

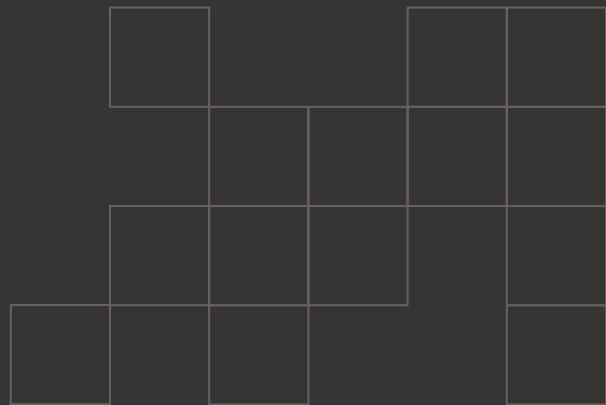
Proyecto IA

Predicción del éxito laboral
en estudiantes de diferentes
áreas del conocimiento

Daniel Alejandro Ayala Vallejo 2220084

Nelson Felipe Moreno 2220064

David Fernando Naranjo 2220046



Introducción

Este proyecto aborda un problema de clasificación utilizando el dataset Education & Career Success, que contiene información de 5.000 estudiantes, como edad, género, universidad, promedio, pasantías, habilidades y ofertas laborales.

El objetivo es predecir una variable categórica, es decir, conocer el cargo laboral que tendrá un estudiante en su primer trabajo, aplicando modelos clásicos de clasificación: Árbol de Decisión (DT), Random Forest (RF) y SVM, y evaluando su desempeño con métricas y curvas de aprendizaje.



Limpieza de los datos

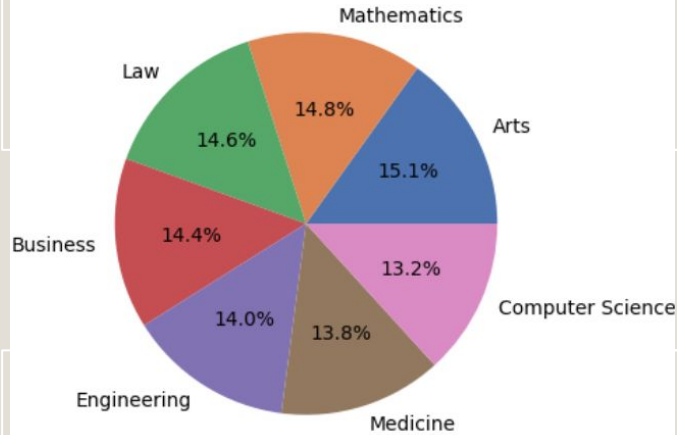


- Se elimina cualquier dato que sea nulo.
- Se elimina el dato “Other” para la columna Género.
- Se reemplazan “Male” y “Female” por 0 y 1 respectivamente.

