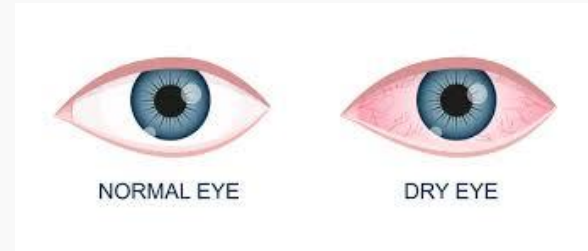


Proyecto :

Detección de síndrome del ojo seco empleando atributos y variables medicas

- Greyler Jose Chacon Chaparro (E2)
- Jonk Keyler Sanchez Pabon (E2)
- Samuel David Traslaviña Mateus (F1)



Contenido de la presentación

01

Presentación
del problema

02

Preprocesamiento
de los datos

03

Modelos y su
evaluación

04

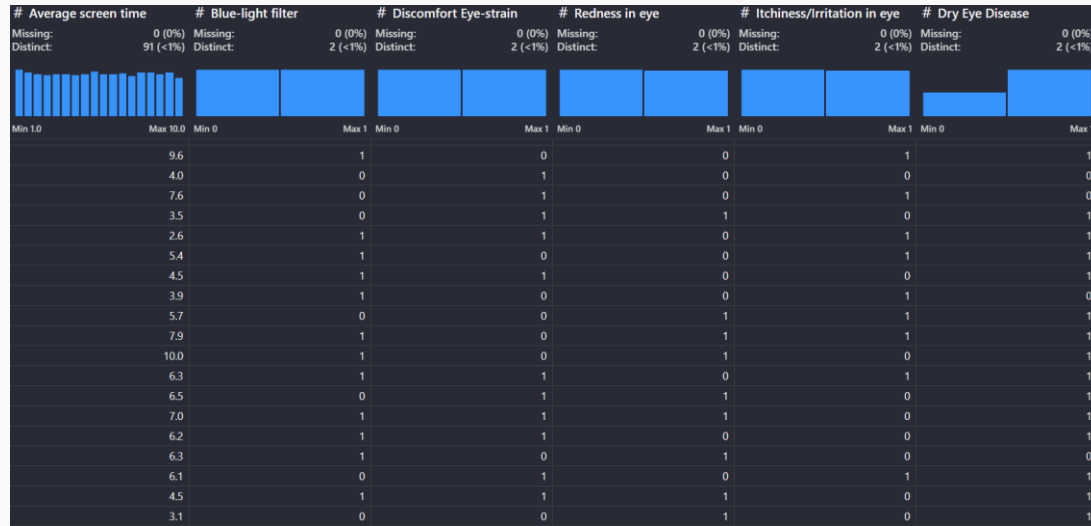
Curvas de
aprendizaje

01

Presentación del problema

Dataset

Se cuenta con el dataset *Dry Eye Disease* , que se encuentra disponible en kaggle. El *dataset* contiene información sobre los hábitos de sueño, actividad física y factores de salud de una muestra de persona, con el objetivo de analizar su relación con la enfermedad del ojo seco (*Dry Eye Disease*).





Problema

La enfermedad puede tener diferentes factores de riesgos que pueden estar relacionados a los hábitos que tiene la persona , como tiempo de exposición a pantallas , actividad física , edad , antecedentes médicos ,entre otros. Se cuenta con el *dataset* con un total de 26 *features* , entonces , se decidió entrenar modelos de *Machine learning* para **clasificar** si un paciente , dadas todas las variables , tiene o no la enfermedad (**clasificación binaria**) bajo un **enfoque supervisado**.

02

Preprocesamiento de los datos

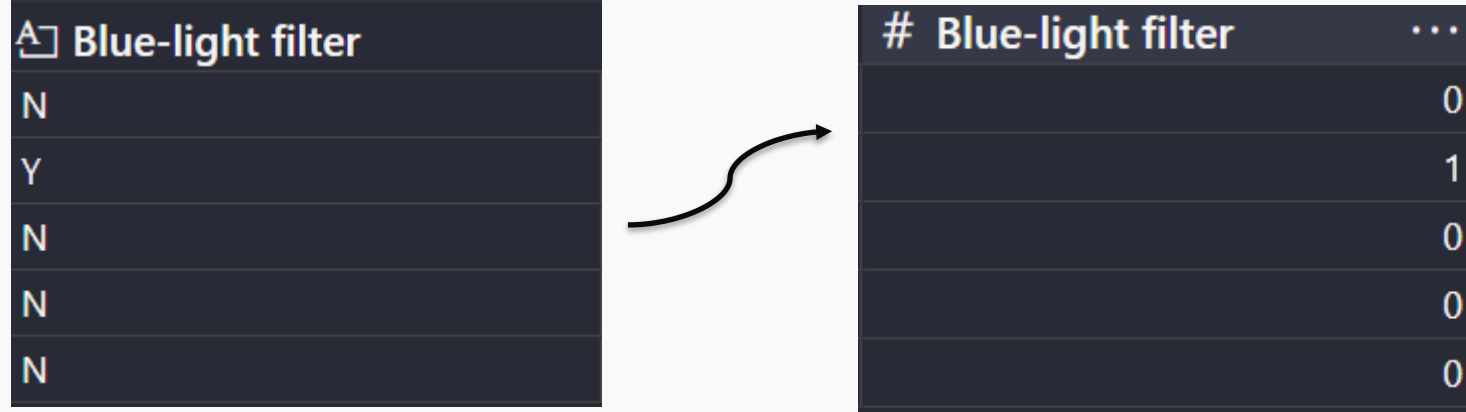
¿Por qué el preprocesamiento?

En la mayoría de las ocasiones que se cuenta con un proyecto de *machine learning* , los datos con los cuales se desea que el modelo aprenda y generalice para la tarea que se desea atacar , cuando están en bruto , pueden contener valores faltantes , algunas columnas se encuentran en un tipo de dato que los modelos no pueden utilizar , como texto y también existen modelos que son sensibles al escalado de los datos.

Para este proyecto se implementó la estrategia de ***label encoder*** y ***feature scaling***, ya que no se contaban con valores faltantes en los datos.

Label Encoder

Esta transformación identifica todos los distintos valores que se tienen para esa *feature*, que se denominan clases, y luego se les asignan un número entero entre el rango de $[0 - \text{\#clases}-1]$.

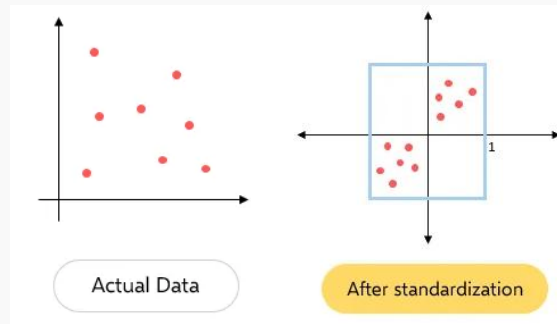


Estandarización

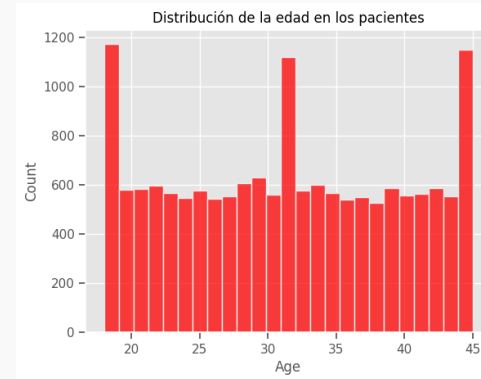
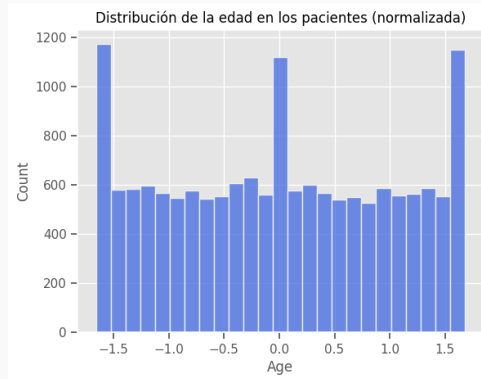
Con muy pocas excepciones , los algoritmos de *machine learning* , no funcionan bien cuando los datos numéricos se encuentran en diferentes escalas , La estandarizacion se hace necesaria cuando los conjuntos de datos contienen diferentes unidades de medición, rangos y magnitud , Si se usa, el modelo conducirá a un funcionamiento sesgado , La estandarización se define así :

$$z_i = \frac{x_i - \bar{x}_i}{s_i}, i = 1, \dots, n$$

Donde \bar{x}_i , es la media de la i-ésima *feature* y s_i su respectiva desviación



Realmente para los modelos que se usaron , solo la *SVM* , se le hace necesario la estandarización de las *features* antes de ajustar el modelo , ya que , sin una escala homogénea, características de gran varianza sesgan el espacio de características y el hiperplano óptimo.



