

Modelos Estadísticos para Clasificación y Regresión

Práctico 5 - Clasificación

IMERL - FIng

13 de setiembre de 2023



1 Repaso

2 Práctico 5

1 Repaso

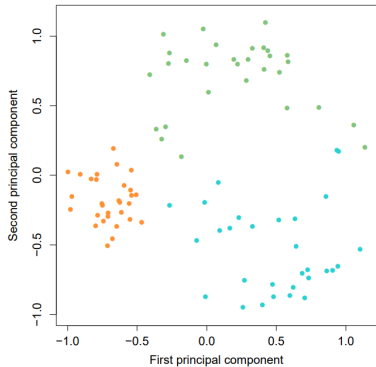
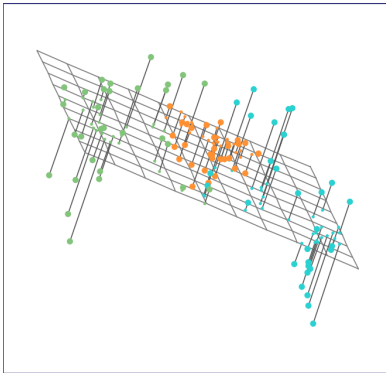
2 Práctico 5

Análisis de Componentes Principales (*PCA*)

Cuando tenemos un dataset grande con atributos correlacionados, PCA nos permite 'resumirlo' utilizando una menor cantidad de atributos que colectivamente explican la mayor parte de la varianza del dataset.

Cada observación vive en un espacio p -dimensional, pero no todas las dimensiones son igual de 'interesantes'. PCA busca un número de dimensiones que sean lo más interesantes posible, en el sentido de la varianza de los datos.[1]

Análisis de Componentes Principales (PCA)



1

¹Imagen extraída de [1].

Problema de Clasificación

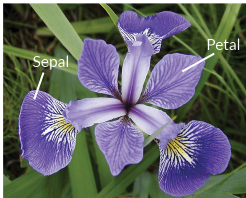
Tenemos un conjunto de N datos de entrenamiento x_1, x_2, \dots, x_N de dimensión D . Para cada dato tendremos su etiqueta correspondiente y_1, y_2, \dots, y_N que describe a cuál de las C clases pertenece. Típicamente estas etiquetas corresponden a valores enteros (por ej: $\{1, 2, 3\}$, $\{0, 1\}$, $\{-1, 1\}$).

Objetivo: dado un nuevo dato x_{nuevo} poder predecir su clase y_{nuevo} .

[4]



Ejemplos de clasificación



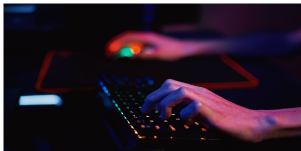
Iris Versicolor



Iris Setosa

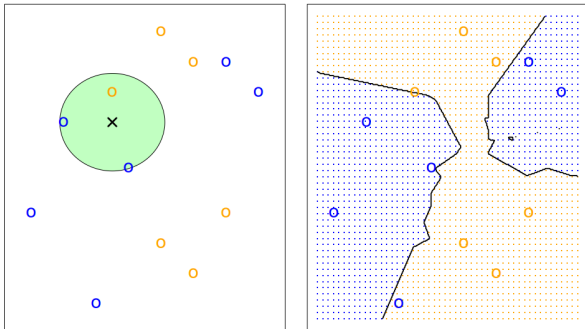


Iris Virginica



Vecinos más Cercanos (kNN)

Para clasificar x_{nuevo} con kNN , hallamos los K puntos de entrenamiento que están más cerca de x_{nuevo} . Luego tomamos t_{nuevo} como la etiqueta de la mayoría de esos K puntos. [4].

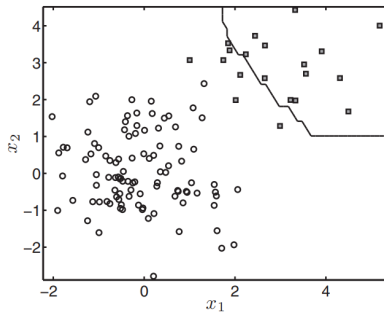
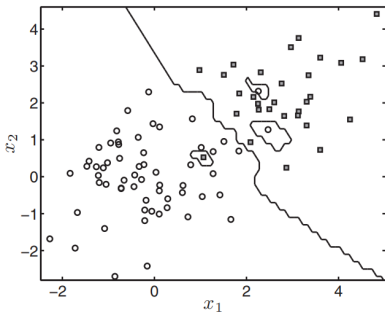


2

²Imagen extraída de [1].

k NN: elección de K

- ¿Qué pasa si tomamos K muy chico?
- ¿Y muy grande?



