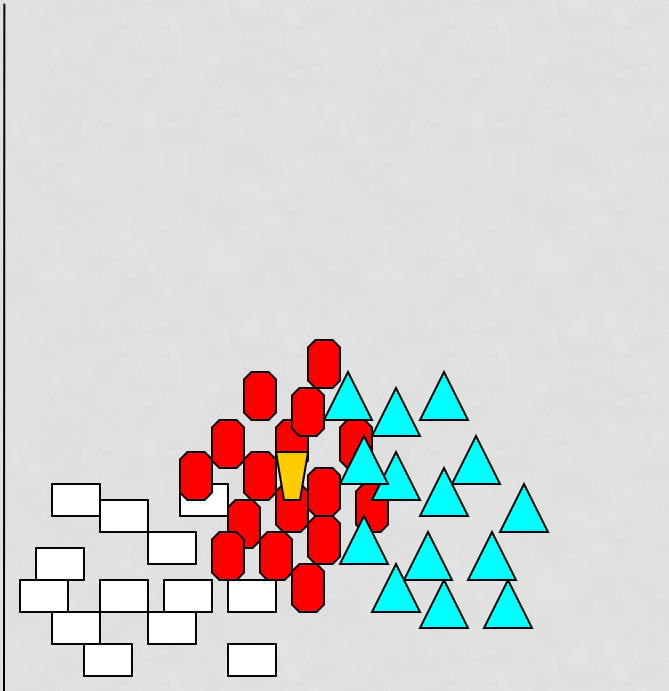
# RESULTADOS

<u>Database</u>	<u>kNN</u> (size)	<u>RT1</u> (size)	<u>RT2</u> (size)	<u>RT3</u> (size)	<u>H-IB3</u> (size)
Anneal	<b>93.11</b> 100	<b>87.85</b> 9.11	<b>95.36</b> 11.42	<b>93.49</b> 8.63	<b>94.98</b> 7.81
Australian	<b>84.78</b> 100	<b>82.61</b> 7.67	<b>84.64</b> 15.41	<b>84.35</b> 5.93	<b>85.99</b> 6.48
Breast Cancer WI	<b>96.28</b> 100	<b>94.00</b> 2.56	<b>96.14</b> 5.79	<b>96.14</b> 3.58	<b>95.71</b> 2.56
Bridges	<b>66.09</b> 100	<b>55.64</b> 20.86	<b>59.18</b> 24.11	<b>58.27</b> 18.66	<b>59.37</b> 38.67
Crx	<b>83.62</b> 100	<b>81.01</b> 6.70	<b>84.93</b> 14.11	<b>85.80</b> 5.46	<b>83.48</b> 6.86
Echocardiogram	<b>94.82</b> 100	<b>93.39</b> 9.01	<b>85.18</b> 7.51	<b>93.39</b> 9.01	<b>93.39</b> 14.85
Flag	<b>61.34</b> 100	<b>58.13</b> 24.51	<b>62.34</b> 32.30	<b>61.29</b> 20.45	<b>51.50</b> 39.18
Glass	<b>73.83</b> 100	<b>60.30</b> 26.11	<b>64.98</b> 31.52	<b>65.02</b> 23.88	<b>67.77</b> 32.92
Heart	<b>81.48</b> 100	<b>79.26</b> 12.96	<b>81.11</b> 21.60	<b>83.33</b> 13.62	<b>76.30</b> 10.33
Heart.Cleveland	<b>81.19</b> 100	<b>77.85</b> 14.26	<b>79.87</b> 20.61	<b>80.84</b> 12.76	<b>74.23</b> 10.78
Heart.Hungarian	<b>79.22</b> 100	<b>78.92</b> 11.38	<b>79.22</b> 15.98	<b>79.95</b> 10.43	<b>74.83</b> 8.88
Heart.Long Beach VA	<b>70.00</b> 100	<b>73.00</b> 11.78	<b>72.00</b> 16.33	<b>73.50</b> 4.22	<b>69.50</b> 11.67
Heart.More	<b>74.17</b> 100	<b>73.20</b> 11.20	<b>74.50</b> 16.98	<b>76.25</b> 9.10	<b>74.75</b> 13.97
Heart.Swiss	<b>92.69</b> 100	<b>91.15</b> 2.08	<b>93.46</b> 2.89	<b>93.46</b> 1.81	<b>84.62</b> 4.79
Hepatitis	<b>80.62</b> 100	<b>76.21</b> 8.67	<b>82.00</b> 13.98	<b>81.87</b> 7.81	<b>72.79</b> 8.03
Horse-Colic	<b>57.84</b> 100	<b>65.09</b> 10.89	<b>66.17</b> 17.98	<b>71.08</b> 7.42	<b>61.82</b> 17.64
Image.Segmentation	<b>93.10</b> 100	<b>84.76</b> 10.21	<b>92.38</b> 13.76	<b>92.62</b> 10.98	<b>90.24</b> 14.79
Ionosphere	<b>84.62</b> 100	<b>84.91</b> 5.67	<b>88.32</b> 12.09	<b>87.75</b> 7.06	<b>88.32</b> 13.61
Iris	<b>94.00</b> 100	<b>89.33</b> 11.70	<b>95.33</b> 16.89	<b>95.33</b> 14.81	<b>92.00</b> 10.96
...					
Average	<b>80.78</b> 100	<b>76.93</b> 14.55	<b>80.09</b> 20.16	<b>80.31</b> 14.32	<b>77.41</b> 16.67