03

Modelos y su evaluación

Métricas

La **matriz de confusión** , mas que ser una métrica en si , es una herramienta que permite contar el número de veces que las instancias de la clase A se clasifican como clase B, para todos los pares A/B

		Predicted	
		0	1
Actual	0	TN	FP
	1	FN	TP

El **accuracy** es una métrica fundamental para evaluar el rendimiento de un modelo de clasificación, proporcionando una vista rápida de qué tan bien se está desempeñando el modelo en términos de predicciones correctas. Se calcula como la relación de predicciones correctas al número total de muestras de entrada :

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Se debe tener la precaución cuando se trabaja con clases desbalanceadas.

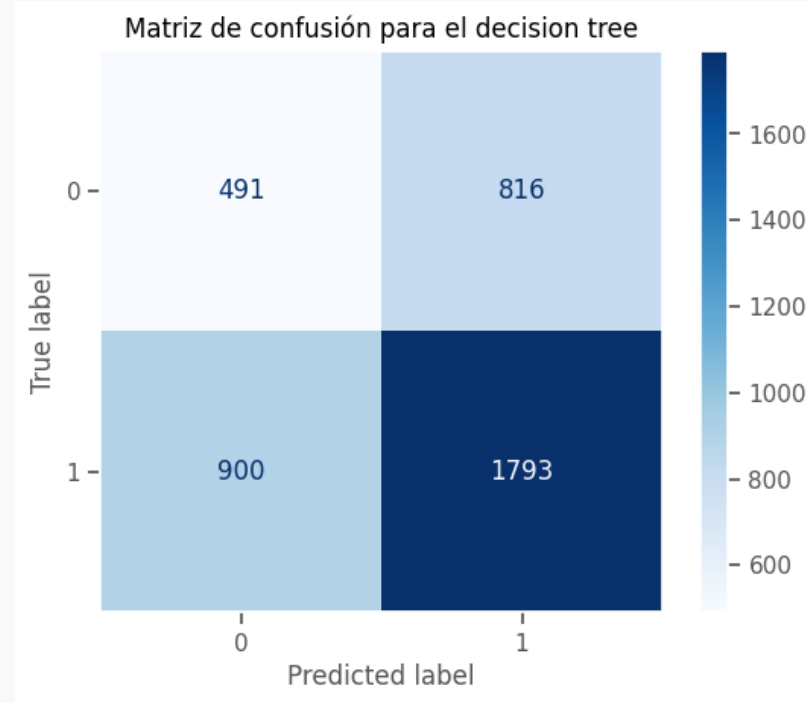
Otra métrica que fue empleada fue el **recall** , en el problema de detección del síndrome del ojo seco, porque se trató de identificar correctamente a los pacientes que sí tienen la condición, incluso si eso significa que el modelo ocasionalmente prediga algunos falsos positivos.

$$recall = \frac{TP}{TP + FN}$$

Gracias a la implementación de *scikit-learn* de `sklearn.metrics.classification_report`, nos permite también conocer la especificidad o tasa de verdaderos negativos :