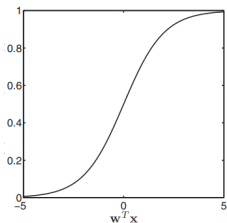
3

³Imagen extraída de [4].

Regresión Logística

El modelo lineal no permite modelar la probabilidad de que un dato X pertenezca a una cierta categoría, ya que podríamos obtener $p(X) > 1$ o $p(X) < 0$. Para solucionar esto, a la salida del modelo lineal $s = w^T X$ se aplica la función *sigmoide*:

$$p(X) = \frac{e^s}{1 + e^s}$$

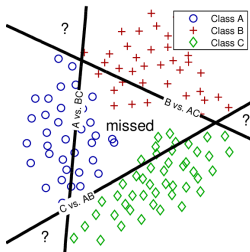


4

⁴Imagen extraída de [4].

Regresión Logística - esquema OvR

Si bien la Regresión Logística está pensada para problemas de clasificación binaria, su uso puede extenderse a clasificación multiclase usando esquemas OvO o OvR .
SkLearn implementa OvR :

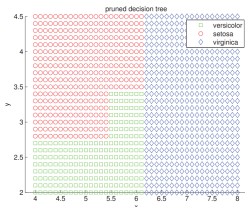
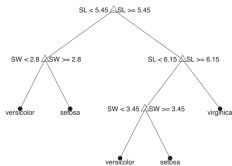


5

⁵Imagen extraída de [3].

Árboles de Decisión

Los árboles de decisión permiten particionar el espacio en regiones simples. El nombre proviene de que las reglas que se usan para segmentar el espacio se pueden resumir en un árbol.



Una vez entrenado el árbol, se predice la clase de un nuevo dato x_{nuevo} a partir de la clase mayoritaria de los datos de entrenamiento que caen en la misma hoja que x_{nuevo} .

⁶Imagen extraída de [2].

