

Tema 3: Aprendizaje Automático en Máquinas

Universidad Pontificia de Salamanca

Manuel Martín-Merino

Contenido

- Aprendizaje estadístico: Introducción
- Aprendizaje estadístico: Sesgo-Varianza
- Problemas dentro del aprendizaje automático
- Análisis discriminante: Modelos lineales
- Análisis discriminante probabilístico
- Redes Neuronales: Modelos no lineales
- Estimación de errores y comparación modelos
- Tópicos avanzados
- Resumen: Discusión

Aprendizaje estadístico: Introducción (I)

- Sea $\{(x_i, y_i)\}_{i=1}^N$ una muestra finita de objetos donde $x_i \in \mathcal{X} \subset \mathbb{R}^d$ es la representación vectorial del objeto i y y_i es la salida óptima.
- El problema de **aprendizaje automático** trata de encontrar una función $f : \mathcal{X} \rightarrow \mathbb{R}$ que permita *predecir la salida para nuevas observaciones* con error mínimo.
- $y_i \in \mathbb{R}$ en problemas de predicción mientras que $y_i \in \{\pm 1\}$ en problemas de clasificación.

Aprendizaje estadístico: Introducción (II)

Matemáticamente, dada una muestra de datos (\mathbf{x}_i, y_i) independientes e idénticamente distribuidos según $p(\mathbf{x}, y)$ desconocida, tratamos de encontrar f .

Esto requiere:

- Proponer una familia de funciones aproximadoras suficientemente flexible $f(\mathbf{x}; w)$.
- Proponer una función de error $L(y, f(\mathbf{x}; \omega))$ que mida la calidad de la aproximación.
- Minimizar el error promedio $E(L(y, f(\mathbf{x}; \omega)))$.

Aprendizaje estadístico: Introducción (III)

Objetivo: Aprender la función óptima $f(\mathbf{x}; \omega)$ que minimiza:

$$E(w) = \int L(y, f(\mathbf{x}; \omega)) p(\mathbf{x}, y) d\mathbf{x} dy \quad (1)$$

como $p(\mathbf{x}, y)$ es desconocida, el riesgo funcional se aproxima por el empírico más un término proporcional a la complejidad del modelo:

$$E(\omega) = \frac{1}{n} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i; \omega)) + \lambda \Omega(\omega) \quad (2)$$

Aprendizaje estadístico: Introducción (IV)

Objetivo: Encontrar una función $f(x_i; \omega)$ que obtenga un balance entre error de entrenamiento y complejidad del modelo.

- $f(x_i; \omega)$ muy compleja implica error entrenamiento 0 pero $\Omega(\omega)$ grande.
- $f(x_i; \omega)$ muy sencilla implica minimizar $\Omega(\omega)$ a costa de error de entrenamiento grande.

Aprendizaje estadístico: Sesgo-Varianza (I)

La complejidad del modelo debe buscar un balance entre sesgo-varianza.

Consideramos que existen varios conjuntos de entrenamiento \mathcal{D} con n patrones y obtenidos de la misma distribución de probabilidad $p(\mathbf{x}, y)$.

Por ser la muestra finita $f_{\mathcal{D}}(\mathbf{x}; \omega)$ dependerá de la muestra particular elegida. Por ello el error se calcula como un promedio para todas las muestras de entrenamiento

$$\mathcal{E}_{\mathcal{D}}\{f_{\mathcal{D}}(\mathbf{x}; \omega) - g(\mathbf{x})\}^2 \quad (3)$$

Aprendizaje estadístico: Sesgo-Varianza (II)

Este error se puede descomponer como:

$$\mathcal{E}_{\mathcal{D}}[\{f_{\mathcal{D}}(\mathbf{x}; \omega) - g(\mathbf{x})\}^2] = \underbrace{\{\mathcal{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x}; \omega)] - g(\mathbf{x})\}^2}_{\text{sesgo}^2} + \underbrace{\mathcal{E}_{\mathcal{D}}[\{f_{\mathcal{D}}(\mathbf{x}; \omega) - \mathcal{E}_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x}; \omega)]\}^2]}_{\text{varianza}} \quad (4)$$

- *Sesgo*: Mide la desviación de la función sobre el valor óptimo del promedio de las estimaciones para diferentes conjuntos de entrenamiento.
- *Varianza*: Mide la sensibilidad del estimador de la muestra a la muestra particular elegida.