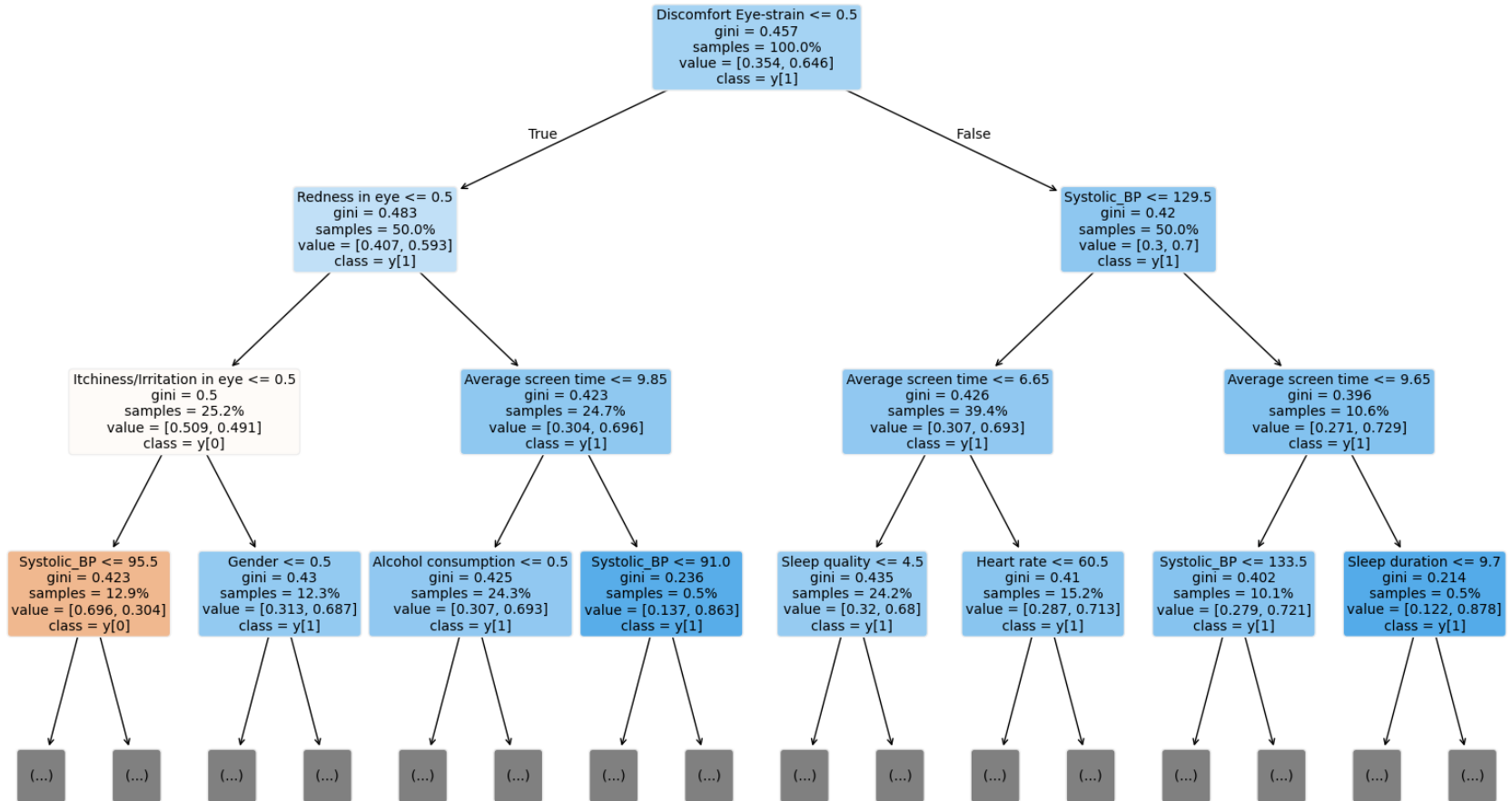
$$specificity = \frac{TN}{TN + FP}$$

Desicion Tree

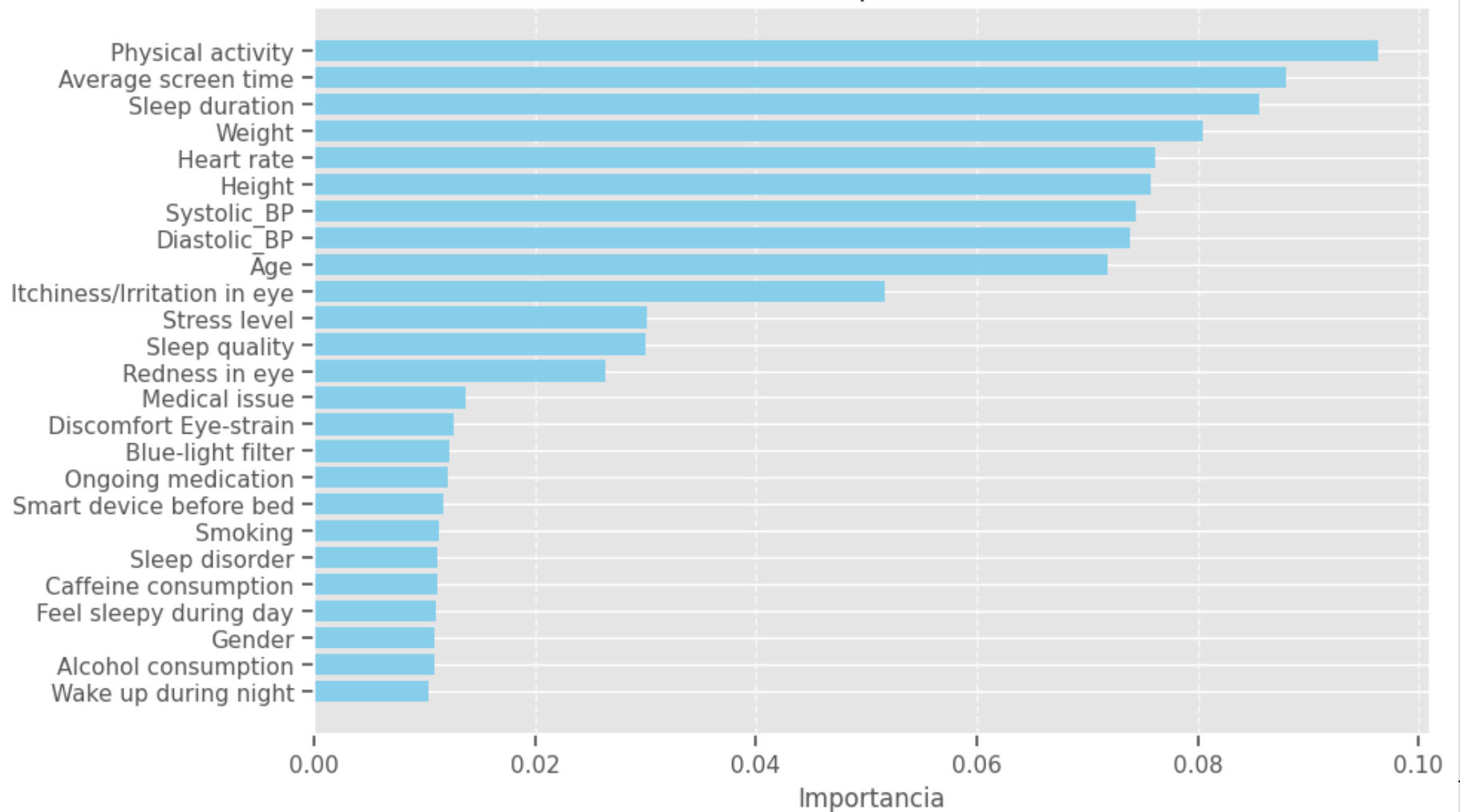


Accuracy: 0.571 , Recall : 0.38 , Specificity : 0.38

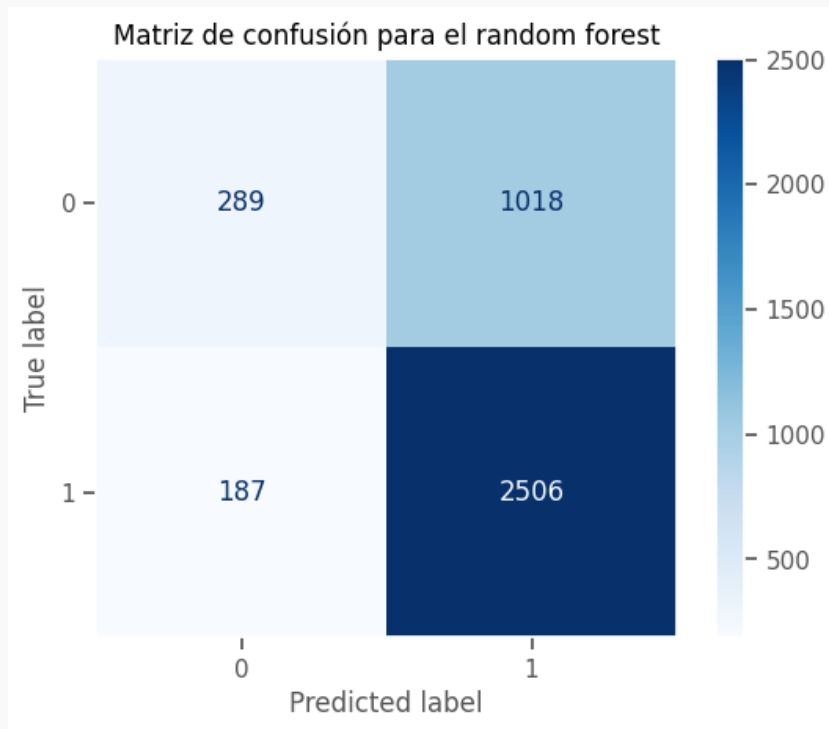
Árbol de Decisión



Características más Importantes (Decision Tree)

