

Técnicas Selectas de Machine Learning y Aplicaciones

Presentado: Juan Manuel Hurtado Restrepo

Supervisor: José Ángel González Prieto

Facultad de Ciencias Matemáticas
Universidad Complutense de Madrid

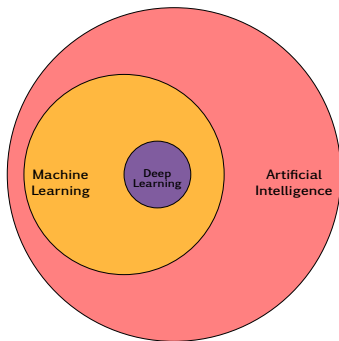
Defensa pública del TFG, 10 de julio de 2025

- 1 Introducción
- 2 Técnicas de regresión
- 3 Técnicas de clasificación
 - Clustering
- 4 Aplicaciones
 - Regresión
 - Clasificación
- 5 Novedades y desafíos
- 6 Conclusiones
- 7 Referencias

- 1 **Introducción**
- 2 Técnicas de regresión
- 3 Técnicas de clasificación
 - Clustering
- 4 **Aplicaciones**
 - Regresión
 - Clasificación
- 5 Novedades y desafíos
- 6 Conclusiones
- 7 Referencias

Machine learning e IA

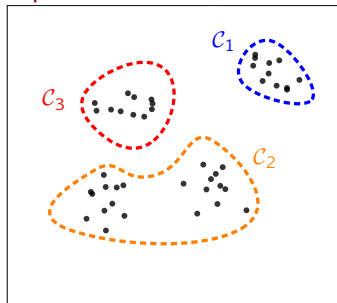
Aprendizaje supervisado y no supervisado



- El **aprendizaje supervisado** se basa en un conjunto de datos, de los que se conocen los correspondientes resultados esperados.

El objetivo es entonces entrenar un sistema de aprendizaje automático con la información anterior, para predecir los resultados cuando se presenten nuevas entradas.

- El enfoque de **aprendizaje no supervisado**:



Problema de regresión y problema de clasificación

Problema de regresión



Figura: Francis Galton.

Del modelo de Galton originado tras el análisis de 205 familias:

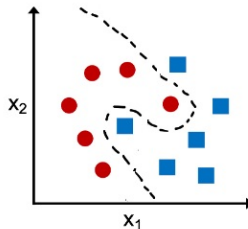
$$E[Y|X] = \frac{2}{3}X.$$

A la formulación general:

$$Y = f(X, \Omega) + \varepsilon$$

Problema de clasificación

Una colección de N objetos O_i con $i \in \{1, \dots, N\}$, cada uno representado por un vector real de dimensión d , $\mathbf{x}_i \in \mathbb{R}^d$; un conjunto $\{1, \dots, M\}$ de clases o categorías; una variable Y etiquetadora que toma valores en dicho conjunto; y una función clasificadora f que mapea vectores del espacio \mathbb{R}^d al conjunto de etiquetas de clase $\{1, \dots, M\}$.



Métodos de estimación de parámetros

Mínimos cuadrados



Figura: Carl F. Gauss.

Para $f(x_k, \boldsymbol{\Omega}) = \sum_{j=1}^m \omega_j \phi_j(x_k) \forall k = 1, \dots, n$ donde $\boldsymbol{\Omega}^T = (\omega_1, \dots, \omega_m)$, tenemos las *ecuaciones normales de Gauss* $\forall l = 1, \dots, m$:

$$\sum_{k=1}^n \left(\sum_{j=1}^m \omega_j \phi_j(x_k) \right) \phi_l(x_k) = \sum_{k=1}^n y_k \phi_l(x_k)$$