Aprendizaje estadístico: Sesgo-Varianza (III)

Ejemplo: Consideramos el problema de regresión $y = (\sin 2\pi x)^2 + \epsilon$, donde ϵ es una variable aleatoria de media 0 y $\sigma_\epsilon = 0,2$. El tamaño del conjunto de entrenamiento es $n = 30$

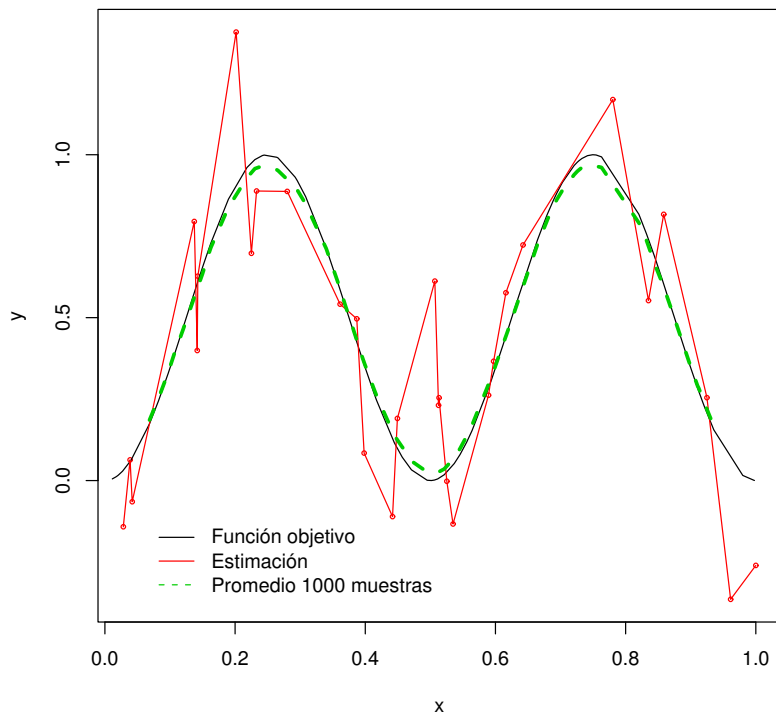


Fig. 1: Grados de libertad 30. Varianza elevada.