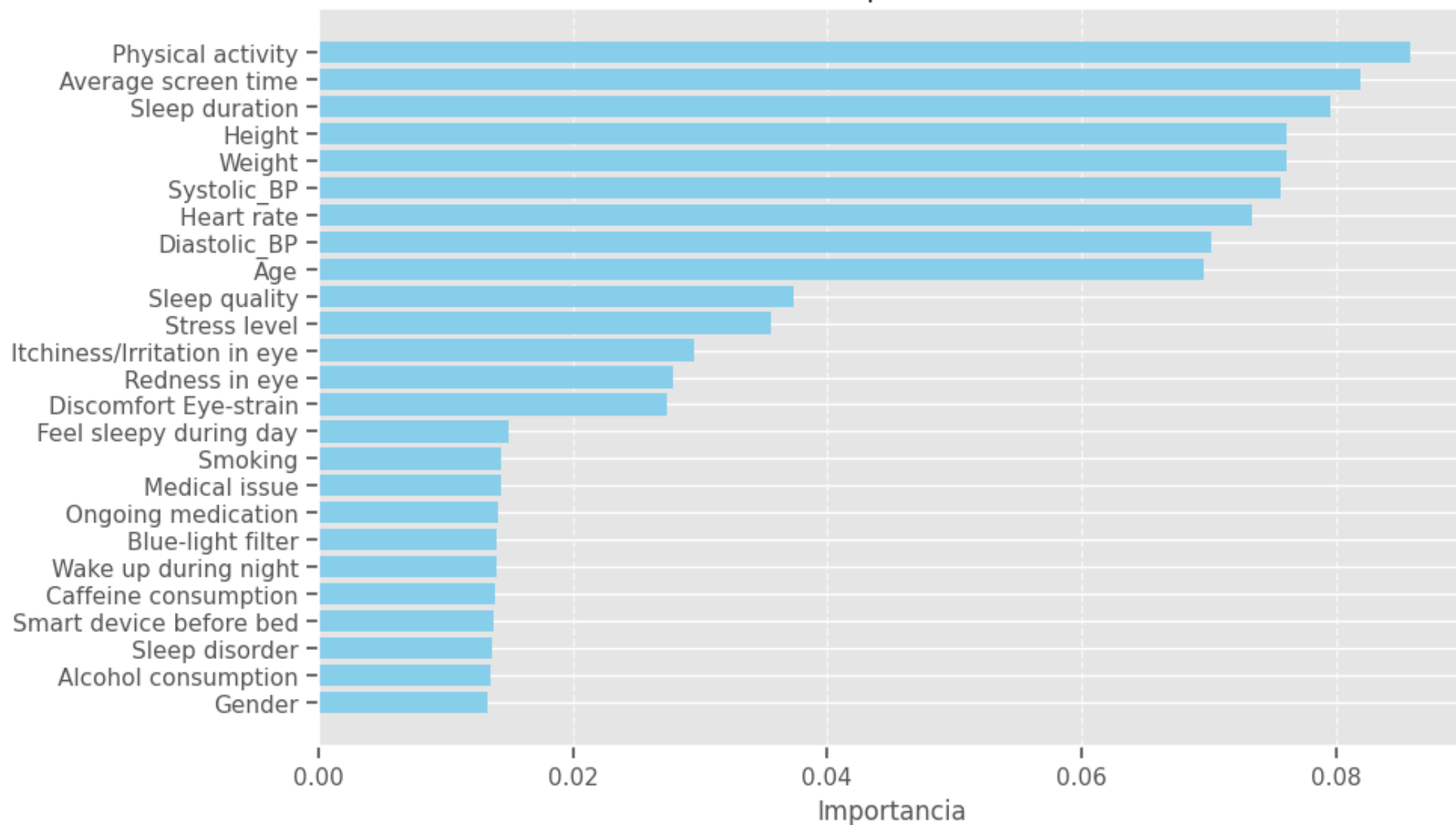


Random Forest

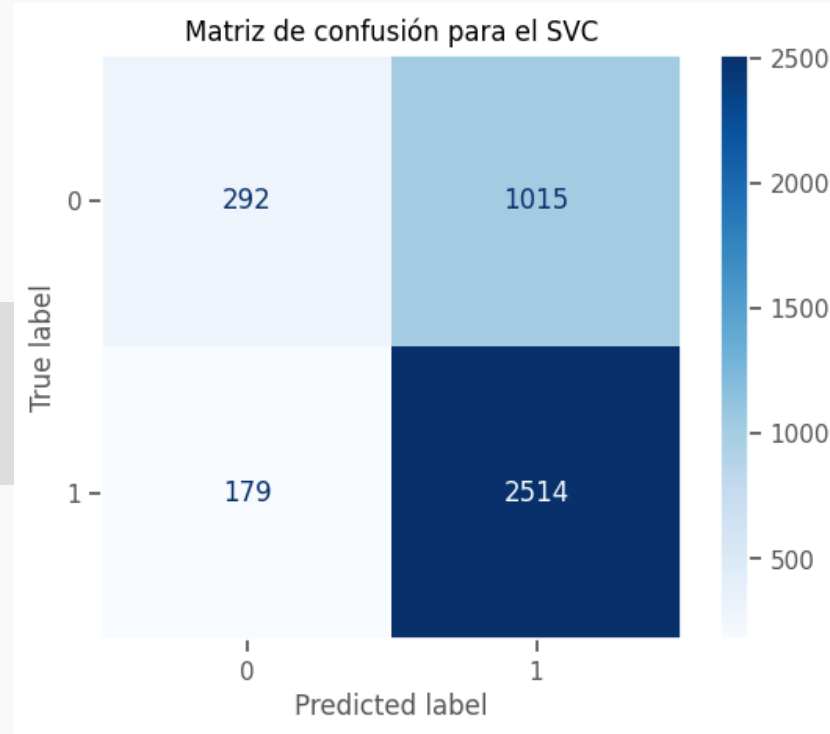


Accuracy: 0.699 , Recall : 0.93, Specificity : 0.22

Características más Importantes (Random Forest)



Support Vector Machine

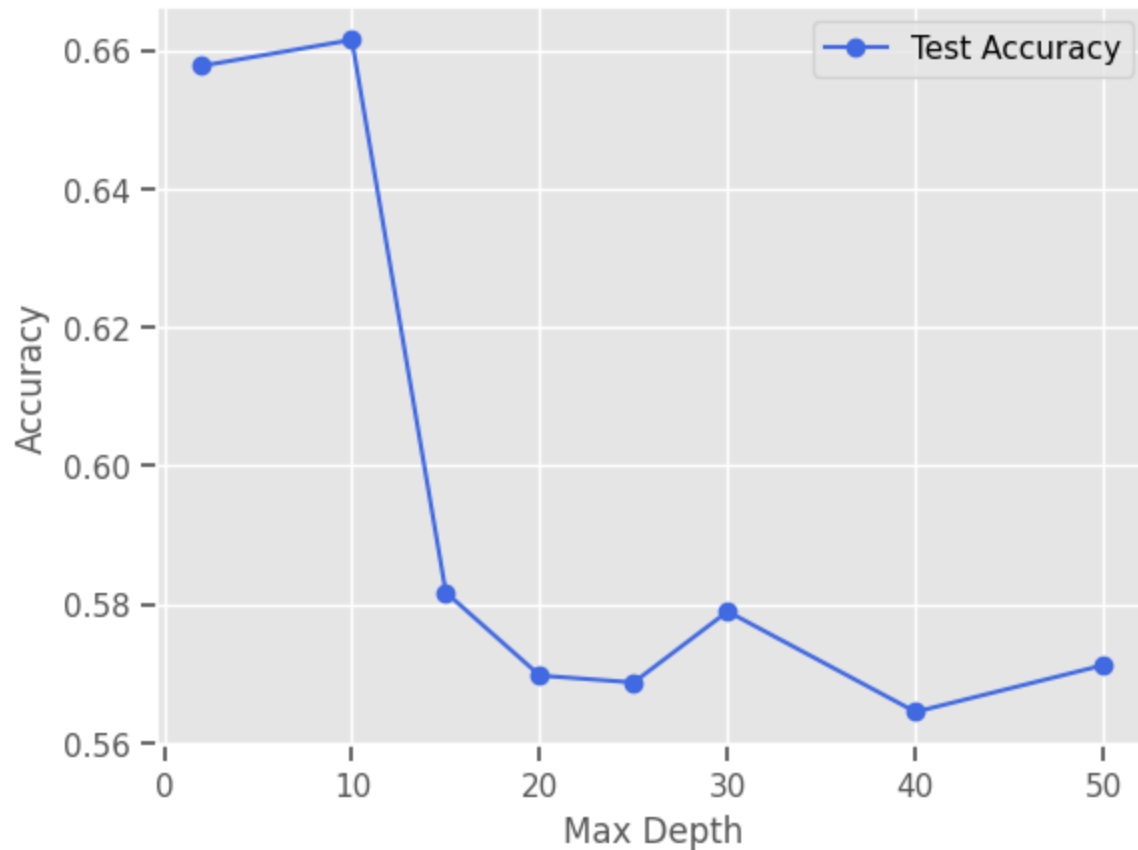


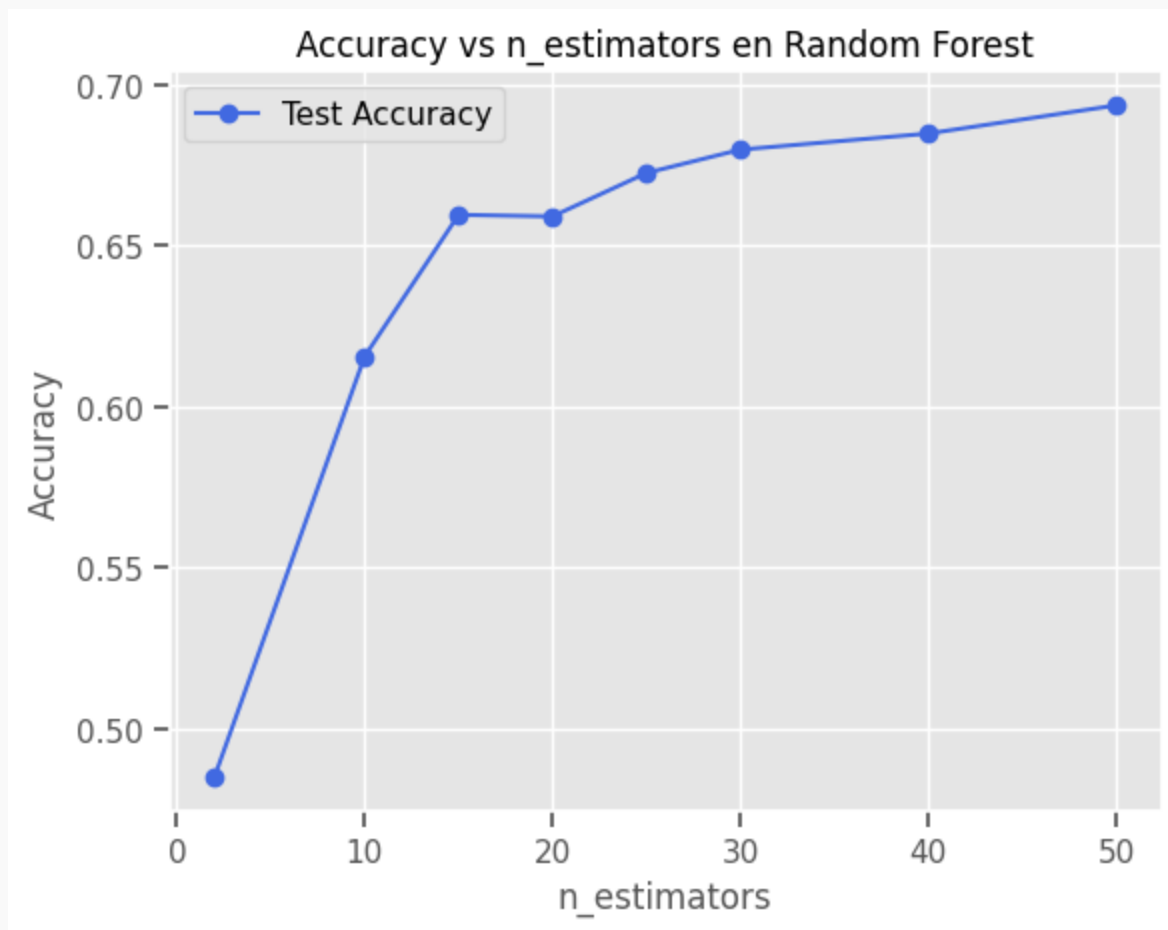
Accuracy: 0.70, Recall : 0.93, Specificity : 0.22