Máxima verosimilitud

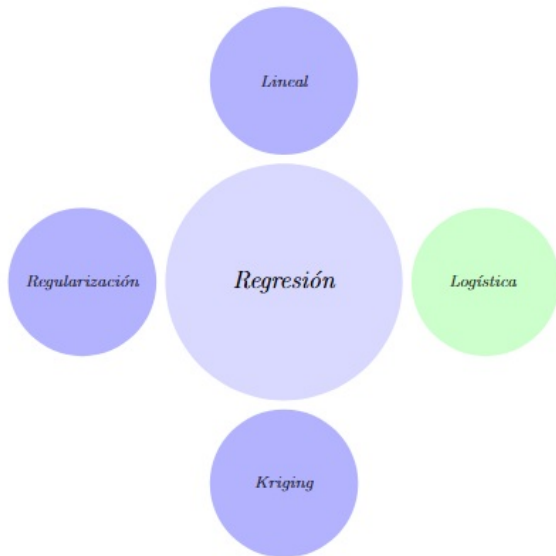


Figura: Ronald A. Fisher

De la realización de una muestra aleatoria X_1, \dots, X_n , proveniente de una población X con distribución de probabilidad $f_{\boldsymbol{\omega}}$ que depende de un parámetro desconocido $\boldsymbol{\omega} \in \Omega \subseteq \mathbb{R}^m$, la función protagonista es la *verosimilitud* $L(\boldsymbol{\omega}) = f_{\boldsymbol{\omega}}(x_1, \dots, x_n)$.

- 1 Introducción
- 2 **Técnicas de regresión**
- 3 **Técnicas de clasificación**
 - Clustering
- 4 **Aplicaciones**
 - Regresión
 - Clasificación
- 5 **Novedades y desafíos**
- 6 **Conclusiones**
- 7 **Referencias**

Técnicas de regresión



Regresión lineal

Modelo de regresión lineal

Siendo ε la discrepancia, la variable respuesta es:

$$Y = \Phi^T(\mathbf{X}) \cdot \Omega + \varepsilon$$

donde $\mathbf{X}^T = (X_1, \dots, X_p)$, $\Phi^T = (\phi_0, \dots, \phi_n)$ con $\phi_0(\mathbf{X}) = 1$ y $\Omega^T = (\omega_0, \dots, \omega_n)$.

Estimación de parámetros

Considerando una muestra de entrenamiento $D = \{(\mathbf{x}_k, y_k)\}_{k=1}^N \subseteq \mathbb{R}^p \times \mathbb{R}$ y siendo $\mathbf{y}^T = (y_1, \dots, y_N)$, tenemos lo siguiente:

$$\hat{\Omega} = \underbrace{(\Psi^T \Psi)^{-1}}_{\Psi^+} \Psi^T \mathbf{y}, \quad \text{donde} \quad \Psi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \dots & \phi_n(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \dots & \phi_n(\mathbf{x}_N) \end{pmatrix}$$

Regularización

Elastic net

Para $0 \leq \alpha \leq 1$ y $\Omega^T = (\omega_0, \dots, \omega_m)$, formalmente el problema es:

$$\hat{\Omega}^{e.n.} = \underset{\Omega}{\operatorname{argmin}} \left\{ \sum_{k=1}^{n-1} \left(y_k - \omega_0 - \sum_{l=1}^m x_{kl} \omega_l \right)^2 + \lambda \sum_{l=1}^m \left(\alpha |\omega_l|^2 + (1 - \alpha) |\omega_l| \right) \right\}$$

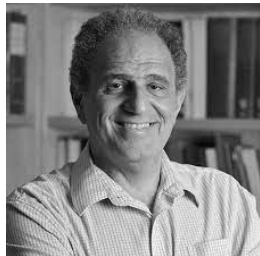


Figura: Robert Tibshirani.