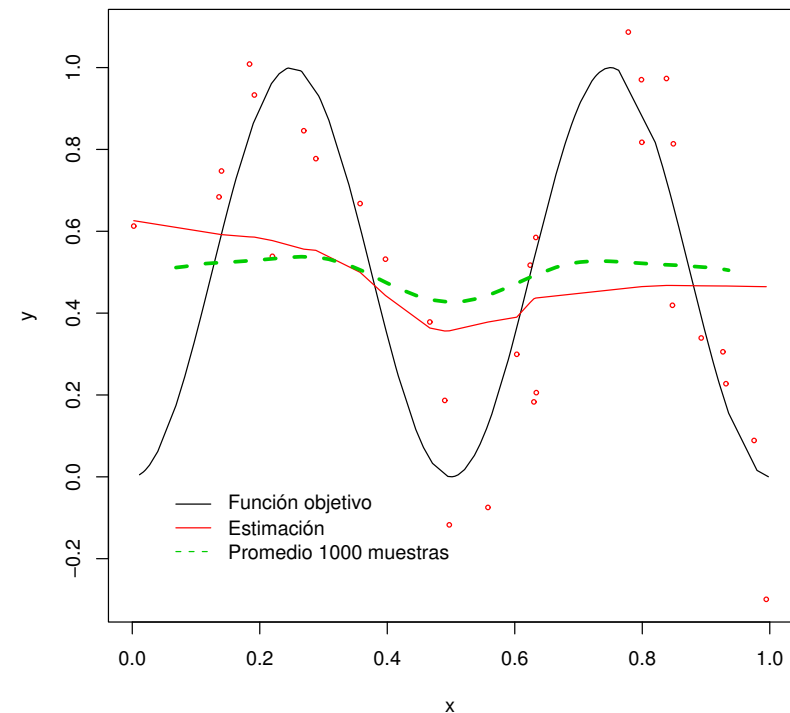
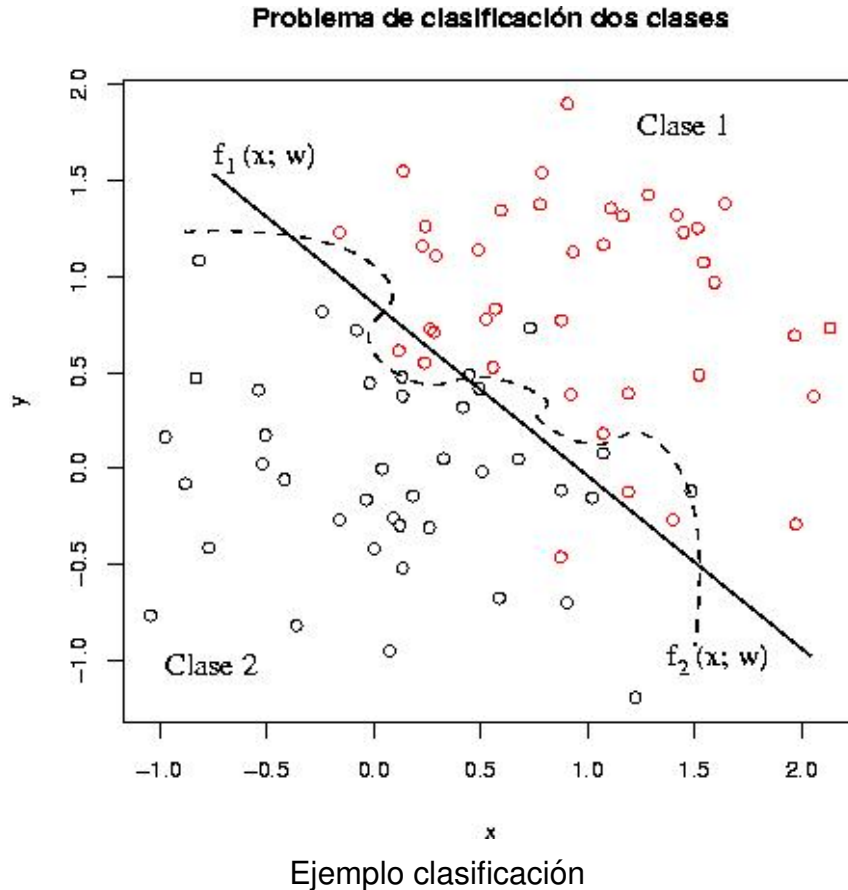


Fig. 2: Grados de libertad 2,39. Sesgo elevado.

Problemas dentro del aprendizaje automático (I)

Clasificación



- Proponer una familia de funciones (por ej. polinomios)
- Proponer una función de error (por ej. cuadrático)
- Estimar los parámetros:
Optimizar

Clasificación

Objetivo: Estimar la $p(C_j|\mathbf{x})$ para nuevos patrones no utilizados en la fase de entrenamiento.