04

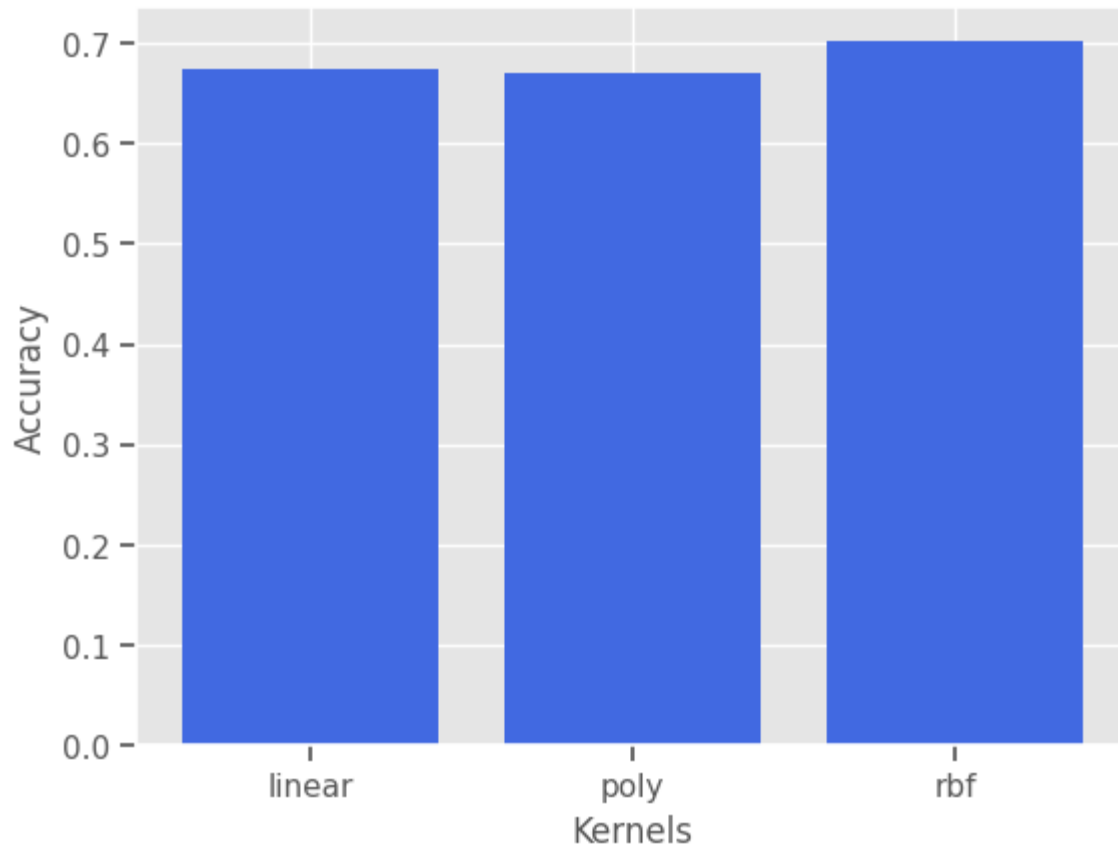
Curvas de aprendizaje

Accuracy vs Max Depth en Decision Tree

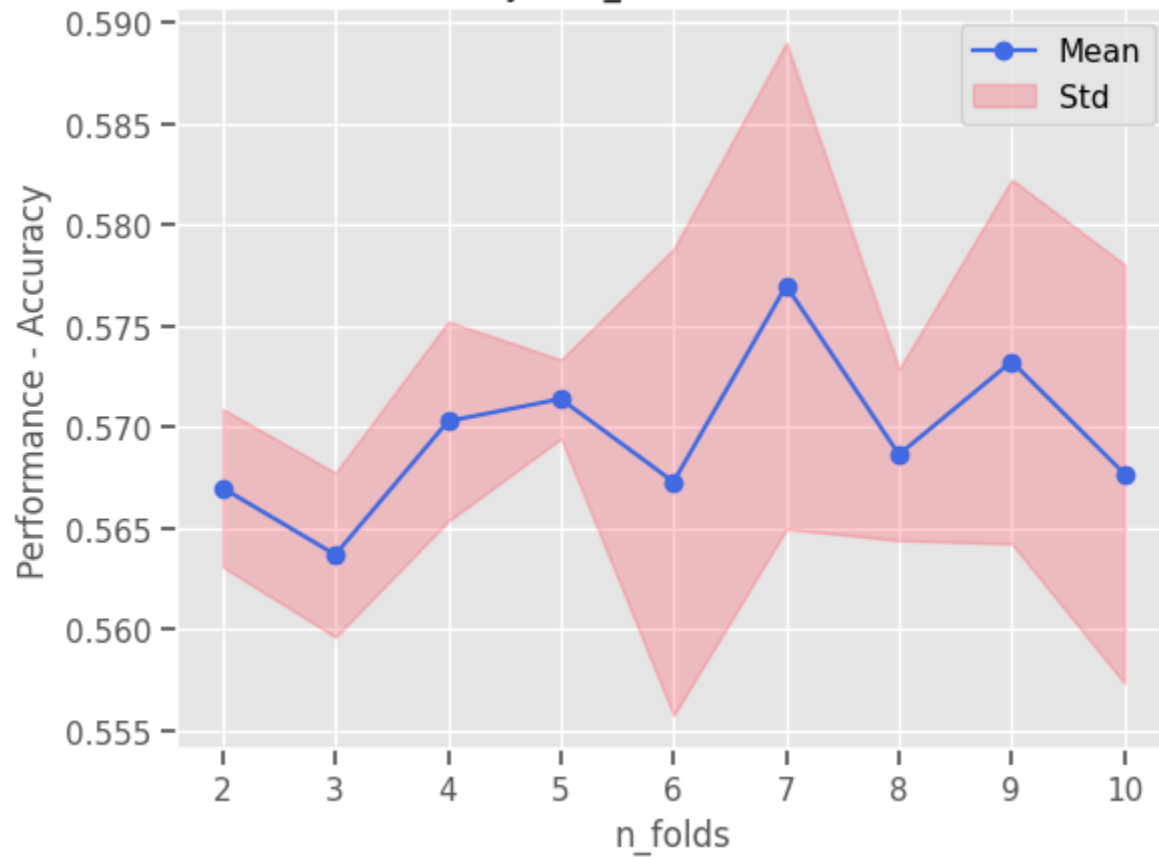




Accuracy vs Kernels in SVC



Accuracy vs n_folds in Decision Tree

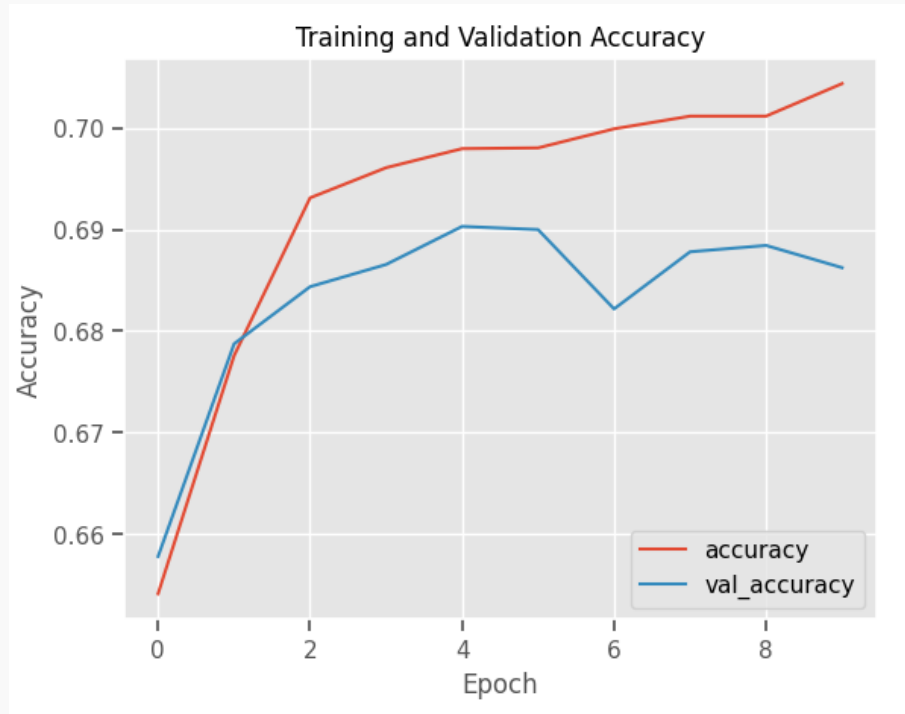


05

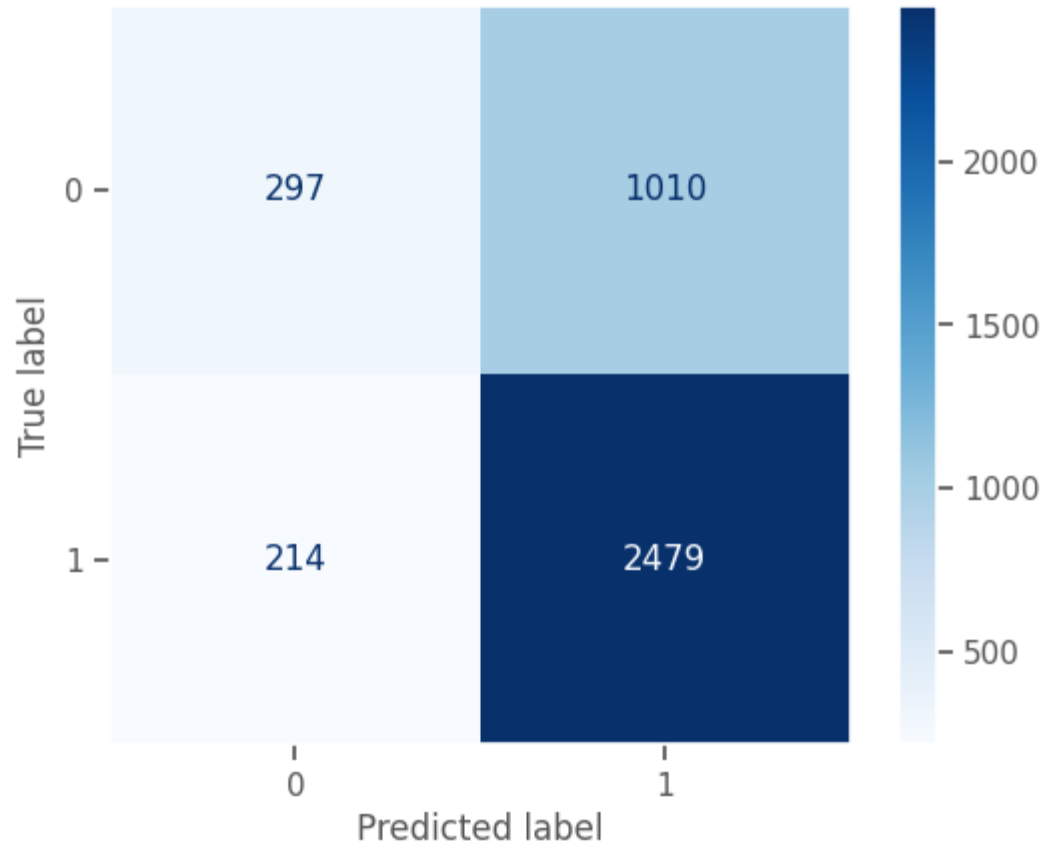
Deep Learning

La primera configuración fue :

1. Entrada (n_registros,25)
2. Tres capas ocultas (128 neuronas)
3. Una capa de salida – Función de activación sigmoide
