

Memoria – Selección de modelo_empleabilidad_recien_titulados

1. Introducción

En este proyecto, el objetivo ha sido aplicar los conocimientos adquiridos sobre MLFlow en un escenario más realista. MLFlow es una herramienta clave para gestionar y hacer seguimiento de experimentos de machine learning, y en este caso, se ha utilizado para entrenar y evaluar modelos de clasificación sobre un dataset relacionado con la empleabilidad de los estudiantes.

Esta práctica busca predecir la empleabilidad de los estudiantes en función de sus habilidades, rendimiento académico y experiencia laboral. Para lograrlo, se han entrenado y comparado dos modelos de clasificación utilizando MLFlow para registrar todas las ejecuciones, parámetros y métricas asociadas a cada experimento.

2. Descripción del Dataset

El conjunto de datos utilizado para este proyecto proviene de Kaggle y contiene información sobre el rendimiento académico, formación y estado laboral de los estudiantes. El objetivo es predecir la empleabilidad de los estudiantes, es decir, si cuando terminan sus estudios salen con un puesto de trabajo o no.

Las principales columnas en el dataset incluyen:

- **CGPA:** Promedio general de calificaciones obtenido por el estudiante.
- **Internships:** Número de prácticas realizadas por el estudiante.
- **Projects:** Número de proyectos en los que el estudiante ha trabajado.
- **Workshops/Certifications:** Cursos y certificaciones obtenidas por el estudiante en línea.
- **AptitudeTestScore:** Puntaje en pruebas de aptitud utilizadas en procesos de reclutamiento.

- **SoftSkillRating:** Calificación de las habilidades blandas del estudiante, como la comunicación y trabajo en equipo.
- **ExtraCurricularActivities:** Actividades extracurriculares en las que participa el estudiante.
- **PlacementTraining:** Programas de formación recibidos por los estudiantes para prepararse para la búsqueda de empleo.
- **SSC_Marks:** Calificaciones obtenidas en la educación secundaria.
- **HCS_Marks:** Calificaciones obtenidas en bachillerato.
- **PlacementStatus:** Columna objetivo, con dos categorías: "Placed" (empleado) y "Not Placed" (no empleado).

3. Inicialización de MLFlow

Para poder registrar y hacer un seguimiento adecuado de los experimentos de Machine Learning realizados, se utilizó **MLFlow**, una plataforma de código abierto que facilita la gestión del ciclo de vida de los modelos de machine learning.

Se configuró un servidor local con la IP `http://127.0.0.1:8080` y se creó un experimento denominado '**Modelo_empleabilidad**' dónde guardaremos todas las ejecuciones de cada modelo.

Inicializamos MLFlow

```

from mlflow import MlflowClient
from pprint import pprint

client = MlflowClient(tracking_uri="http://127.0.0.1:8080")

# Provide an Experiment description that will appear in the UI
experiment_description = {
    "Este es el proyecto de predicción de empleabilidad estudiantil.",
    "Este experimento entrena modelos para clasificar el estado de inserción laboral de los estudiantes."
}

# Provide searchable tags that define characteristics of the Runs that
# will be in this Experiment
experiment_tags = {
    "project_name": "prediccion-empleabilidad",
    "store_dept": "produce",
    "team": "stores-el",
    "project_quarter": "Q3-2023",
    "mlflow.note.content": experiment_description,
}

# Comentamos porque ya lo tenemos creado
Create the Experiment, providing a unique name
produce_empleabilidad_experiment = client.create_experiment(
    name="Modelo_empleabilidad", tags=experiment_tags)

# Sets the current active experiment to the "Modelo_empleabilidad" experiment and
# returns the Experiment metadata
apple_experiment = mlflow.set_experiment("Modelo_empleabilidad")

# Define a run name for this iteration of training.
# If this is not set, a unique name will be auto-generated for your run.
run_name = "empleabilidad_rf_test"

# Define an artifact path that the model will be saved to.
artifact_path = "rf_empleabilidad"

```

4. Preprocesamiento y análisis de los datos

En esta parte, analizamos los datos y nos dimos cuenta de que algunas variables tenían escalas muy diferentes. Esto podría hacer que los modelos no funcionaran de la mejor manera, por lo que decidimos estandarizar los datos. De esta forma, los modelos pueden trabajar mejor con ellos.

