

Daniel Yaruro Contreras
Juan Toloza Parada
Oscar Ortega Lozano



PROYECTO MACHINE LEARNING CLASSIFICATION



Prediction Placement Status



INTRODUCTION

PLACEMENT PREDICTION

+

Este dataset se centra en la predicción de la empleabilidad de estudiantes en función de su desempeño académico, experiencia práctica y habilidades blandas. Contiene información sobre calificaciones, participación en pasantías, proyectos, certificaciones y actividades extracurriculares, así como puntajes en pruebas de aptitud y entrenamientos específicos para la colocación laboral. Su objetivo es analizar qué factores influyen en la posibilidad de que un estudiante consiga un empleo tras finalizar sus estudios.

+





CLASSIFICATION PLACEMENT STATUS

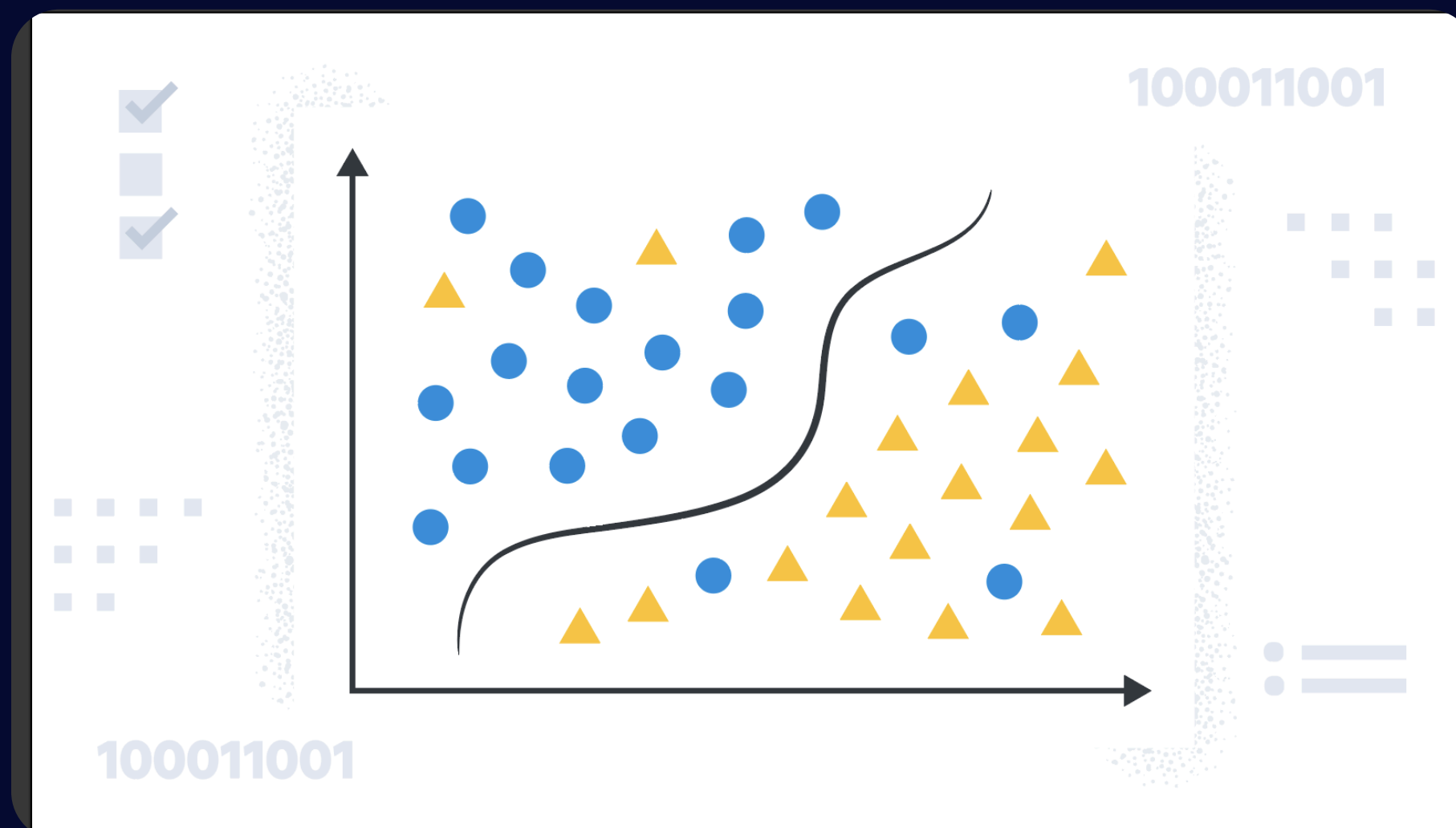
Este campo (Placement Status) indica si un estudiante logró o no conseguir un trabajo después del proceso de reclutamiento universitario. Los valores típicos de esta columna son:

Placed

El estudiante fue contratado por una empresa.