4. Optimizador : Adam – lr : $1e-4$
5. Loss : binary cross entropy
6. 10 Epocas , batch-size : 16

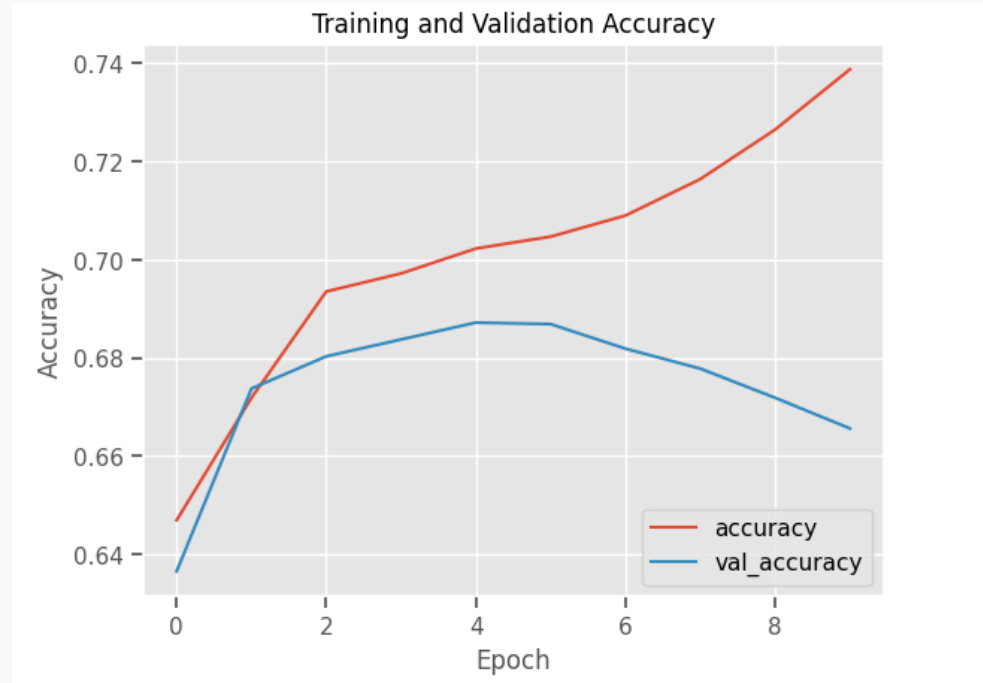


Matriz de confusión para la red neuronal con tres capas ocultas

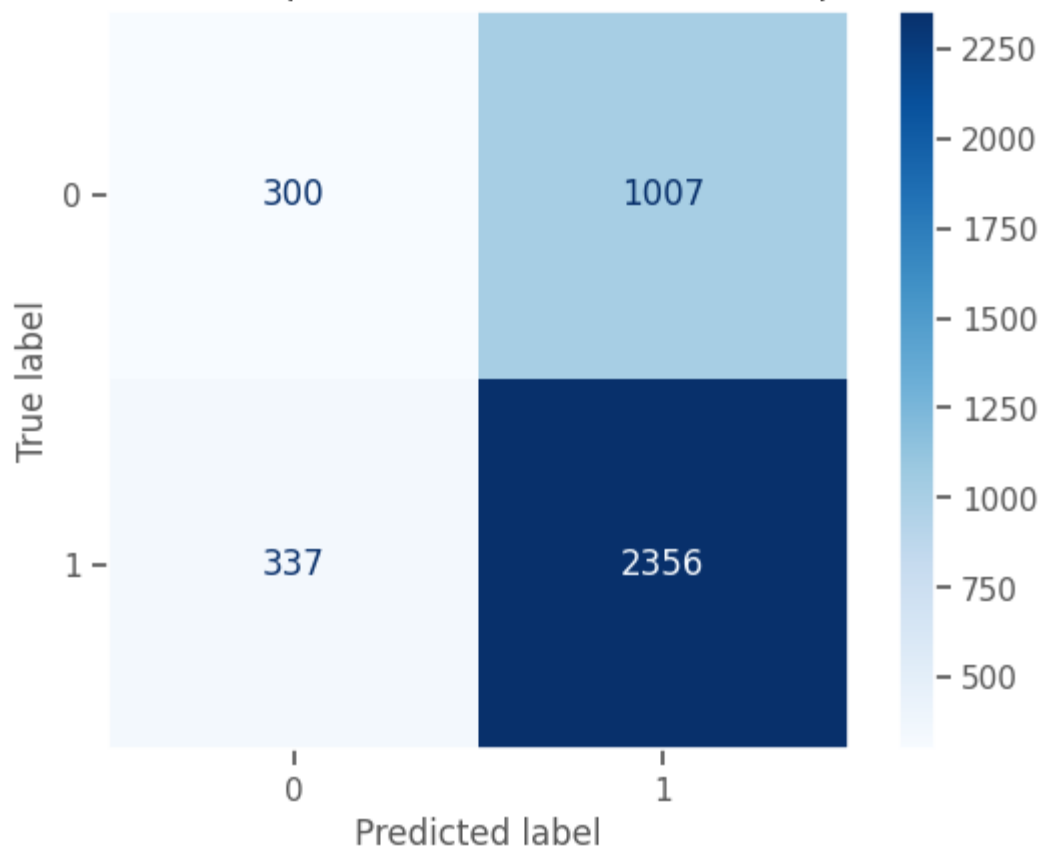


La segunda configuración fue :

1. Entrada (n_registros,25)
2. Seis capas ocultas (128 neuronas)
3. Una capa de salida – Función de activación sigmoide
4. Optimizador : Adam – lr : $1e-4$
5. Loss : binary cross entropy
6. 10 Epocas , batch-size : 32