5. Modelos de clasificación

Para la clasificación, utilizamos cuatro modelos: DecisionTreeClassifier, RandomForestClassifier, KNeighborsClassifier (KNN) y XGBoostClassifier. Realizamos entre 10 y 15 ejecuciones por modelo, variando los hiperparámetros con el objetivo de obtener el mayor accuracy posible. Todas las ejecuciones fueron registradas en MLFlow; en los dos primeros modelos, el registro se realizó manualmente, mientras que en los otros dos, se hizo de manera automática.

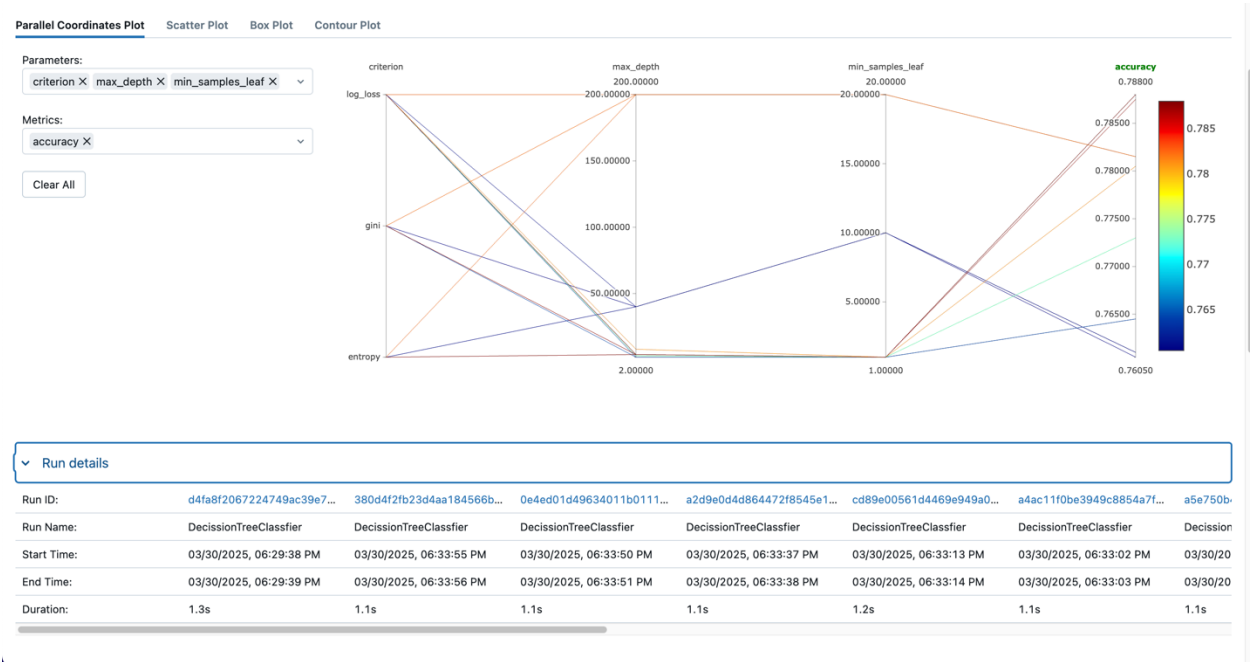
<input type="checkbox"/>	Run Name	Created	Dataset	Duration	Source	Models
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.1s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.1s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.2s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.1s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.1s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.2s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.2s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.1s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.2s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.1s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.1s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.1s	ipykern...	sklearn
<input type="checkbox"/>	DecisionTreeClassifier	1 hour ago	-	1.3s	ipykern...	sklearn
<input type="checkbox"/>	Run Name	Created	Dataset	Duration	Source	Models
<input type="checkbox"/>	thoughtful-squid-284	38 minutes ago	dataset (594874af) Eval , dataset (c...	1.5s	ipykern...	xgboost
<input type="checkbox"/>	stylish-swan-709	38 minutes ago	dataset (594874af) Eval , dataset (c...	1.2s	ipykern...	xgboost
<input type="checkbox"/>	useful-slug-903	38 minutes ago	dataset (dabee8d2) Train , dataset ...	1.3s	ipykern...	xgboost
<input type="checkbox"/>	handsome-bat-778	38 minutes ago	dataset (dabee8d2) Train , dataset ...	1.4s	ipykern...	xgboost
<input type="checkbox"/>	serious-smelt-341	39 minutes ago	dataset (dabee8d2) Train , dataset ...	1.5s	ipykern...	xgboost
<input type="checkbox"/>	defiant-panda-523	51 minutes ago	dataset (594874af) Eval	2.1s	ipykern...	sklearn
<input type="checkbox"/>	monumental-hound-976	51 minutes ago	dataset (594874af) Eval	2.1s	ipykern...	sklearn
<input type="checkbox"/>	useful-snipe-537	51 minutes ago	dataset (594874af) Eval	2.0s	ipykern...	sklearn
<input type="checkbox"/>	unique-vole-403	52 minutes ago	dataset (594874af) Eval	2.0s	ipykern...	sklearn
<input type="checkbox"/>	unequaled-midge-134	52 minutes ago	dataset (594874af) Eval	2.0s	ipykern...	sklearn
<input type="checkbox"/>	dazzling-whale-102	52 minutes ago	dataset (594874af) Eval	1.7s	ipykern...	sklearn
<input type="checkbox"/>	puzzled-smelt-395	52 minutes ago	dataset (594874af) Eval	2.3s	ipykern...	sklearn
<input type="checkbox"/>	RandomForestClassifier	58 minutes ago	-	1.2s	ipykern...	sklearn