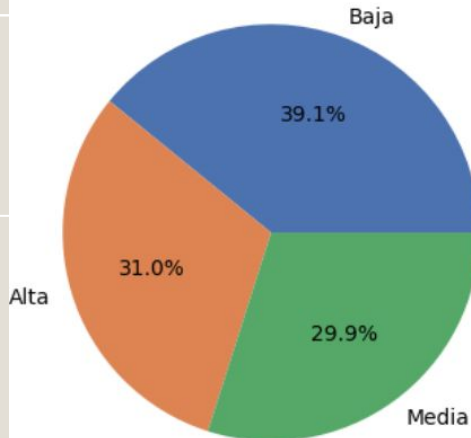


Distribución de algunas columnas

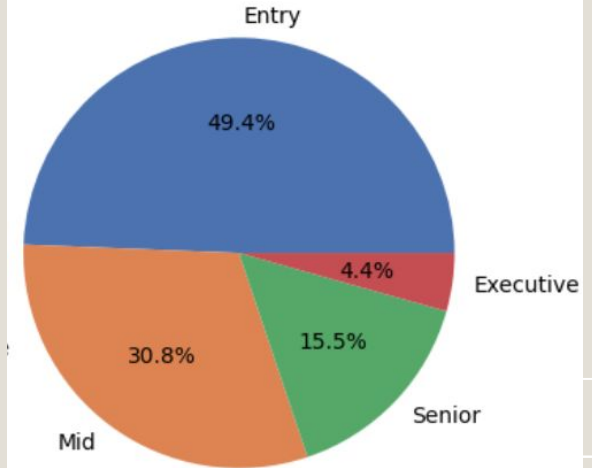
Distribucion de disciplinas



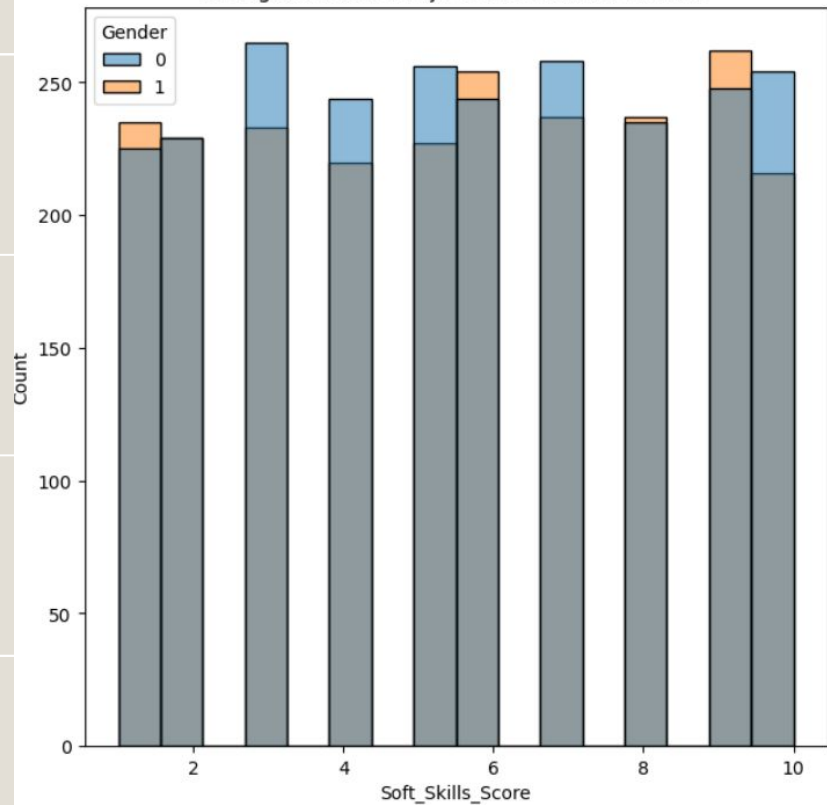
Distribucion de satisfaccion en la carrera