Enfoques:

- Probabilístico. Estima $p(\mathbf{x}|C_j)$ y aplica la regla de decisión Bayesiana:

$$\mathbf{x} \in C_k \quad \text{si} \quad p(\mathbf{x}|C_k)p(C_k) > p(\mathbf{x}|C_j)p(C_j) \quad \forall j \neq k \quad (5)$$

- Redes Neuronales y regresión logística. Estima directamente $p(C_j|\mathbf{x})$
- Análisis discriminante. Estima directamente la ecuación de la frontera $f(\mathbf{x}; \omega)$

Problemas dentro del aprendizaje automático (II)

Predicción: Identificación funciones

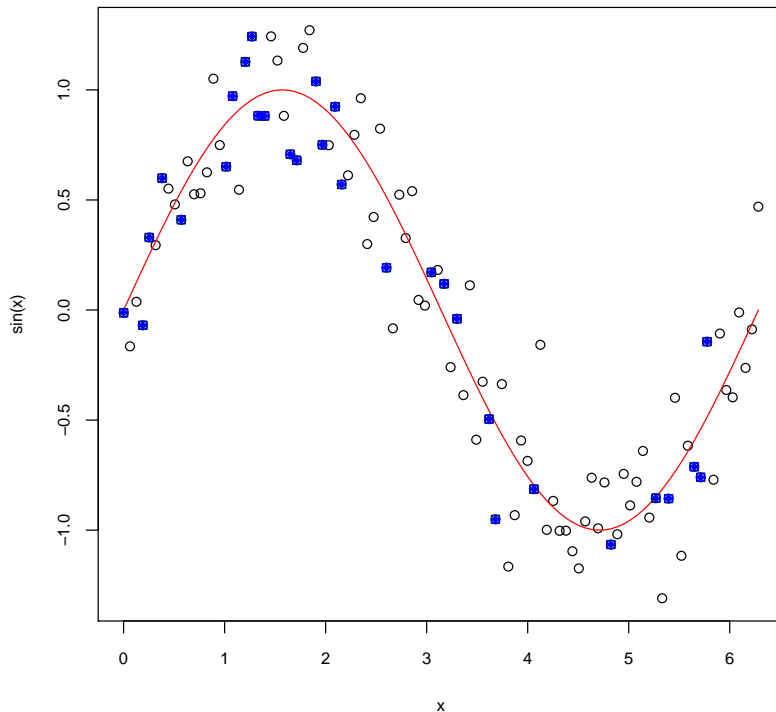
- Proponer una familia de funciones (por ej. polinomios)

- ▷ Modelos lineales: Regresión lineal, logística.
- ▷ Modelos no lineales: Redes Neuronales.

- Proponer una función de error (por ej. cuadrático)