Casos especiales

- $\alpha = 0 \longrightarrow \text{LASSO.}$
- $\alpha = 1 \longrightarrow \text{ridge.}$

Generalización

Para $p \in \mathbb{N} \cup \{\infty\}$, la técnica se puede generalizar considerando:

$$\|\cdot\|_p^p : \mathbb{R}^m \rightarrow \mathbb{R}, \mathbf{z} \mapsto \|\mathbf{z}\|_p^p = \sum_{l=1}^m |z_l|^p.$$

Funciones kernel

Para una aplicación $\phi : \mathcal{X} \rightarrow \mathcal{V}, \mathbf{x} \mapsto \phi(\mathbf{x})$, con \mathcal{X} el espacio de entrada y \mathcal{V} un espacio con producto interior, una función kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ se define como:

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{V}} = \phi(\mathbf{x})^T \phi(\mathbf{x}') = \sum_{i=1}^m \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

donde $\phi^T = (\phi_1, \dots, \phi_m)$.

Estimación en regresión de procesos gaussianos

Dado el conjunto de entrenamiento $D = \{(\mathbf{x}_k, y_k)\}_{k=1}^N \subseteq \mathbb{R}^p \times \mathbb{R}$, el kernel gaussiano $k(\mathbf{x}, \mathbf{x}') = \exp(-\theta \|\mathbf{x} - \mathbf{x}'\|)$ y una nueva observación \mathbf{x}_{N+1} , los parámetros de y_{N+1} son:

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T C_N^{-1} \mathbf{y}; \quad \sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T C_N^{-1} \mathbf{k}$$

Regresión logística

Estimación de parámetros en la versión binaria

Teniendo presente que $\boldsymbol{\Omega}^T = (\omega_0, \boldsymbol{\omega}^T)$ y $p(\mathbf{x}; \boldsymbol{\Omega}) = \frac{1}{\exp\{-(\omega_0 + \boldsymbol{\omega}^T \mathbf{x})\} + 1}$, si partimos de un conjunto de entrenamiento $D = \{(\mathbf{x}_k, y_k)\}_{k=1}^N \subseteq \mathbb{R}^P \times \mathbb{R}$, tenemos que:

$$\begin{aligned}\boldsymbol{\Omega}^{new} &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \\ \text{donde} \quad \mathbf{z} &= \mathbf{X} \boldsymbol{\Omega}^{old} + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p}),\end{aligned}$$

$$\mathbf{W} = \begin{pmatrix} p(\mathbf{x}_1; \boldsymbol{\Omega}^{old})(1 - p(\mathbf{x}_1; \boldsymbol{\Omega}^{old})) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & p(\mathbf{x}_N; \boldsymbol{\Omega}^{old})(1 - p(\mathbf{x}_N; \boldsymbol{\Omega}^{old})) \end{pmatrix},$$
$$\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}, \quad \mathbf{p} = \begin{pmatrix} p(\mathbf{x}_1; \boldsymbol{\Omega}^{old}) \\ \vdots \\ p(\mathbf{x}_N; \boldsymbol{\Omega}^{old}) \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} - & \mathbf{x}_1^T & - \\ & \vdots & \\ - & \mathbf{x}_N^T & - \end{pmatrix}$$

- 1 Introducción
- 2 Técnicas de regresión
- 3 Técnicas de clasificación**
 - Clustering
- 4 Aplicaciones
 - Regresión
 - Clasificación
- 5 Novedades y desafíos
- 6 Conclusiones
- 7 Referencias

Técnicas de clasificación



Formulación general

La hipótesis de independencia estocástica, nos permite escribir la probabilidad a posteriori como:

$$P(Y | \mathbf{X}) = \frac{P(Y)}{Z} \prod_{l=1}^n P(X_l | Y)$$

donde $Z = P(\mathbf{X}) = \sum_{k=1}^m P(Y = k) P(\mathbf{X} | Y = k), k \in \{1, \dots, m\}.$