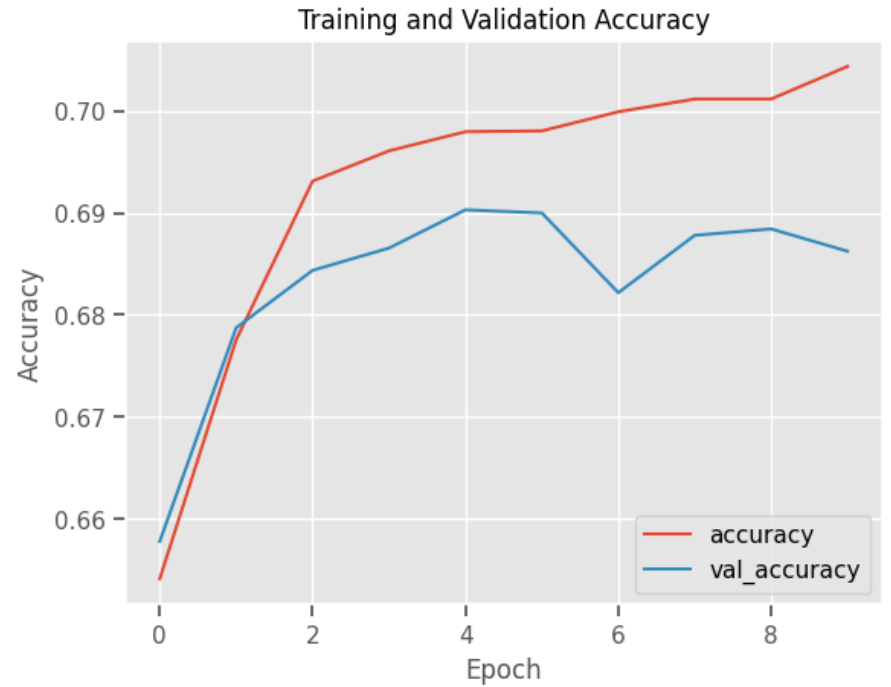


Matriz de confusión para la red neuronal con seis capas ocultas

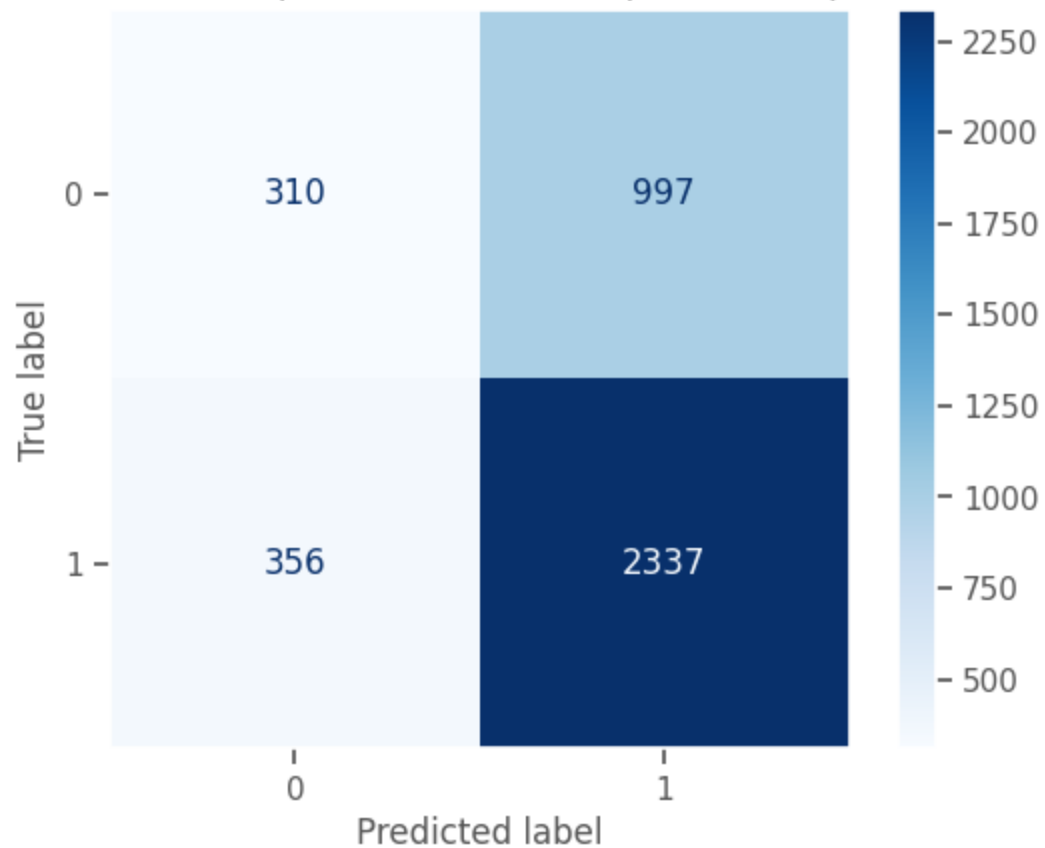


La segunda configuración fue :

1. Entrada (n_registros,25)
2. Diez capas ocultas (128 neuronas)
3. Una capa de salida – Función de activación sigmoide
4. Optimizador : Adam – lr : $1e-4$
5. Loss : binary cross entropy
6. 10 Epocas , batch-size : 32



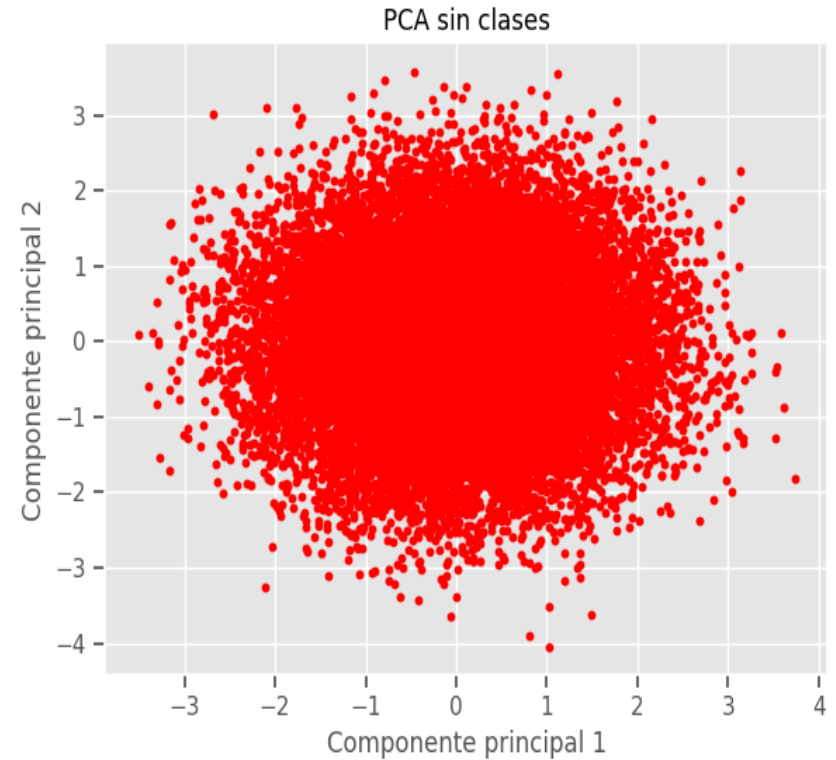
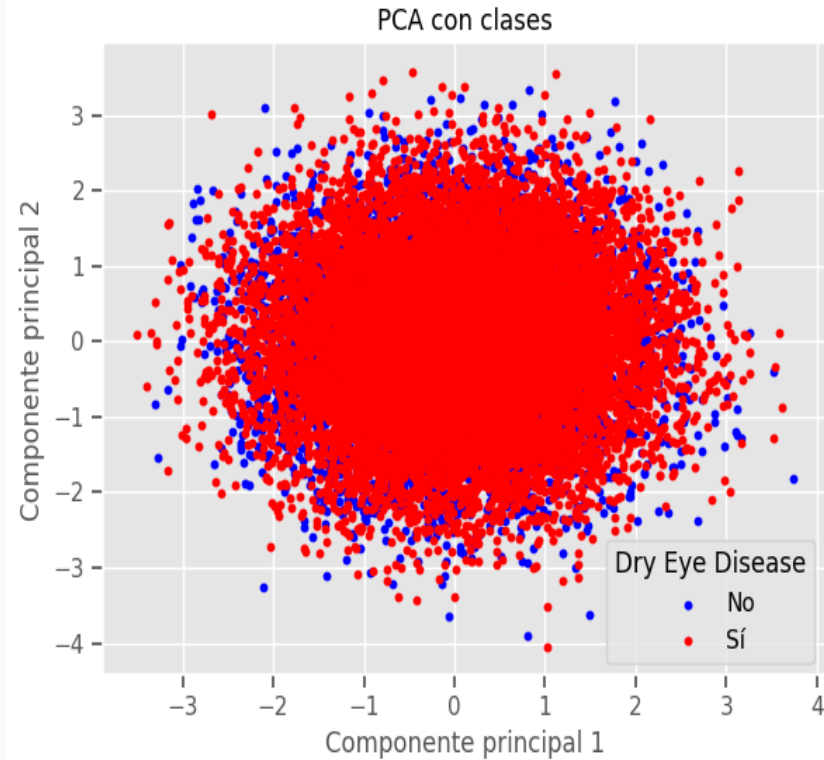
Matriz de confusión para la red neuronal para diez capas ocultas