- Estimar los parámetros:
Optimizar

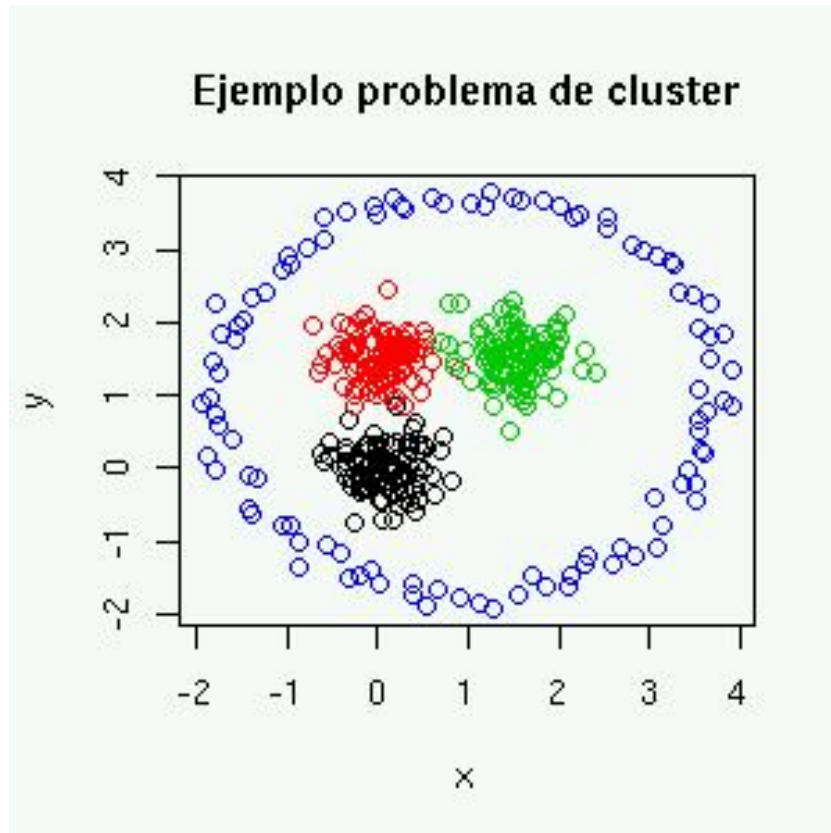
Ejemplo predicción



Ejemplo identificación funciones

Problemas dentro del aprendizaje automático (III)

Análisis de Cluster

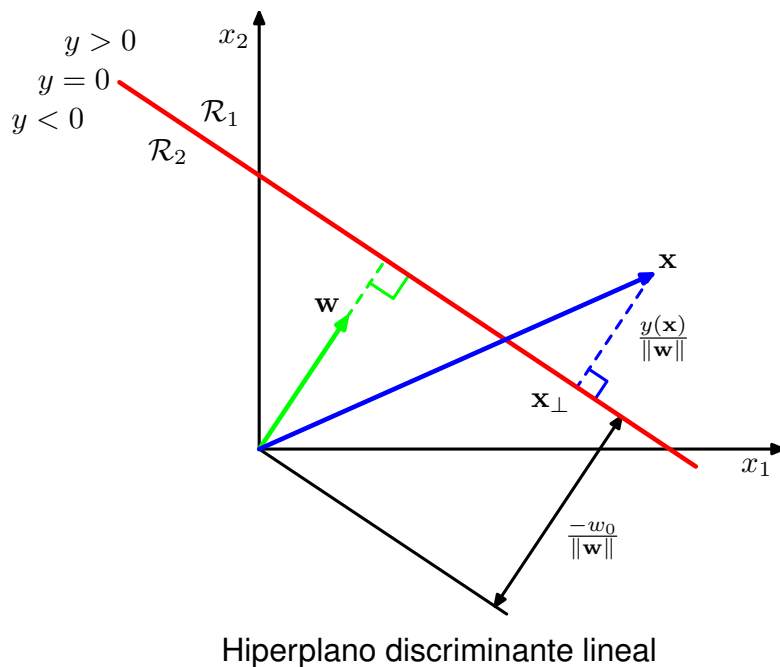


Ejemplo cluster

- Proponer un modelo
- Definir función de error que mida calidad cluster
- Estimar los parámetros:
Optimizar

Análisis discriminante: Modelos lineales (I)

- Sea $\{x_i, t_i\}_{i=1}^N$ una muestra finita de patrones de entrenamiento y $\{C_j\}_{j=1}^M$ el conjunto de clases. $t_{ij} = 1$ si x_i pertenece a la clase C_j .

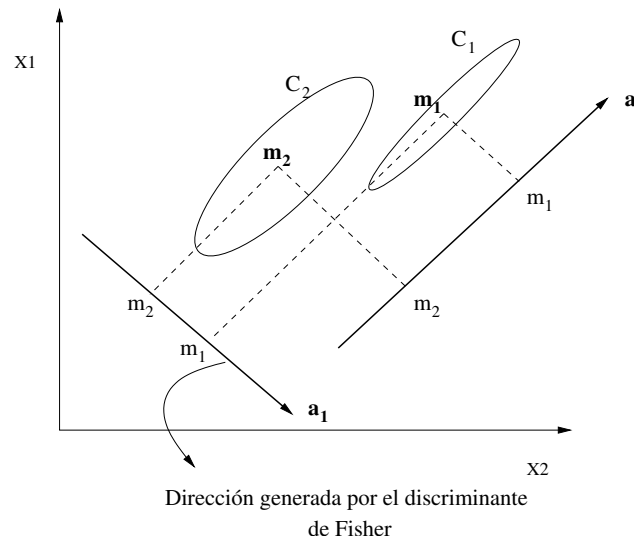


- La función discriminante lineal se escribe como:
$$f(x; w) = w^T x + w_0$$

 $f(x) \geq 0$ entonces $x \in C_1$,
en otro caso $x \in C_2$.
- ¿ Función de error a optimizar ?