El argumento de la técnica desemboca en maximizar la probabilidad a posteriori, sabiendo que $P(Y = k) = \frac{N_k}{N}$:

$$\hat{y} = \operatorname{argmax}_{y \in \{1, \dots, m\}} P(Y = y) \prod_{l=1}^n P(X_l | Y = y)$$

k-Nearest neighbors

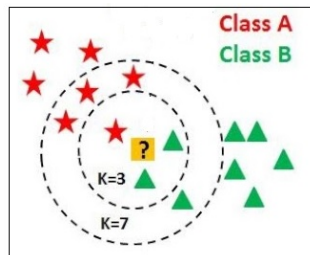
Planteamiento general

Partimos de un conjunto de datos de entrenamiento, $\mathbb{R}^n \times \mathcal{Y} \supseteq D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ con \mathcal{Y} es un conjunto finito. Dada una nueva observación $\mathbf{x} \in \mathbb{R}^n$, la predicción de clase será:

$$\hat{y} = \operatorname{argmax}_{c \in \mathcal{Y}} \sum_{(\mathbf{x}_i, y_i) \in D_k(\mathbf{x})} \chi_c(y_i), \text{ donde } c \equiv \{c\}, D_k(\mathbf{x}) \subseteq D$$

Observación

- $\chi_A(z) = \begin{cases} 1, & z \in A \\ 0, & z \notin A \end{cases}$
- El concepto cercanía está asociado a una distancia, $d_p(\mathbf{r}, \mathbf{q}) = (\sum_{i=1}^n |q_i - r_i|^p)^{\frac{1}{p}}$.



Decision trees

CART

Partimos de un conjunto de datos de entrenamiento, $D = \{(\mathbf{x}_k, y_k)\}_{k=1}^N \subseteq \mathbb{R}^n \times \mathcal{Y}$, donde \mathcal{Y} es un conjunto finito. Dividimos el conjunto de datos en dos subconjuntos:

$$\begin{cases} D_1 &= \{(\mathbf{x}_k, y_i) \in D \mid x_k^j \leq s\} \\ D_2 &= \{(\mathbf{x}_k, y_i) \in D \mid x_k^j > s\} \end{cases}, \quad x_k^j, s \in \mathbb{R}, j \in \{1, \dots, n\}.$$

Criterio

Reducción de la impureza, medida con el índice de Gini $G(t) = 1 - \sum_{k=1}^m p_k^2$ donde $p_k = \frac{1}{|t|} \sum_{k=1}^{|t|} \chi_c(y_k)$:

$$\Delta G = G(t) - \left(\frac{|D_1|}{|t|} G(D_1) + \frac{|D_2|}{|t|} G(D_2) \right)$$



Figura: Leo Breiman.

Support vector machines

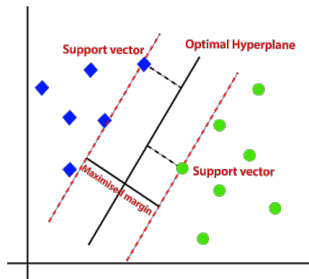
Caso binario linealmente separable

Sea el conjunto de entrenamiento $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, con $\mathbf{x}_k \in \mathbb{R}^n$ e $y_k \in \{-1, +1\}$. Maximizar el margen $\frac{2}{\|\boldsymbol{\Omega}\|}$, equivale a minimizar $\frac{1}{2}\|\boldsymbol{\Omega}\|^2$.

$$\begin{aligned} \text{mín :} & \quad \frac{1}{2}\|\boldsymbol{\Omega}\|^2 \\ \text{sujeto a:} & \quad y_k(\boldsymbol{\Omega}^T \mathbf{x}_k + b) \geq 1, \quad \forall k = 1, \dots, N. \end{aligned}$$

Generalización

$$\begin{aligned} \text{mín :} & \quad \frac{1}{2}\|\boldsymbol{\Omega}\|^2 + C \sum_{i=1}^n \xi_k \\ \text{s.a:} & \quad y_k(\boldsymbol{\Omega}^T \mathbf{x}_k + b) \geq 1 - \xi_k \\ & \quad \xi_k \geq 0, \quad \forall k = 1, \dots, N \end{aligned}$$