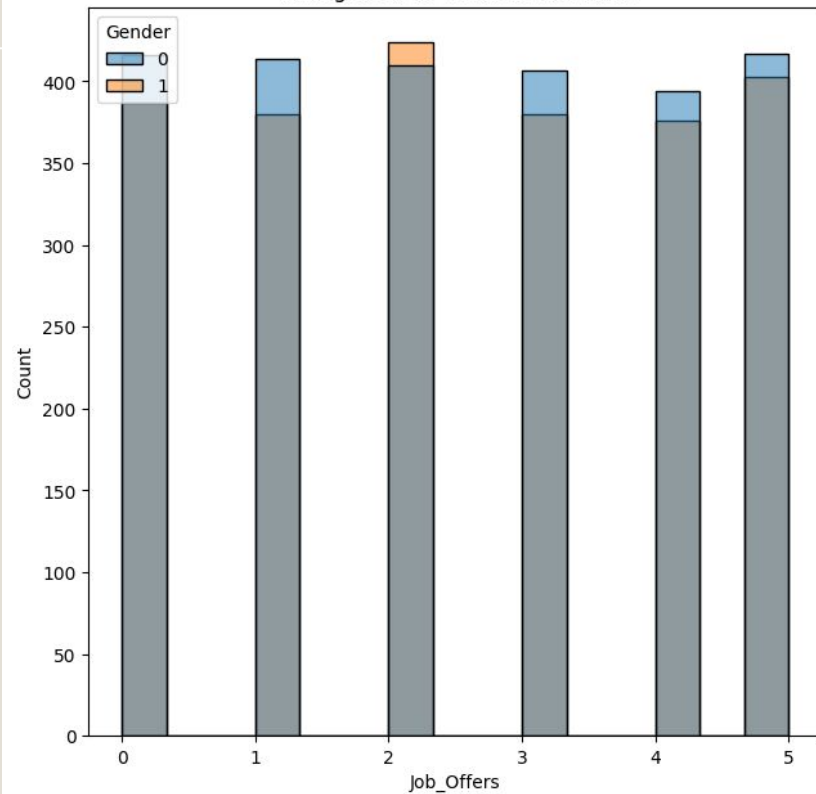
Distribucion nivel actual de trabajo



Histograma de Puntaje de Habilidades blandas



Histograma de Ofertas Laborales

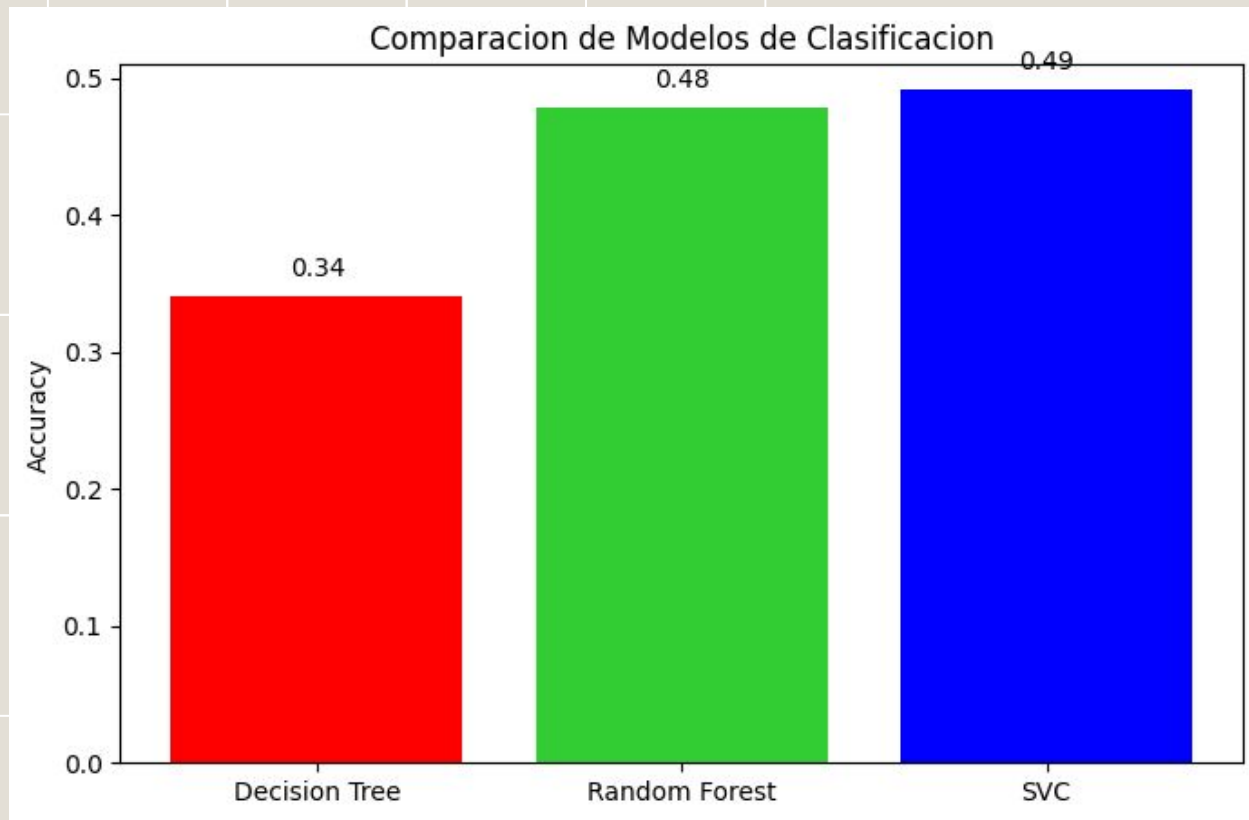


CLASIFICACIÓN

Para este conjunto de datos se realizó un proceso de clasificación a través de Machine Learning, con el fin de predecir la clase o ground truth “Current Job Level” para los estudiantes universitarios. Para ello se implementaron los métodos:

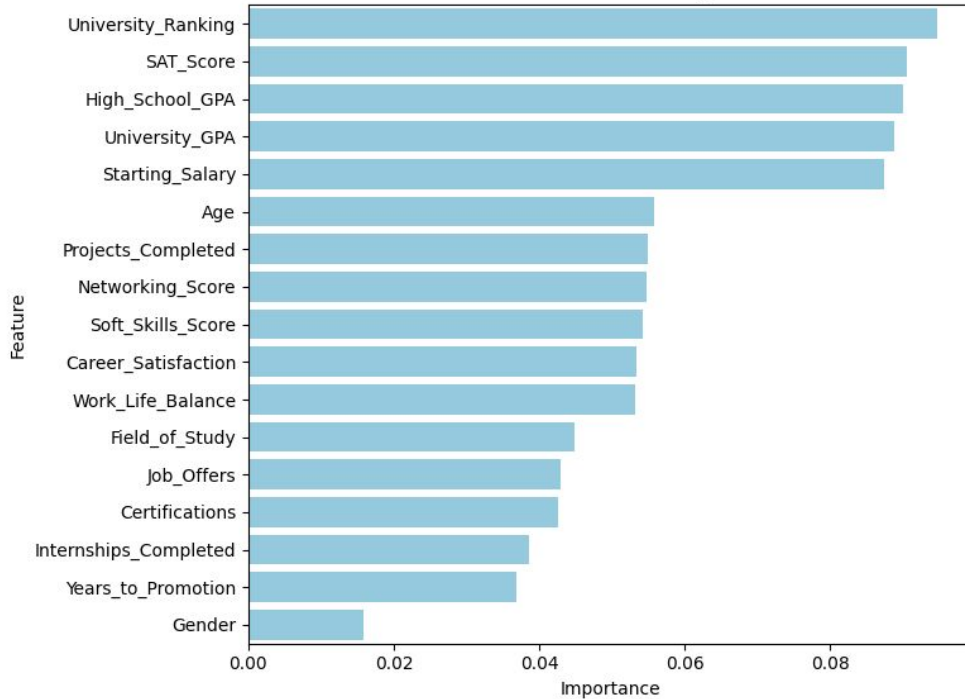
- Decision tree
- Random forest
- SVC (Support Vector Machine)

Adicionalmente como pre procesamiento se empleó el método “train_test_split” para seguir el pipeline recomendado por el docente y tener mejor desempeño en el proceso de clasificación.

