

		Clase Estimada	
		Positiva	Negativa
Clase Real	Positiva	TP <b>A</b>	FN <b>B</b>
	Negativa	FP <b>C</b>	TN <b>D</b>

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$F - \text{Measure} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + FN + FP + TN}$$

$$P(\text{Positiva}) = \frac{TP + FN}{TP + FN + FP + TN} + \frac{TP + FP}{TP + FN + FP + TN}$$

$$P(\text{Positiva}) = \frac{TN + FN}{TP + FN + FP + TN} + \frac{TN + FP}{TP + FN + FP + TN}$$

$$\text{Kapa} = \frac{\text{Accuracy} - P(\text{acuerdo})}{1 - \text{Accuracy M2}}$$

$$P(\text{acuerdo}) = P(\text{Positiva}) + P(\text{Negativa})$$

Las instancias se evalúan individualmente y se compara la predicción con la clase real

- **A TP** → Verdaderos Positivos → instancias que eran de la clase positiva y se predijeron como positivas
- **D TN** → Verdaderos Negativos → instancias que eran de la clase negativa y se predijeron como negativo
- **B FN** → Falsos Negativos → instancias que eran de la clase positiva pero se predijeron como negativas
- **C FP** → Falsos Positivos → instancias que eran de la clase negativa pero se predijeron como positivas
- **Kapa** Mide desacuerdo entre dos estimadores Kapa = 0 Si no hay acuerdo entre los clasificadores que no sea lo que se esperaría por casualidad

## Análisis de ROC

Permite caracterizar la calidad de los clasificadores en base a su rendimiento

Es una representación gráfica de la sensibilidad frente a la especificidad para un sistema clasificador binario según se varía el umbral de discriminación

Proporciona herramientas para seleccionar los modelos posiblemente óptimos y descartar modelos subóptimos independientemente del coste de la distribución

En el punto ( 0 , 0 ) se predice todo clase negativa En el punto ( 1 , 1 ) se predice todo clase positiva En el punto ( 0 , 1 ) Se situaría el mejor método de predicción

- Si TPR = FPR obtenemos un clasificador aleatorio en el que los atributos no aportan información sobre la clase
- Podemos obtener un resultado peor que el del clasificador aleatorio si estamos por debajo de la línea TPR = FPR. En este caso la predicción del clasificador guarda cierta correlación con la realidad, pero la correlación es negativa. La solución es tomar la decisión contraria obteniendo un resultado simétrico sobre la diagonal del espacio ROC
- El mejor clasificador será aquel que tenga mayor área bajo la curva ROC

Razón de Falsos Positivos → eje X

$$FPR = \frac{FP}{FP + TN} \rightarrow (1 - \text{Specificity})$$

Razón de Verdaderos Positivos → eje Y

$$TPR = \frac{TP}{TP + FN} \rightarrow \text{Sensitivity}$$

Para calcular el área de la curva de ROC lo descomponemos en formas trigonométricas básicas y calculamos el área de cada

$$\text{Area Cuadrado} = \text{Base} * \text{altura} \quad \text{Area Triangulo} = \frac{1}{2} * \text{Base} * \text{altura}$$

## Comparación de clasificadores

Se evalúan las instancias por parejas:

- **A 11** → Las instancias pertenecen al mismo clúster y al mismo grupo
- **D 00** → Las instancias No pertenecen al mismo clúster ni al mismo grupo
- **B 01** → Las instancias pertenecen al mismo clúster pero no al mismo grupo
- **C 10** → Las instancias pertenecen al mismo grupo pero no al mismo clúster
- **Coefficiente de correlacion** permite medir lo similares que son dos clasificadores en función del número de aciertos y fallos que cometen.
  - o Los clasificadores que más correlan son aquellos que predijeran siempre lo mismo.  
En este caso B y C serán siempre cero por lo que la ecuación quedaría como  $1/A * D$  quesaría siempre positivo
  - o Los clasificadores que menos correlan son aquellos que predijeran siempre lo opuesto.  
En este caso A y D serian cero por lo que la ecuación quedaría como  $-1/B * C$  quesaría siempre negativo

$$\text{Jaccard} = \frac{A}{A + B + C}$$

$$\text{Fowlkes - Mallows} = \sqrt{\frac{A}{A+B} * \frac{A}{A+C}}$$

$$\text{RandIndex} = \frac{A + D}{A + B + C + D}$$

		Clasificadas en el mismo Clúster	
		Si	No
Pertenecen al mismo grupo experto	Si	11 <b>A</b>	01 <b>B</b>
	No	10 <b>C</b>	00 <b>D</b>