Linear discriminant analysis

Planteamiento general

Partiendo de un conjunto de clases $\{1, \dots, m\}$ y siendo

$f_k(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right\}$, se quiere determinar la probabilidad a posteriori:

$$P(Y = k \mid \mathbf{X} = \mathbf{x}) = \frac{f_k(\mathbf{x})\pi_k}{\sum_{l=1}^m f_l(\mathbf{x})\pi_l}, \quad \mathbf{x}^T = (x_1, \dots, x_n).$$

Estimación de parámetros

Utilizando los datos de entrenamiento $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^N \subseteq \mathbb{R}^n \times \mathbb{R}$, las estimaciones son:

- $\hat{\pi}_k = \frac{N_k}{N}$, $N_k :=$ número de observaciones de clase k .
- $\hat{\boldsymbol{\mu}}_k = \sum_{g_i=k} \frac{\mathbf{x}_i}{N_k}$, $g_i \in \{1, \dots, m\}$.
- $\hat{\Sigma} = \sum_{k=1}^m \sum_{g_i=k} \frac{(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^T}{N - m}$.

k-Means clustering

Planteamiento general

Dado el conjunto de datos $D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ con $\mathbf{x}_m \in \mathbb{R}^n$, el objetivo del algoritmo *k-means* es dividir estos datos en k grupos disjuntos C_1, \dots, C_k de tal forma que se minimice:

$$G_{k\text{-means}} = \sum_{m=1}^k \sum_{\mathbf{x} \in C_m} \|\mathbf{x} - \boldsymbol{\mu}_m\|^2, \text{ donde } \boldsymbol{\mu}_m = \frac{1}{|C_m|} \sum_{\mathbf{x} \in C_m} \mathbf{x}$$

El algoritmo puede inicializarse, eligiendo aleatoriamente k puntos como centroides, $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$. Cada punto \mathbf{x}_m se asigna al cluster cuyo centroide esté más cerca, lo que es equivalente a ir construyendo los clústeres a partir de los datos de D :

$$C_l = \left\{ \mathbf{x}_m \in D \mid l = \underset{j \in \{1, \dots, k\}}{\operatorname{argmin}} \|\mathbf{x}_m - \boldsymbol{\mu}_j\|, \forall m = 1, \dots, N \right\}, \quad \forall l \in \{1, \dots, k\}$$

- 1 Introducción
- 2 Técnicas de regresión
- 3 Técnicas de clasificación
 - Clustering
- 4 **Aplicaciones**
 - Regresión
 - Clasificación
- 5 Novedades y desafíos
- 6 Conclusiones
- 7 Referencias

Regresión

El dataset fue publicado originalmente por la Comisión de Taxis y Limusinas de Nueva York (TLC), disponibilizado en Kaggle (New York City Taxi Trip Duration) para una competencia cuyo objetivo era construir un modelo que permitiera predecir la duración total de los viajes en taxi en la ciudad de Nueva York.

Técnica	RMSLE (%)
linear regression	61,4
ridge regression	61,4
LASSO regression	62,9
elastic net regression	62,4

Cuadro: Comparativa de técnicas de regresión.

RMSLE

Para p_k la predicción de la duración del viaje y a_k la duración real del viaje, tenemos:

$$\text{RMSLE} = \sqrt{\frac{1}{N} \sum_{k=1}^N (\log \{p_k + 1\} - \log \{a_k + 1\})^2}$$

Clasificación

El origen de datos en este caso es el dataset MNIST, ampliamente utilizado por investigadores y estudiantes para reconocimiento de patrones y aprendizaje automático. La base de datos MNIST se derivó de datos originales del *National Institute of Standards and Technology* (NIST) y fue disponibilizada en Kaggle (MNIST Dataset) gracias a Yann A. LeCun Chief AI Scientist en Meta.