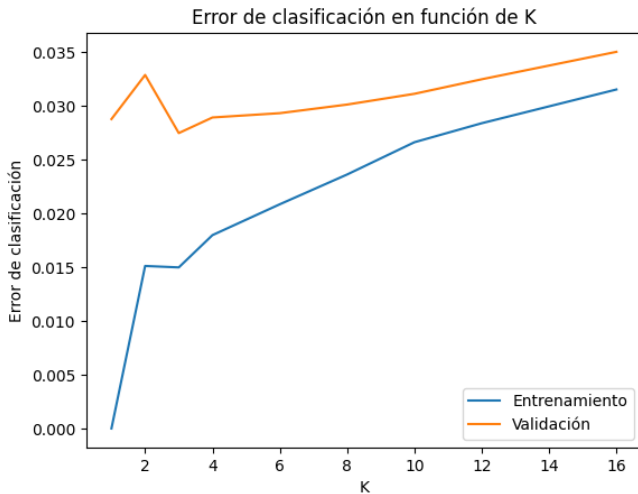
Preguntas



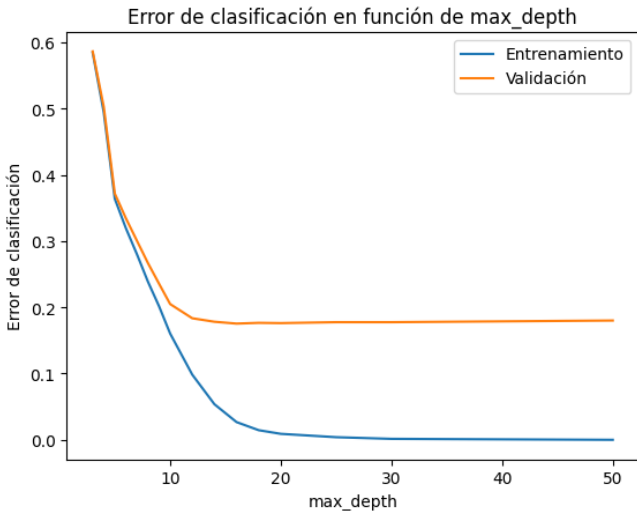
1 Repaso

2 Práctico 5

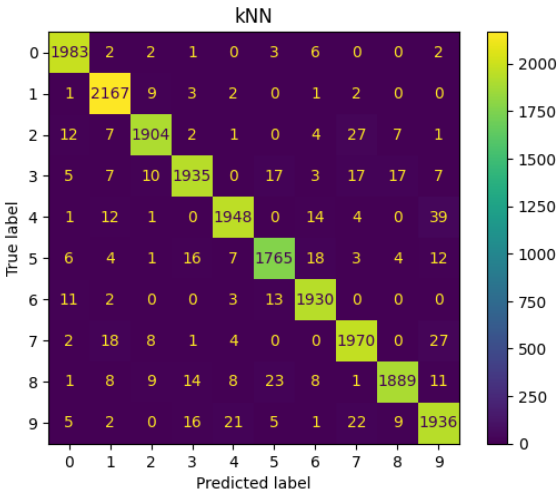
Práctico 5



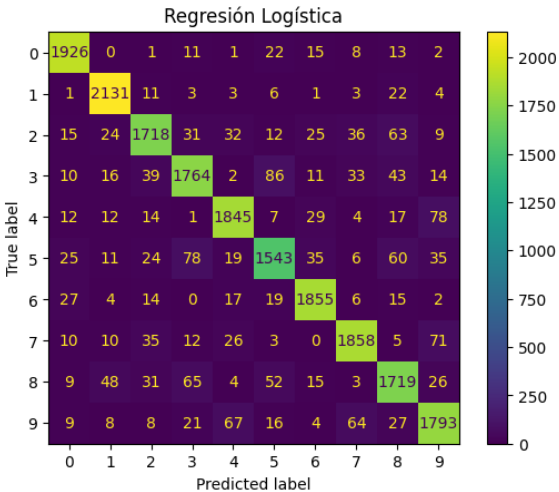
Práctico 5



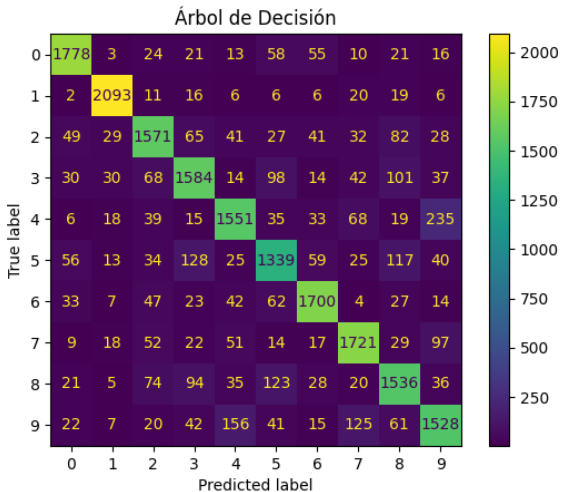
Práctico 5



Práctico 5



Práctico 5



Referencias

- [1] Gareth James et al. *An Introduction to Statistical Learning*. Springer International Publishing, 2023. DOI: [10.1007/978-3-031-38747-0](https://doi.org/10.1007/978-3-031-38747-0). URL: <https://doi.org/10.1007/978-3-031-38747-0>.
- [2] Kevin P. Murphy. *Machine learning : a probabilistic perspective*. Cambridge, Mass. [u.a.]: MIT Press, 2013. ISBN: 9780262018029 0262018020.
- [3] Sebastian Nusser. “Robust Learning in Safety-Related Domains : machine learning methods for solving safety-related application problems”. Tesis doct. Jul. de 2009.
- [4] Simon Rogers y Mark Girolami. *A First Course in Machine Learning*. Chapman y Hall/CRC, oct. de 2016. DOI: [10.1201/9781315382159](https://doi.org/10.1201/9781315382159). URL: <https://doi.org/10.1201/9781315382159>.