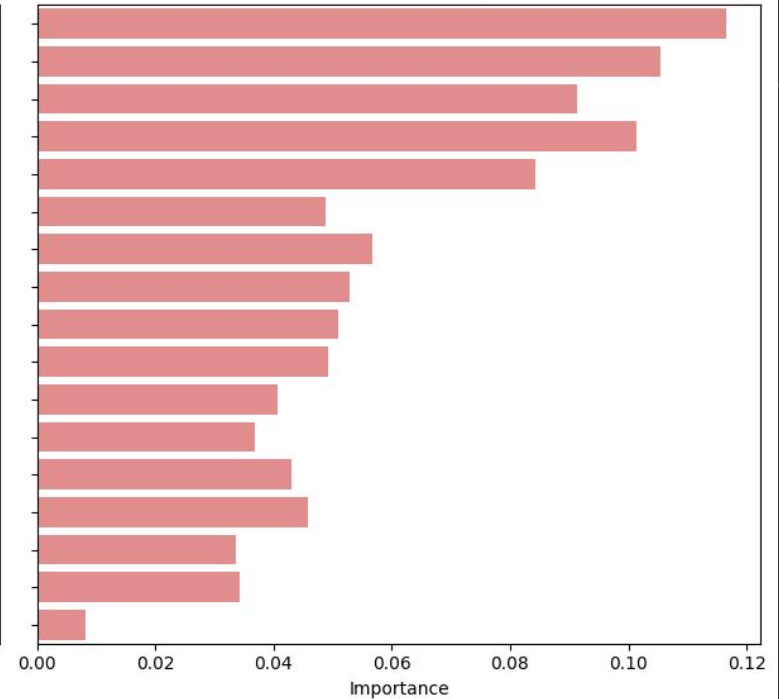


| Feature importances

Random Forest - feature importances



Decision Tree - feature importances



Curvas de aprendizaje

A continuación se realizó un proceso de curvas de aprendizaje para optimizar los procesos de clasificación de los 3 métodos mostrados anteriormente encontrando el MEJOR valor para el hiper parámetro clave a la hora de entrenar con cada uno

Árbol de Decisión (DT):

Se varió el valor de “max_depth” que indica que tan profundo ira el Decision tree para entrenar con el conjunto de datos.

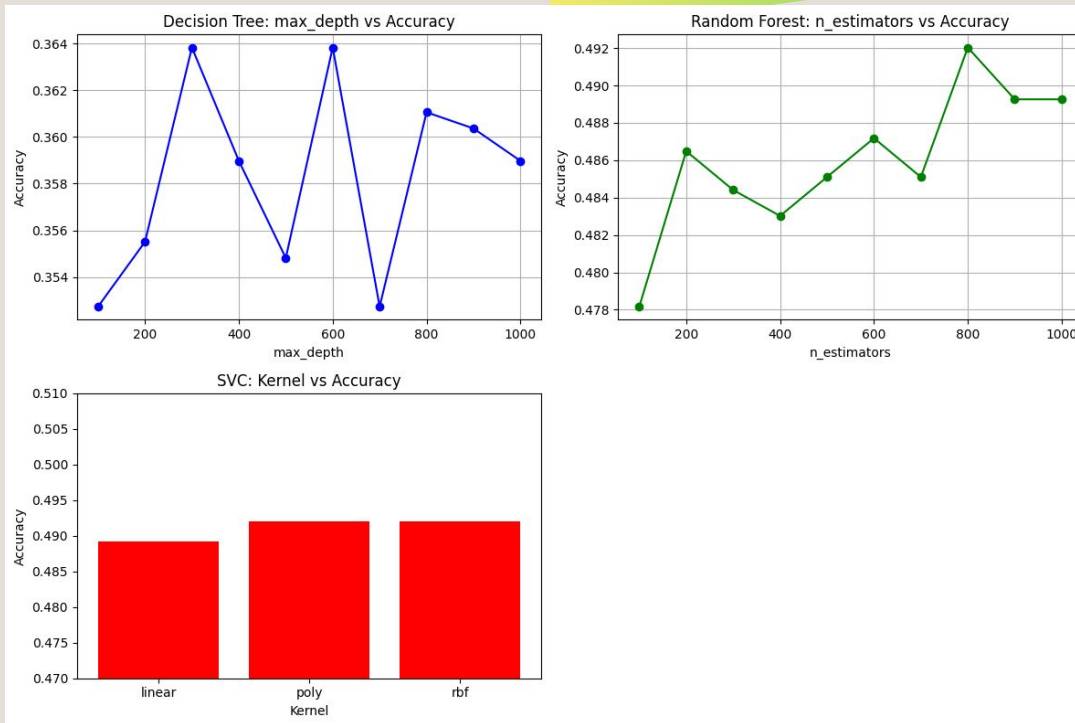
Random Forest (RF):

Se varió el parámetro “n_estimators” que indica el número de árboles de decisión que poseerá el método de RF.