44 matching runs

5. Comparación de resultados

Para la comparación de los resultados, comenzamos seleccionando todas las ejecuciones de un mismo modelo y comparando las métricas obtenidas en cada una. De entre todas esas ejecuciones, elegimos la que mostró las mejores métricas. Una vez que seleccionamos la mejor ejecución de cada modelo, comparamos los resultados entre los cuatro modelos utilizados. Finalmente, decidimos cuál de ellos ofrecía el mejor rendimiento global.

DecissionTreeClassifier



Los mejores hiperparámetros para este modelo fueron los siguientes:

criterion	gini	accuracy	0.788
max_depth	4	f1	0.741
min_samples_leaf	1	precision	0.759
min_samples_split	2	recall	0.723
random_state	42		

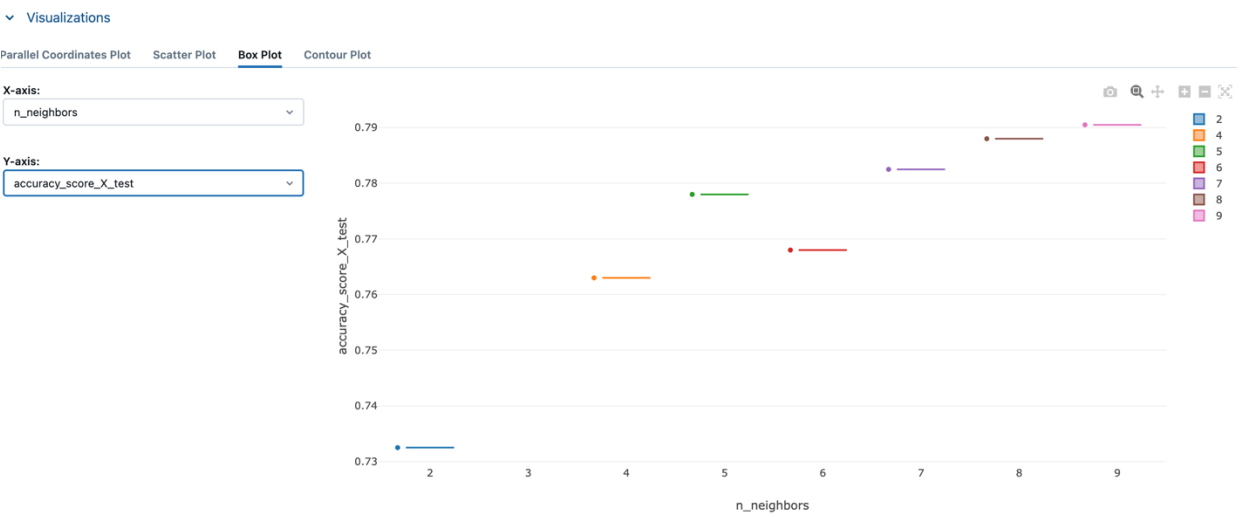
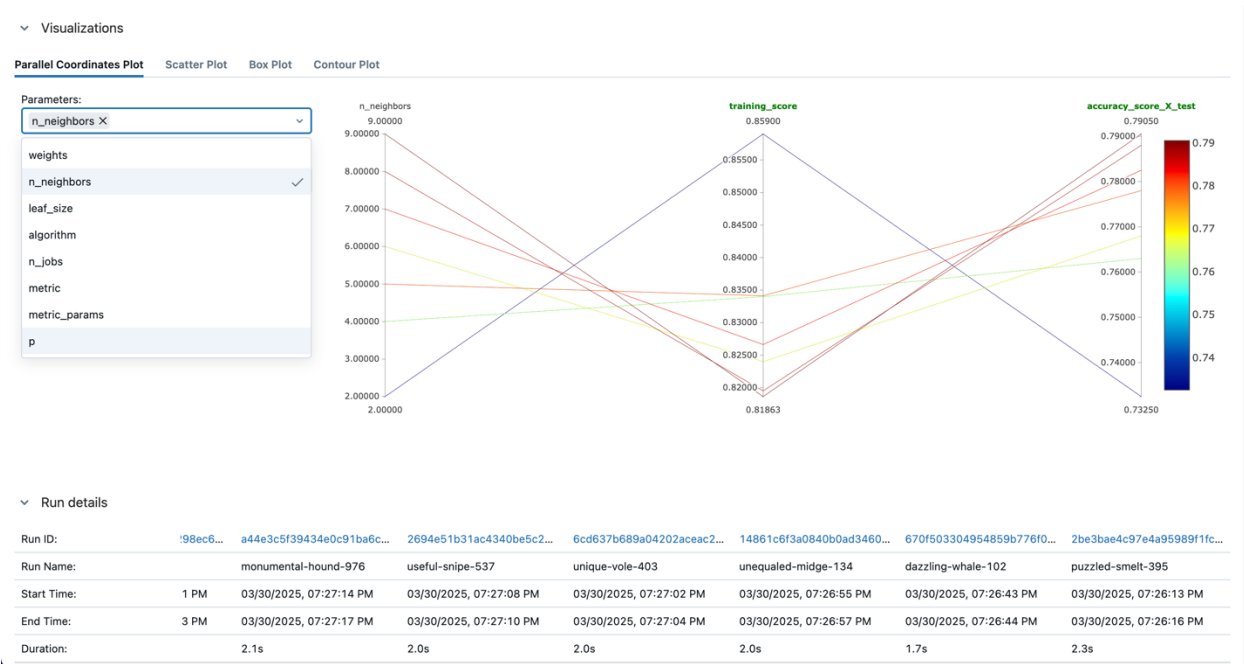
RandomForestClassifier



Los mejores hiperparámetros para este modelo fueron los siguientes:

criterion	entropy	accuracy	0.796
min_samples_leaf	1	f1	0.75
min_samples_split	2	precision	0.769
n_estimators	5	recall	0.732

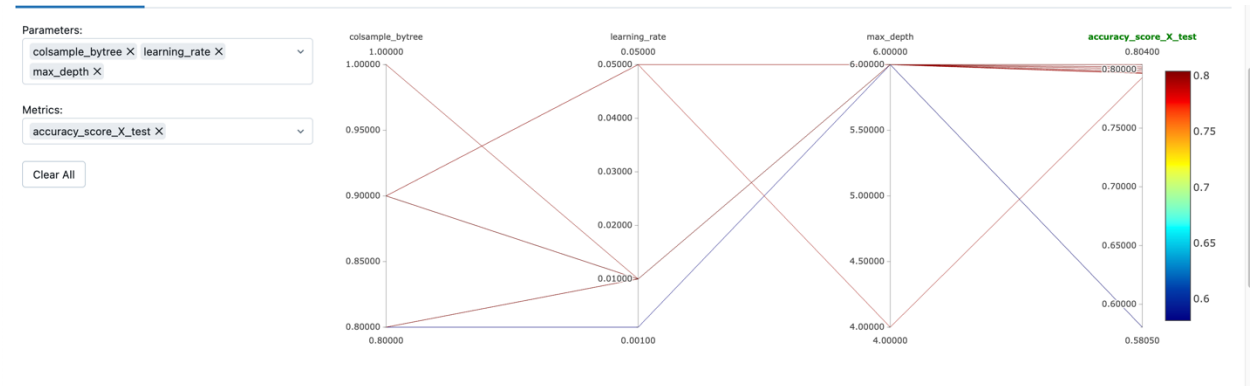
KNeighborsClassifier



Los mejores hiperparámetros para este modelo fueron los siguientes:

accuracy_score_X_test	0.791
f1_score_X_test	0.746
precision_score_X_test	0.759
recall_score_X_test	0.734
training_accuracy_score	0.819
training_f1_score	0.818
training_log_loss	0.373
training_precision_score	0.818
training_recall score	0.819

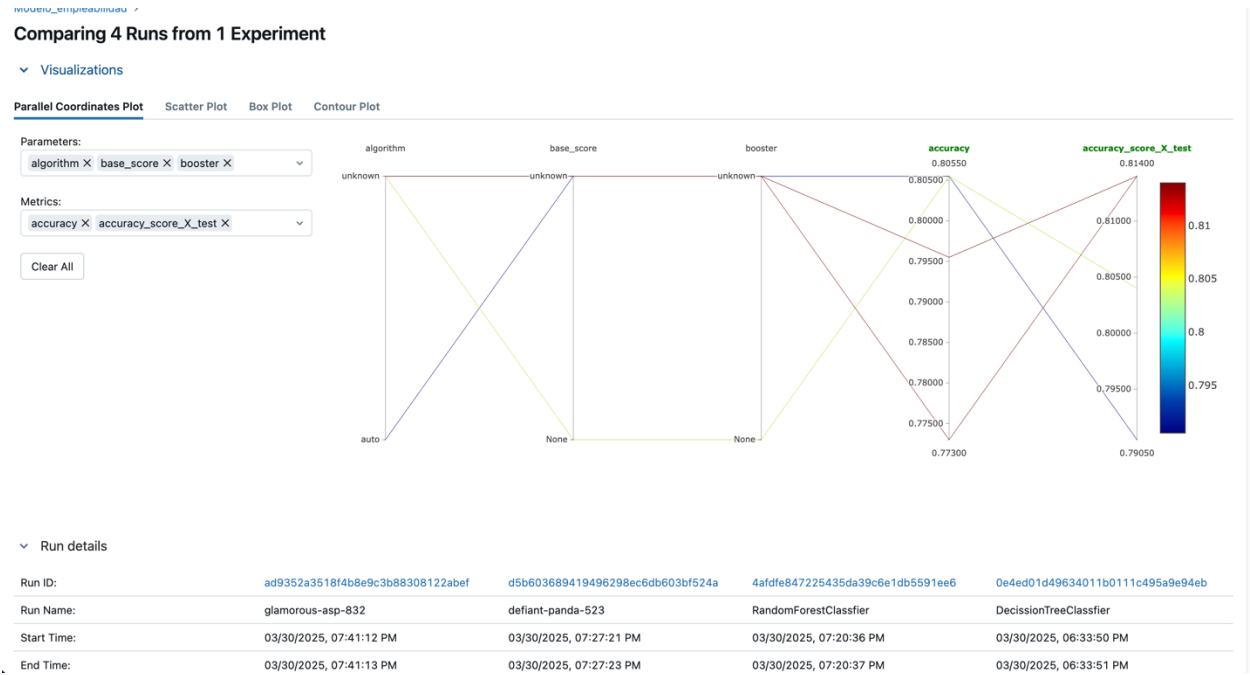
XGBoost



Los mejores hiperparámetros para este modelo fueron los siguientes:

colsample_bytree	0.9	accuracy_score_X_test	0.804
learning_rate	0.01	f1_score_X_test	0.755
max_depth	6	precision_score_X_test	0.793
num_boost_round	100	recall_score_X_test	0.721

Comparación de las mejores ejecuciones de cada modelo



6. Conclusión

Nos decidimos por el mejor modelo, en este caso, la ejecución identificada como (ad9352a3518f4b8e9c3b88308122abef) – glamorous-asp-832, que alcanzó un accuracy de 0.804. Esta ejecución corresponde al modelo de XGBoost, el cual resultó ser el más eficiente en cuanto a rendimiento.