$$\text{coeficiente decorrelacion} = \frac{A * D - B * C}{(A + B) * (D + C) * (A + C) * (B + D)}$$

$$\text{Estadística Q} = \frac{A * D - B * C}{A * D + B * C}$$

$$\text{Desacuerdo} = B + C$$

$$\text{Doble Falta} = D$$

## Orden de complejidad

$$S(N, K) = \frac{1}{K!} * \sum_{i=1}^K [ (-1)^K * \binom{K}{i} * i^N ] \quad \text{siendo} \quad \binom{K}{i} = \frac{K!}{i! * (K - i)!}$$

- N Numero de instancias
- K Numero de subconjuntos
- Casos básicos
  - o 1 Clúster con N instancias solo admite una agrupación. Todas las instancias en dicho clúster **S(N,1) = 1**
  - o N Clusters con N instancias solo admite una agrupación. Una instancia por cada clúster **S(N,N) = 1**
- Caso general **S(N,K) → K \* S(N-1) S( N-1, K-1)**

### Orden Complejidad RedNeuronal

d = inputs H = Neuronas Ocultas K = Salidas

n = Instancias de entrenamiento

$$[ [ W = (d + 1) * H ] + [ b = (H + 1) * K ] ] * n$$

Parametros W = Capa 0 \* Capa 1 \* ... \* Capa i

Bias b = Capa 1 + Capa 2 + ... Capa i

## Clustering Jerárquico

Asume una topología en árbol que define dependencias entre instancias

Descubre agrupamientos con un número de grupos que va desde 1 hasta N

Ofrece relación entre particiones con distinto número de clusters según una jerarquía

- En el nivel inferior hay tantos clusters como instancias cada uno con una sola muestra
- En cada iteración se agrupan los dos clusters más cercanos anotando la distancia en la que se produjo dicha unión
- El algoritmo termina cuando solo queda un clúster

Este algoritmo ofrece una estructura en árbol que permite distintas agrupaciones en función de por donde se realice el corte

## Clustering Particional Algoritmo K-Medias

Consiste en agrupar los datos en un número de conjuntos pre-determinado y posteriormente se van moviendo las instancias de un clúster a otro en base a las medidas de similitud/disimilitud. Necesita como parámetro el número de clusters a predecir

Se trata de un algoritmo de cuantificación vectorial uniforme que pretende obtener el mínimo o máximo sub optimo local de un conjunto de datos:

- Inicialmente elegimos al azar los valores que formaran parte del CodeBook seleccionando tantas instancias como clusters a predecir.
- Asignamos cada instancia a un único clúster comparando las distancias entre dicho clúster y los centroides del CodeBook
- Recorremos los elementos que han sido asignados a cada clúster con la finalidad de recalcular el centroide y actualizar el CodeBook
- Se repite el proceso hasta que la solución converja. Para ello se calcula el error cometido entre esta iteración y la anterior

### Cálculo de las distancias

**Single Link** distancia entre las instancias más cercanas

**Complete Link** distancia entre las instancias más lejanas

**Abaraje Link** distancia entre los centroides de cada clúster

## Autor: Cuesta Alario David

### Técnicas de Evaluación

Exploran la estructura intrínseca de los datos como compromiso de dos métricas:

**Cohesión:** Mide la cercanía de las instancias que pertenecen a cada clúster mediante indicadores estadísticos como:

- **SSE Clúster:** Suma de los errores cuadráticos de cada instancia respecto al centroide del clúster al que pertenecen
- **SSE Partición:** Suma de la cohesión interna de todos los clusters

**Separabilidad:** Mide la distancia inter-grupal que existe entre los distintos clusters

- **SSE Externo:** Suma de los errores cuadráticos de cada instancia respecto a cada uno de los centroides de los clusters a los que No pertenece dicha instancia
- **BSS:** Suma de los errores cuadráticos del centroide de cada clúster con respecto al centroide global de la partición ([centroide de los centroides](#))

$$\text{SSE Externo: } \sum_{\forall \text{ Instancias}} \left[ \sum_{\forall \text{ Clusters} \neq \text{Cluster de la instancia}} d^2(\text{Instancia}, \text{Centroide Cluster}) \right] \quad \text{BSS: } \sum_{\forall \text{ Cluster}} d^2(\text{Centroide Cluster}, \text{Centroide Particion})$$