SVM (SVC):

Se varió el “kernel” utilizado por el metodo para realizar el proceso de clasificación con el dataset.

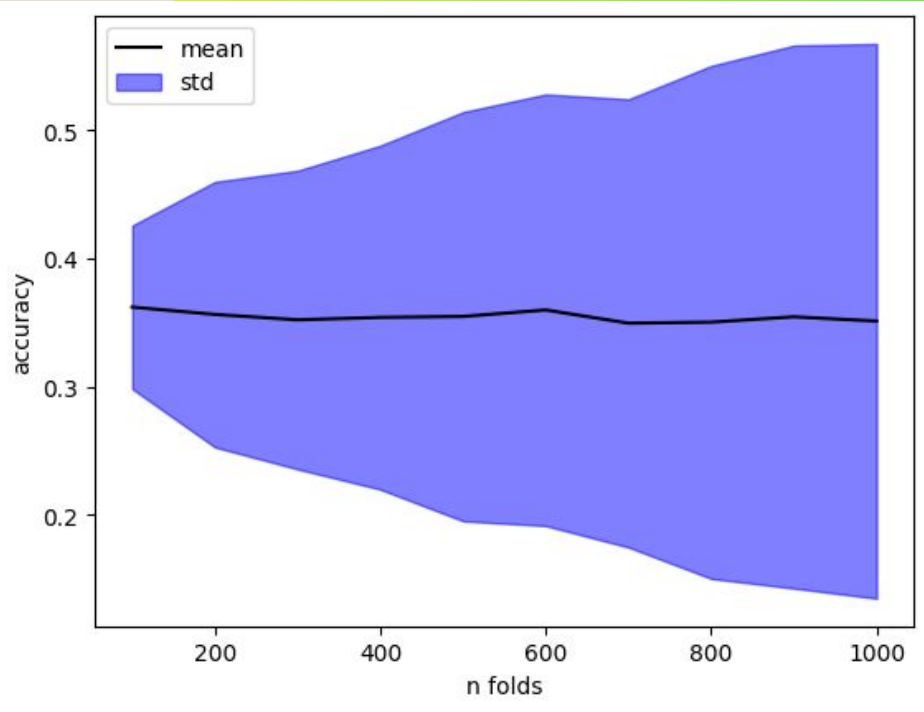
Resultado



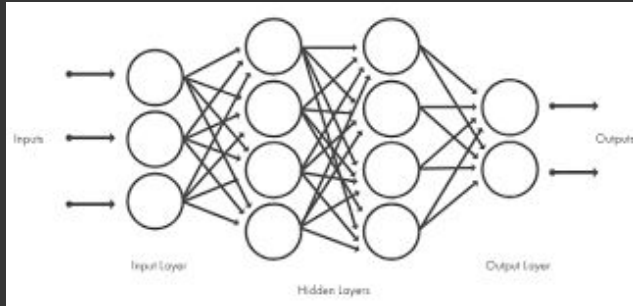
Cross Validation

Para este apartado, se volvera a realizar un procedimiento de curvas de aprendizaje pero unicamente para el metodo Decision tree, utilizando adicionalmente la funcion de validacion cruzada `cross_val_score()` para asi poder obtener el numero de ventana mas apropiado para que al igual se optimice el proceso de clasificacion.