Análisis discriminante: Modelos lineales (II)

Enfoque Fisher: Proyecta los datos sobre un hiperplano lineal que maximiza la separabilidad entre clases.



Criterio de separabilidad (biclase):

- ▷ Maximiza separación entre centros de clases proyectadas $(\bar{m}_2 - \bar{m}_1)^2$.
- ▷ Minimiza la dispersión de cada clase en torno a su media $(s_1^2 + s_2^2)$ donde $s_k^2 = \sum_{i \in C_k} (\bar{x}_i - \bar{m}_k)^2$.

Criterio de Fisher

La función de error a optimizar se puede escribir como:

$$J(\mathbf{w}) = \frac{(\bar{m}_2 - \bar{m}_1)^2}{s_1^2 + s_2^2} = \quad (6)$$

$$= \frac{\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w}}{\sum_{x_i \in C_1} \mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_1) (\mathbf{x}_i - \mathbf{m}_1)^T \mathbf{w} + \sum_{x_i \in C_2} \mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_2) (\mathbf{x}_i - \mathbf{m}_2)^T \mathbf{w}} \quad (7)$$

Definimos las matrices de covarianza interclase e intraclase como:

$$S_B = (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^T \quad (8)$$

$$S_w = \sum_{x_i \in C_1} \mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_1) (\mathbf{x}_i - \mathbf{m}_1)^T \mathbf{w} + \sum_{x_i \in C_2} \mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_2) (\mathbf{x}_i - \mathbf{m}_2)^T \mathbf{w} \quad (9)$$

Análisis discriminante: Modelos lineales (III)

Esto supone optimizar el siguiente criterio:

$$J_F = \frac{|\mathbf{w}^T S_b \mathbf{w}|}{|\mathbf{w}^T S_w \mathbf{w}|}, \quad (10)$$

Esto equivale a resolver el siguiente problema de autovalores y autovectores:

$$S_b \mathbf{w} = S_w \mathbf{w} \Lambda \quad (11)$$

que se puede resolver eficientemente mediante operaciones de álgebra lineal.

Proyectados los datos, se asignan a la clase de centroide más cercano.

Análisis discriminante: Modelos lineales (IV)

Propiedades discriminante Fisher:

- Es atractivo para trabajar en alta dimensión y eficiente computacionalmente.
- La frontera inducida es lineal y no está preparado para trabajar con clases multimodales.
- La dimensión máxima del espacio proyección es $M - 1$.
- Es bastante sensible a atípicos y no maximiza la capacidad de generalización.

Análisis discriminante: Modelos lineales (V)

Discriminante lineal puede **optimizar** también el **error cuadrático medio**:

$$E = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i + w_0 - t_i)^2 \quad (12)$$

Definimos la variable respuesta:

$$t_i = \begin{cases} N/N_1 & \text{si } \mathbf{x}_i \in C_1 \\ -N/N_2 & \text{si } \mathbf{x}_i \in C_2 \end{cases} \quad (13)$$

Calculando las derivadas e igualando a 0 obtenemos:

$$\frac{\partial E}{\partial w_0} = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i + w_0 - t_i) = 0 \quad (14)$$

$$\frac{\partial E}{\partial \mathbf{w}} = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i + w_0 - t_i) \mathbf{x}_i = 0 \quad (15)$$

$$(16)$$

Operando se obtiene:

$$(\mathbf{S}_w + \frac{N_1 N_2}{N} \mathbf{S}_B) \mathbf{w} = N(\mathbf{m}_1 - \mathbf{m}_2) \quad (17)$$

Teniendo en cuenta que $S_B w$ es paralelo a $m_2 - m_1$ el vector w de la función discriminante se puede expresar:

$$w \propto S_w^{-1}(m_2 - m_1) \quad (18)$$

que es equivalente a la expresión obtenida maximizando la separabilidad.