### Silhouette

Es un indicador que combina la cohesión y la separabilidad en una única variable

- **a** Es la distancia media entre una **instancia** y el resto de instancias de su mismo clúster. Un valor pequeño de **a** indica que la instancia analizada está cerca de las demás instancias de su clúster
- **b** Es distancia media entre más cercana entre la **instancia** y cada uno de los demás clusters. Un valor pequeño de **b** indica que la instancia analizada está lejos de los demás clusters (**Tanto a como b deben ser siempre positivos dado que son cálculos de distancias**)

**Silhouette** Se puede utilizar como indicador para encontrar el número óptimo de clusters. Se calcula la Silhouette para cada posible agrupación y se escoge la agrupación que la maximice.

- o Esta comprendido entre 1 y -1 Si es negativo indica que la agrupación es mala (**poca cohesión y/o separabilidad**)

$$\text{Silhouette} = \frac{b(\text{Instancia}) - a(\text{Instancia})}{\text{MAX}[b(\text{Instancia}), a(\text{Instancia})]}$$

$$a(\text{Instancia}) = \frac{\sum_{\forall \text{ Instancias} \in \text{Cluster}} d^2(\text{Instancia}, \text{Resto de instancias del mismo cluster})}{\text{Numero Instancias del cluster} - 1}$$

$$b(\text{Instancia}) = \frac{\text{MIN}}{\forall \text{ Clusters} \neq \text{Cluster de la instancia}} \left[ \frac{\sum_{\forall \text{ Instancias} \in \text{Cluster}} d^2(\text{Instancia}, \text{Instancia})}{\text{Numero Instancias del cluste}} \right]$$

$$\text{Silhouette}(\text{Instancia}) \longrightarrow \begin{cases} 1 - \frac{a}{b} & \text{Si } a < b \\ 0 & \text{Si } a = b \\ \frac{b}{a} - 1 & \text{Si } a > b \end{cases}$$

### Regresión Lineal

Es bastante habitual comparar los errores con respecto de la media para para probar la desviación del modelo. El **coeficiente de determinación** denominado **R<sup>2</sup>** es un estadístico que permite determinar la calidad de un modelo lineal para replicar los resultados.

- Su principal propósito es predecir futuros resultados o probar una hipótesis.
- Adquiere valores entre 0 y 1.

$$R^2 = \frac{SSR}{SST} \quad \text{SSE} = \sum_{\forall \text{ Instancias}} (\text{Salida Real}(\text{Instancia}) - \text{Salida Estimada}(\text{Instancia}))^2$$
$$SSR = \sum_{\forall \text{ Instancias}} (\text{Salida Estimada}(\text{Instancia}) - \text{Media de todas las instancias})^2$$
$$SST = \sum_{\forall \text{ Instancias}} (\text{Salida Real}(\text{Instancia}) - \text{Media de todas las instancias})^2$$

### Entrenamiento del perceptrón

Consiste en buscar los parámetros internos que mejor se ajusten a la muestra utilizando como entradas los atributos de todas las instancias de entrenamiento:

- Inicializamos los pesos de cada perceptron a valores cercanos a cero
- Se analizan las instancias de una en una en orden aleatorio
  - o Se calcula la salida de la red para la instancia actual y se le aplica la función de activación
  - o Se calcula el error y si no es cero Se actualizan los pesos en base al error del modelo actual con respecto a la salida esperada Con la finalidad de contrarrestar el peso cada instancia para aproximar los resultados futuros a la salida deseada

$$Y = W * X \xrightarrow{f(Y)} \text{Clase}$$

$$W \leftarrow W + \delta * (\text{Salida Real} - \text{Salida Estimada}) * X$$

- Este proceso se repite hasta que la variación del error cruce un umbral admisible

**Sigmoide**  $h(Y) = \frac{1}{1 + e^{-Y}}$   $Y = \begin{cases} 0 & \text{Si } h(Y) < 0.5 \\ 1 & \text{Si } h(Y) \geq 0.5 \end{cases}$   $h(Y_i) = \frac{e^{Y_i}}{e^{Y_1} + e^{Y_2} + \dots + e^{Y_i} + \dots + e^{Y_n}}$