Resultado



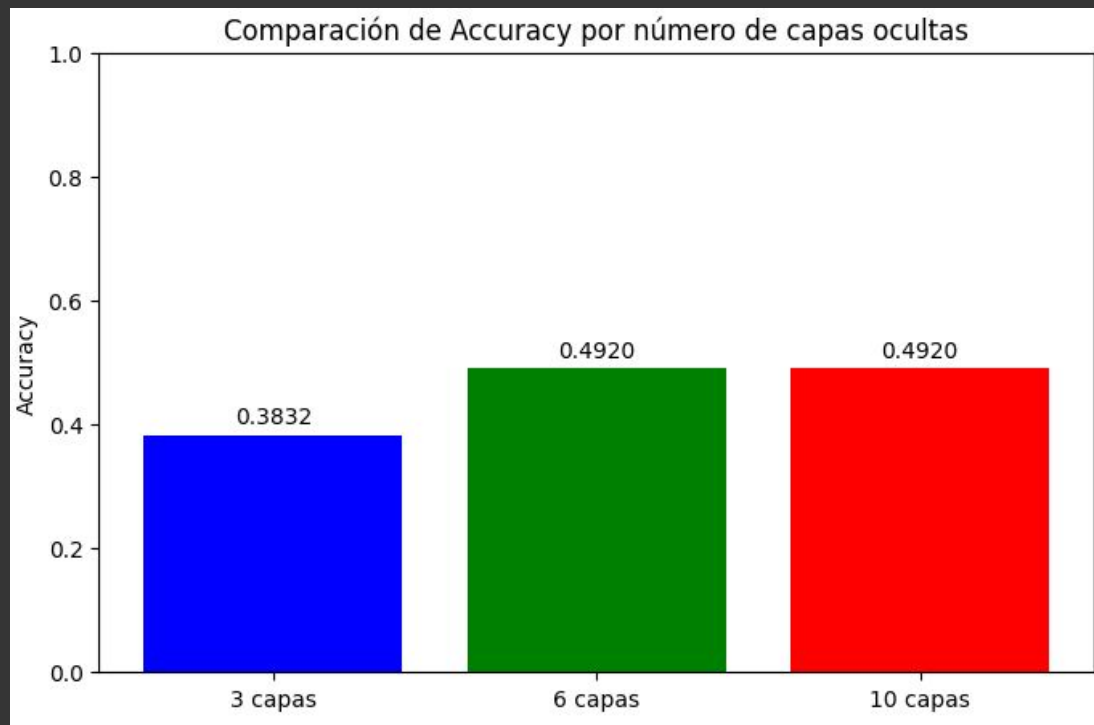
Deep Learning



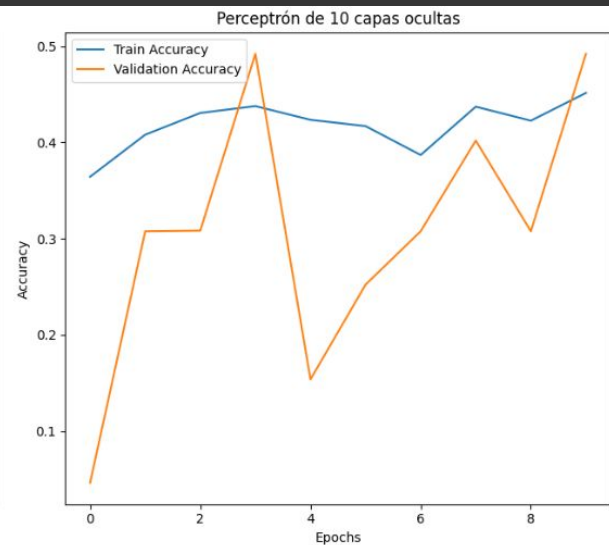
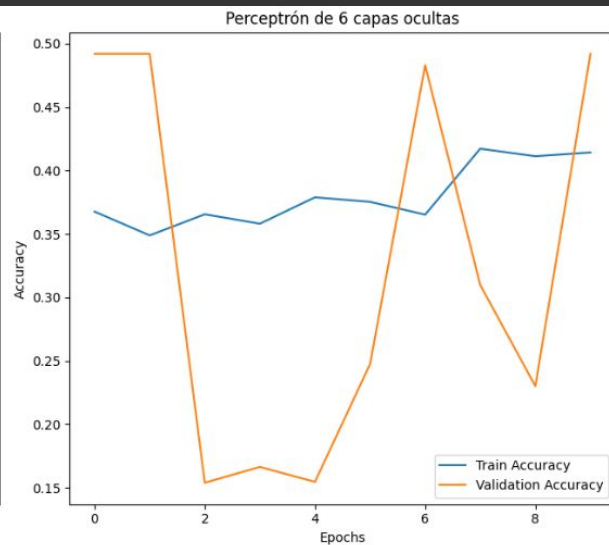
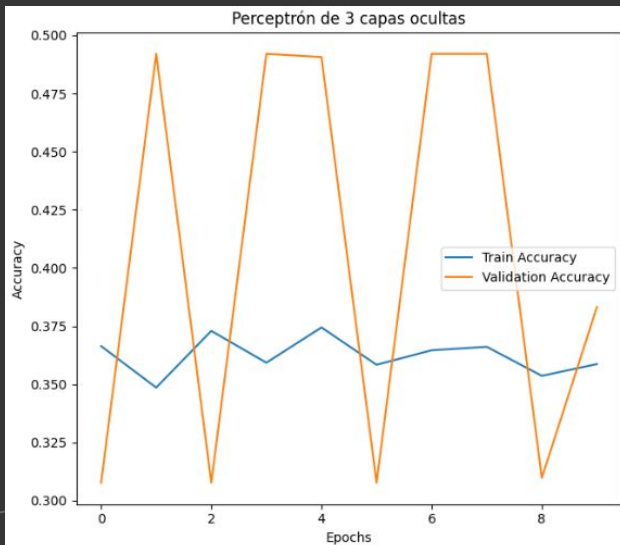
La implementación cuenta con 3 modelos, cada uno con 3, 6 y 10 capas de neuronas ocultas respectivamente. La métrica usada para evaluar el proceso de clasificación fue “accuracy”.