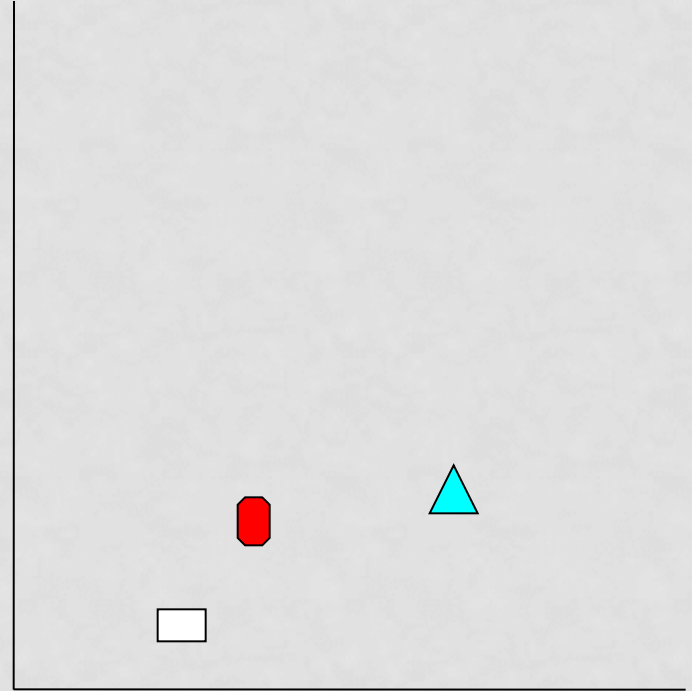
# CLASIFICACIÓN BASADA EN PROTOTIPOS

# PROTOTIPOS

Altura



Altura



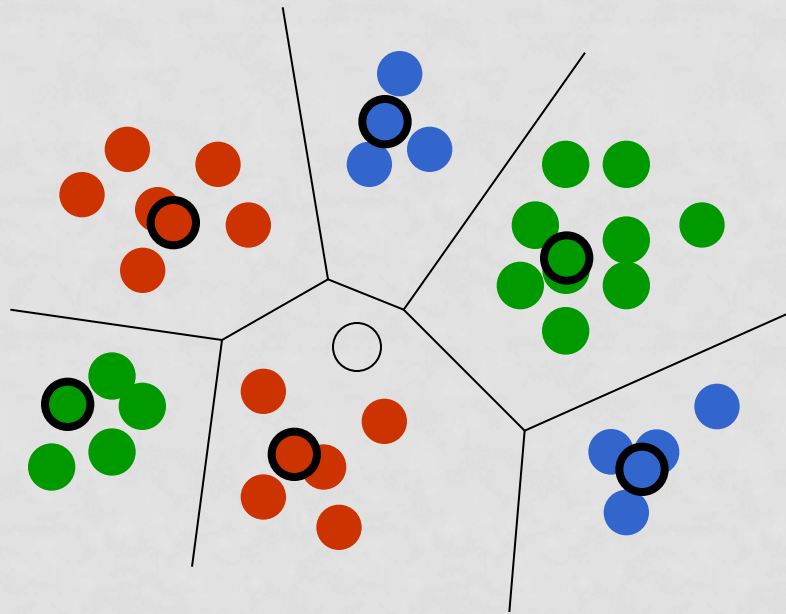
Peso

Peso

Mejora la eficiencia en espacio (sólo se guardan unos pocos prototipos) y en tiempo (se computan muchas menos distancias cuando llega el dato de test)