Técnica	k	Precisión (%)
k-nearest neighbours	1	95
	3	94,2
	5	93,5
	7	93
	9	92,4
SVM	–	95,1
k-means clustering	10	57,8

Cuadro: Comparativa de técnicas de clasificación y clustering.

Precisión

$$\text{Precisión} = \frac{\sum_m TC_m}{\sum_m (TC_m + FC_m)}$$

- 1 Introducción
- 2 Técnicas de regresión
- 3 Técnicas de clasificación
 - Clustering
- 4 Aplicaciones
 - Regresión
 - Clasificación
- 5 Novedades y desafíos**
- 6 Conclusiones
- 7 Referencias

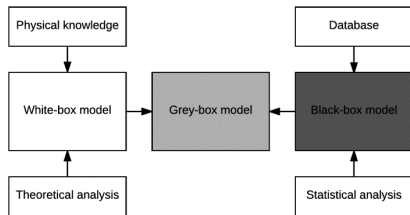


Figura: Enrico Camporeale.

El valor de las predicciones de los algoritmos de machine learning se hace evidente en el **space weather**, en situaciones de **regresión** como: valores de índices geomagnéticos, número de manchas solares, ocurrencia de erupciones solares, tiempo de propagación de un eyección de masa coronal, velocidad del viento solar y el flujo de partículas solares. Encontramos también circunstancias que dan pie a la **clasificación**, como: los distintos grupos de manchas solares y los tipos de viento solar de acuerdo con su origen.

- 1 Introducción
- 2 Técnicas de regresión
- 3 Técnicas de clasificación
 - Clustering
- 4 Aplicaciones
 - Regresión
 - Clasificación
- 5 Novedades y desafíos
- 6 Conclusiones
- 7 Referencias

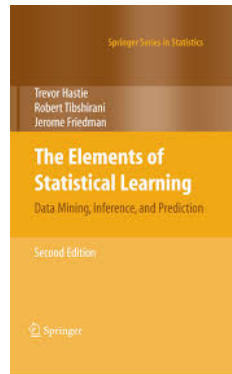
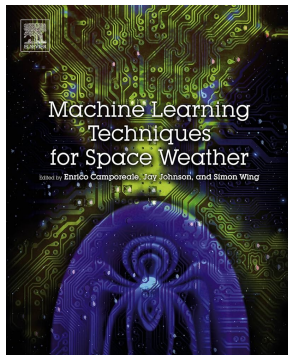
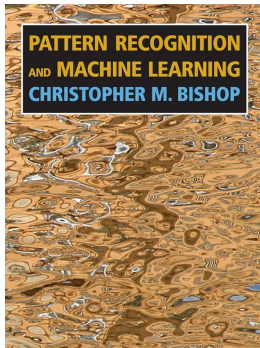
Conclusión

Identificamos dos puntos centrales a destacar que recogen la esencia de este proyecto y de los conocimientos adquiridos en su ejecución, estos son:

- 1 La profundización en algunas de las técnicas más distinguidas del aprendizaje supervisado y en una de las más conocidas de aprendizaje no supervisado.
- 2 La inmersión en el universo de las aplicaciones del aprendizaje automático, me permitió conocer el auténtico potencial transformador de esta herramienta, trascendiendo a las visiones parcializadas e incompletas de los medios de comunicación. A partir de esta probada capacidad revolucionaria, intuyo el enorme impacto que puede tener en disciplinas como la física, a tal punto que su irrupción puede suponer un punto de inflexión en la noble misión que históricamente se ha atribuido esta ciencia de descifrar los misterios del universo.

- 1 Introducción
- 2 Técnicas de regresión
- 3 Técnicas de clasificación
 - Clustering
- 4 Aplicaciones
 - Regresión
 - Clasificación
- 5 Novedades y desafíos
- 6 Conclusiones
- 7 Referencias**

Referencias



¡Gracias!