Para realizar Aprendizaje Profundo, se implementó un perceptrón multicapa con ayuda de la librería TensorFlow, esto para realizar una segunda vez un proceso de clasificación teniendo en cuenta el mismo ground truth “Current Job Level”.

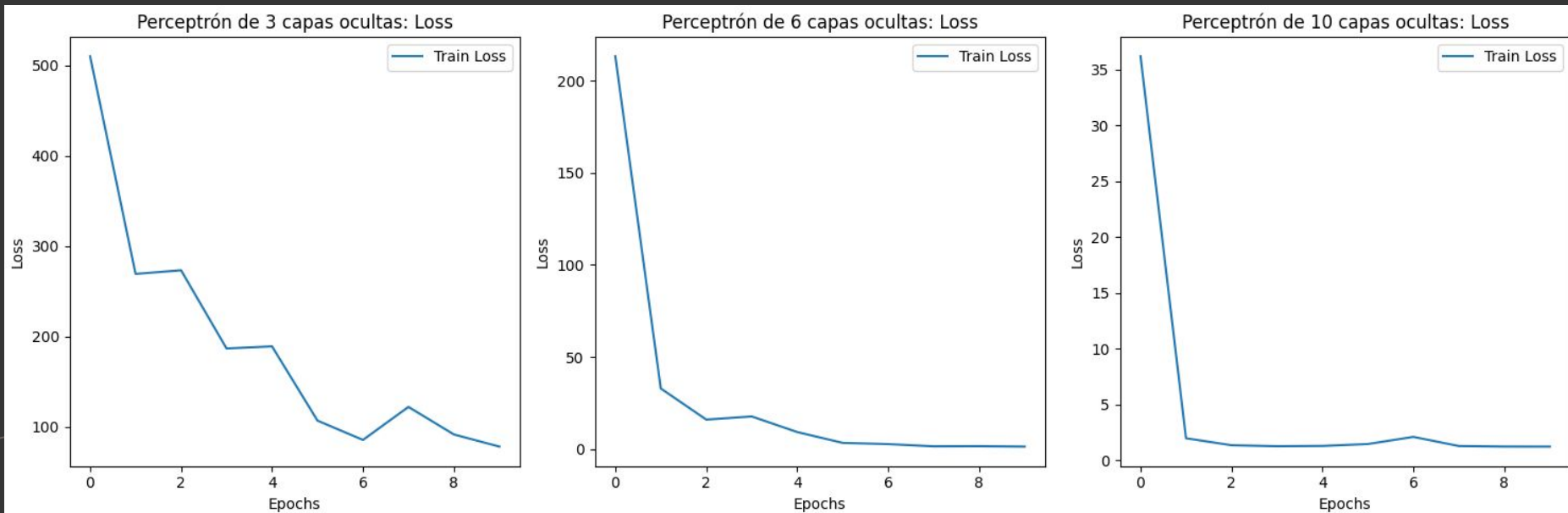