06

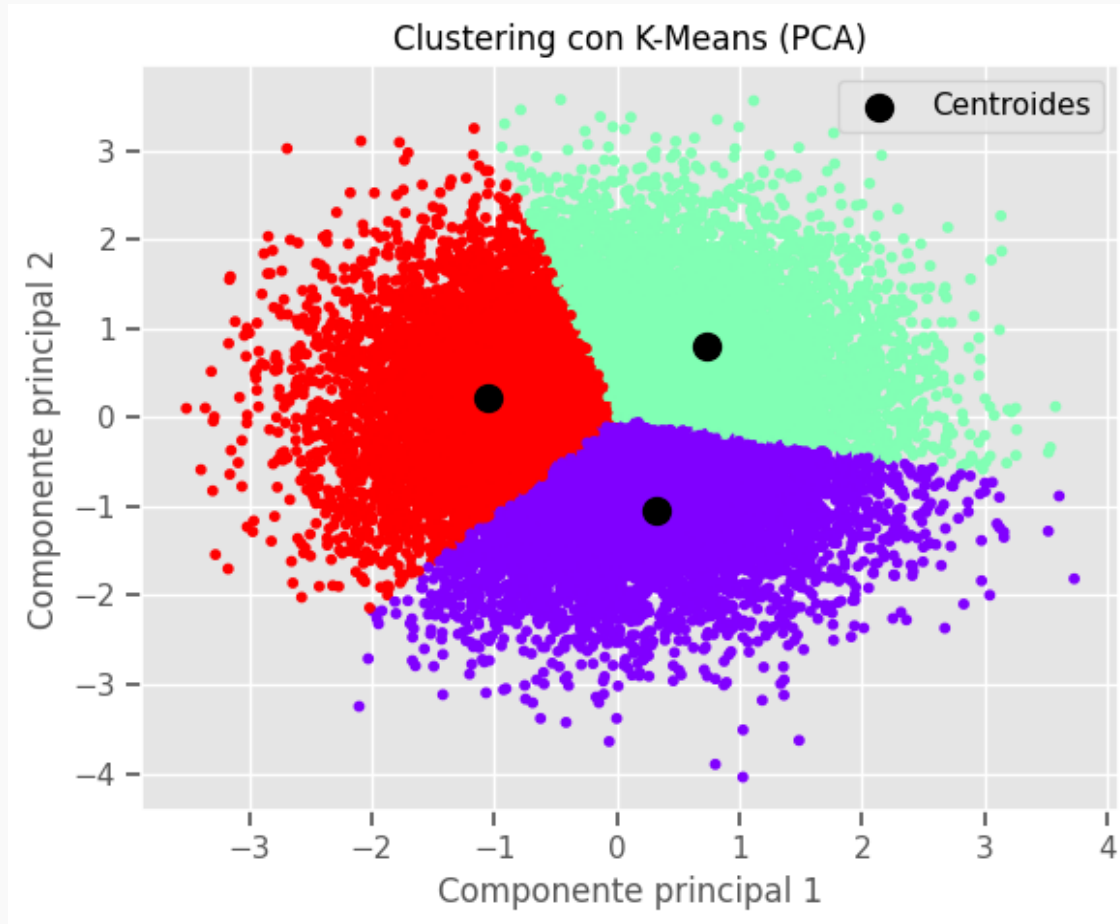
Aprendizaje no supervisado

PCA

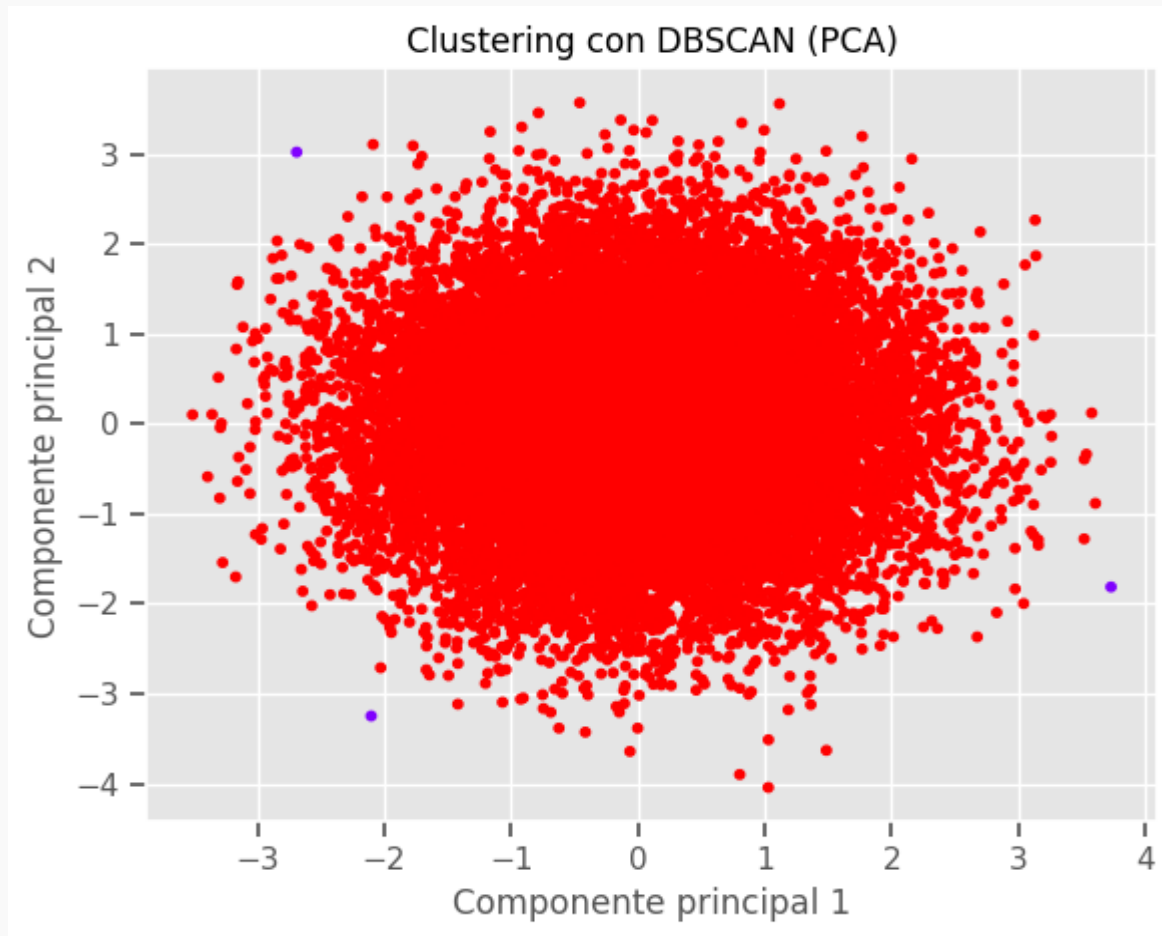


-
- Varianza explicada por componente:
 - Componente 1: 0.0426 (4.26%)
 - Componente 2: 0.0425 (4.25%)
 - Varianza total explicada por los 2 componentes:
0.0851 (8.51%)

K-means



DBSCAN



Conclusiones

- Los modelos de aprendizaje supervisado aplicados permiten identificar patrones relevantes entre variables relacionadas con la salud y la presencia del síndrome del ojo seco, demostrando que el enfoque basado en datos es viable para apoyar el diagnóstico preliminar.
- El modelo de SVM obtuvo el mejor rendimiento (70.1% de precisión), lo que sugiere que este algoritmo logra separar de forma más efectiva los casos positivos y negativos dentro del espacio de características del conjunto de datos.
- El Random Forest también mostró un rendimiento sólido (69.8%), validando la utilidad de técnicas de ensamblaje para mejorar la generalización y mitigar el sobreajuste, en comparación con modelos simples como el Árbol de Decisión (57.1%).
- El uso de preprocesamiento (label encoding y escalado) fue esencial para garantizar el correcto funcionamiento de los modelos, en especial SVM, que es sensible a las diferencias de escala entre variables.
- Un espacio de tan solo 2 componentes , proyectado a partir de los datos originales , no son suficientes para garantizar preservar una cantidad aceptable de varianza del espacio original , por lo cual , no se logra una gran distinción entre clases .

Gracias 👍