**Tangente hiperbólica**  $h(Y) = 1 - \tanh^2(X)$

**Función signo**  $Y = \begin{cases} 1 & \text{Si } Y \geq 0 \\ -1 & \text{Si } Y < 0 \end{cases}$

### Asociación

**Support**  $\longrightarrow$  Probabilidad conjunta. Numero de sucesos en los que se dan X e Y al mismo tiempo

**Confidence**  $\longrightarrow$  Probabilidad condicionada. Probabilidad de que suceda Y sabiendo que ha sucedido X

Support Normalizado en función del número de apariciones de la frecuencia del antecedente

**Lift**  $\longrightarrow$  Mide la dependencia entre dos sucesos (**sucesos independientes si Lift = 1**)

Desglosamos los ítem-sets en dos sub-conjuntos donde un sub-conjunto sea el antecedente y el otro el consecuente

Realizaremos una poda de los ítems en función del Support: Trataremos de analizar aquellos ítems que tengan suficiente Support porque si el ítem aparece pocas veces cualquier regla de asociación resultante tiene una alta probabilidad de ser fruto del azar.

- Para reducir el coste computacional se determina un umbral de modo que si no se supera se corta el proceso de análisis:
  - o Si una asociación no cumple al asociar otro elemento cumplirá menos o igual
  - o Si dos atributos independientes no están asociados al combinarlos lo estarán menos o igual

Después de la poda analizaremos el Confidence para seleccionar los ítems que más apariciones tengan teniendo en cuenta sus apariciones globales. Con esto tenemos el número de veces que ha aparecido el ítem XY en relación al número de veces que ha aparecido el elemento X.

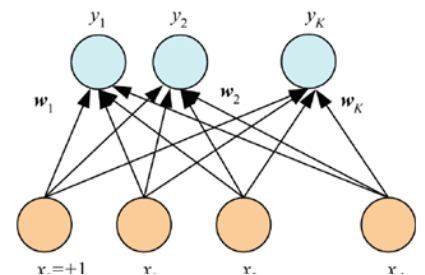
### Combinación de clasificadores

Se basan en los resultados de varios clasificadores débiles para obtener un clasificador más robusto.

- Da lugar a distintos modelos-base haciendo inferencia de distintos conjuntos de entrenamiento
  - Dada una instancia se le asigna la clase más votada por los distintos modelos-base
- La diferencia entre ambas es que Boosting asigna un peso al resultado de cada clasificador en función de su confianza.

$$\varepsilon_j \longrightarrow B_j = \frac{\varepsilon_j}{1 - \varepsilon_j} \xrightarrow{\text{Weighted error}} w_j = \log_2 \left( \frac{1}{B_j} \right) \quad (1 - \text{Accuracy})$$

$$\sum_{\forall \text{ Instancias} \in \text{Cluster}} d^2(\text{Instancia}, \text{Centroide Cluster})$$
$$\sum_{\forall \text{ Clusters}} \text{SSE}_{\text{Clúster}}$$



$$\begin{pmatrix} W_{11} & W_{12} & W_{13} \\ W_{21} & W_{22} & W_{23} \\ W_{31} & W_{32} & W_{33} \end{pmatrix} * \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix} = \begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \end{pmatrix} \xrightarrow{f(Y)} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$
$$Y_i = 1 * b_i + X_1 * W_{1-i} + X_2 * W_{2-i} + \dots + X_n * W_{n-i}$$

To	From	X <sub>1</sub>	X <sub>2</sub>	X <sub>n</sub>	b
Y <sub>1</sub>		W <sub>1-1</sub>	W <sub>2-1</sub>	W <sub>n-1</sub>	b <sub>1</sub>
Y <sub>2</sub>		W <sub>1-2</sub>	W <sub>2-2</sub>	W <sub>n-2</sub>	b <sub>2</sub>
Y <sub>k</sub>		W <sub>1-k</sub>	W <sub>2-k</sub>	W <sub>n-k</sub>	b <sub>n</sub>

**Support**  $(x \Rightarrow y) = P(X, Y)$

**Confidence**  $(x \Rightarrow y) = P(Y|X) = \frac{P(X, Y)}{P(X)}$

**Lift**  $(x \Rightarrow y) = \frac{P(X, Y)}{P(X)P(Y)}$

Se dispone de un conjunto de L modelos-base:

- Media aritmética:  $y_i(x) = \frac{1}{L} \sum_{j=1}^L d_{ji}(x)$
- Máximo:  $y_i(x) = \max_{j=1}^L d_{ji}(x)$
- Producto:  $y_i(x) = \prod_{j=1}^L d_{ji}(x)$