Resultado



Accuracy vs Epochs

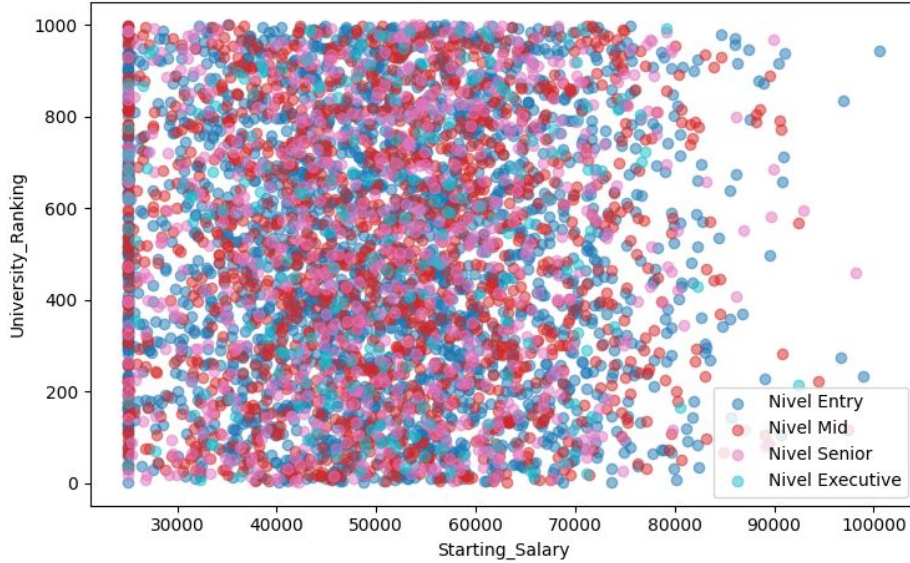


Loss vs Epochs

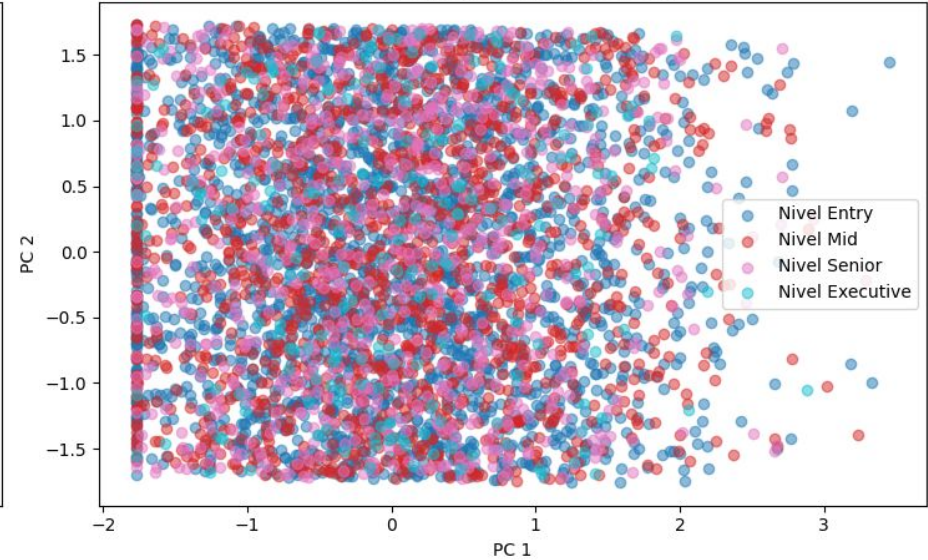


PCA

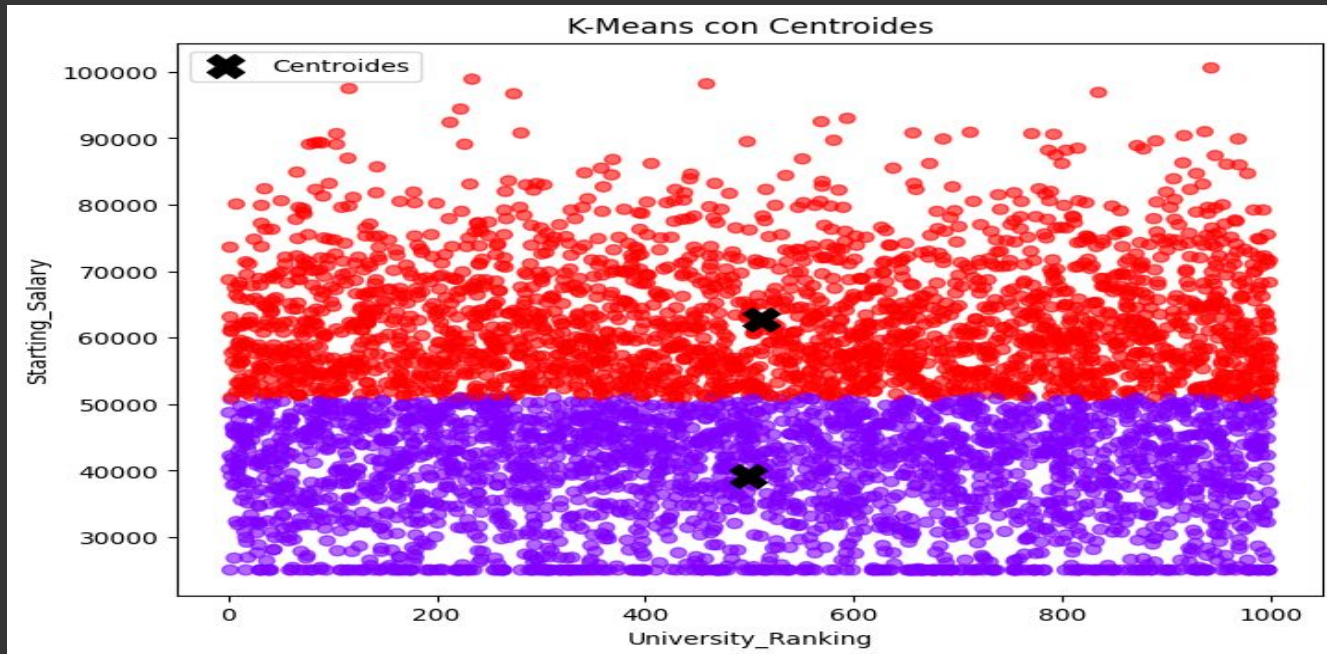
Características originales



Componentes principales (PCA)



K-Means con Centroides



DBSCAN