# Learning vector quantization (LVQ)

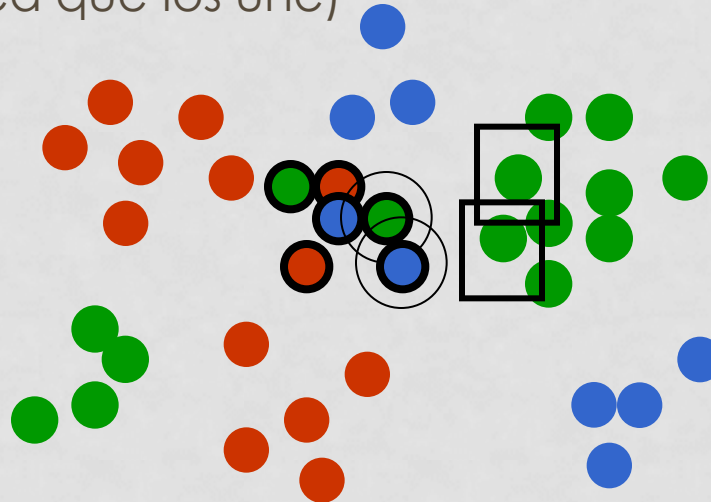
- Cada prototipo tiene una etiqueta
- Se clasifica según la clase del prototipo más cercano (o según sus regiones de Voronoi)



# LVQ 1

## LVQ I entrenamiento:

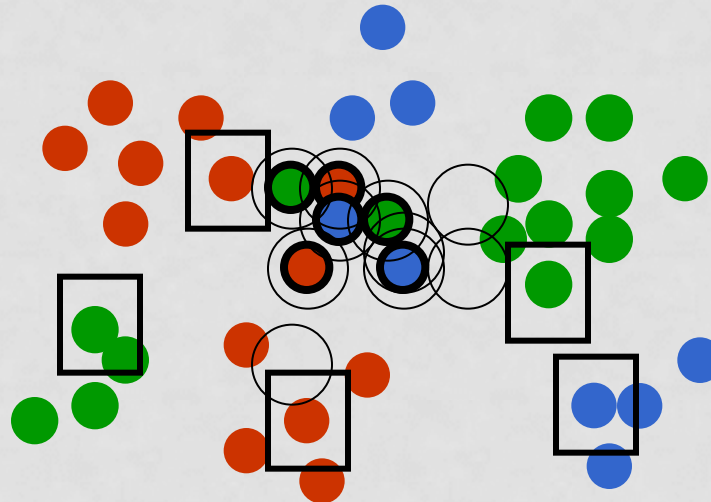
- Inicialización aleatoria
- repetir:
  - presentar instancia de entrenamiento
  - determinar prototipo más cercano
  - Acercarlo o alejarlo de la instancia según la clase (y según la dirección de la línea que los une)



# LVQ 2.1

## LVQ 2.1 entrenamiento:

- Inicialización aleatoria
- repetir:
  - presentar instancia de entrenamiento
  - Determinar:
    - prototipo más cercano correcto
    - prototipo más cercano incorrecto
  - Acercarlo o alejarlo si está dentro de una ventana



# REGLA DE ACTUALIZACIÓN DE LVQ 2.1

$$\begin{aligned}w_i(t+1) &= w_i(t) - \alpha(t)(x - w_i(t)), \\w_j(t+1) &= w_j(t) + \alpha(t)(x - w_j(t)),\end{aligned}$$

Ventana: que no ocurra que un prototipo está muy cerca y otro muy lejos. Se puede forzar haciendo que el mas cercano esté suficientemente lejos ( $d_{\text{cerca}} > s * d_{\text{lejos}}$ ) ;  $s < 1$

$$\min \left( \frac{d_i}{d_j}, \frac{d_j}{d_i} \right) > s,$$