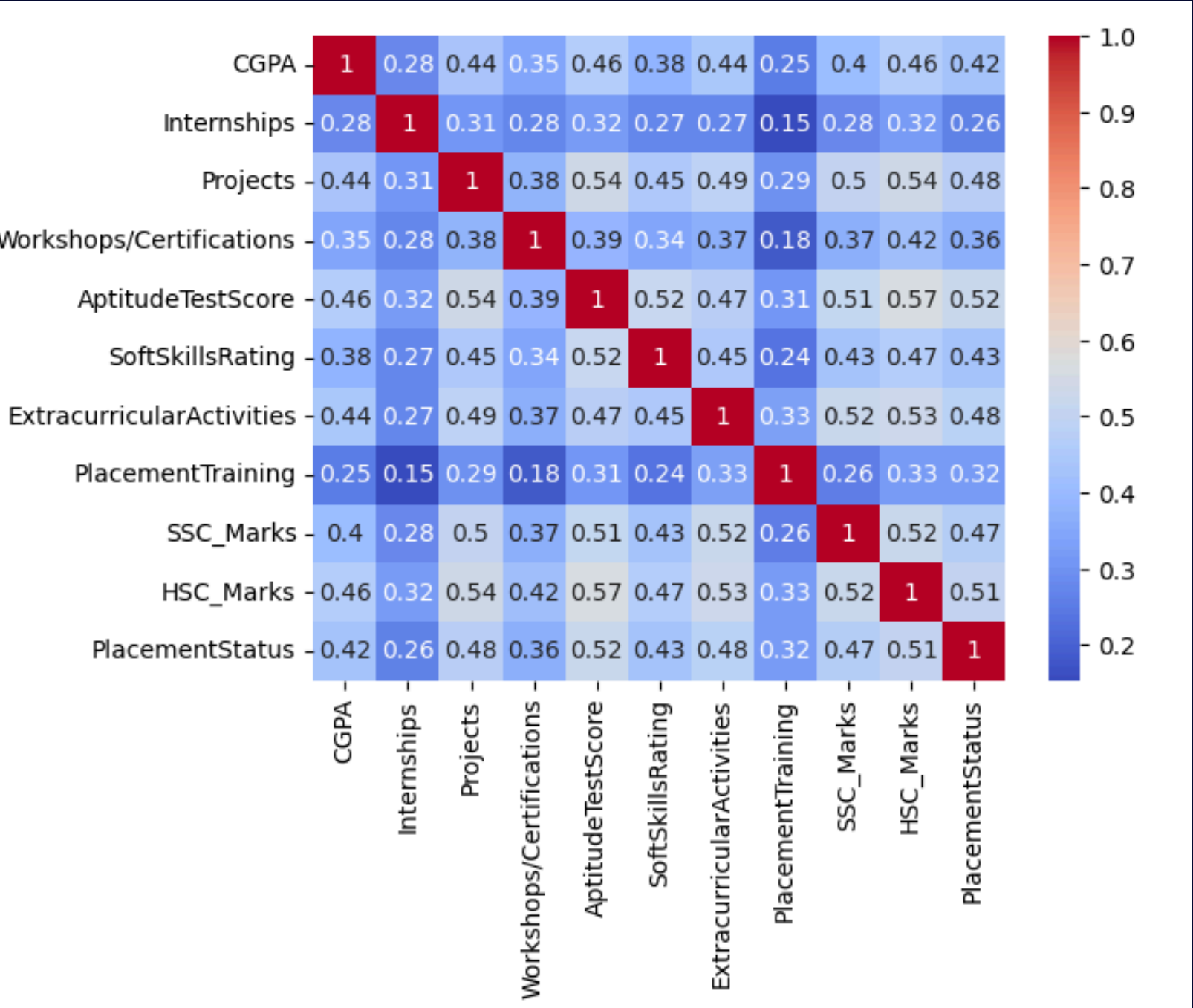
Not Placed

El estudiante no consiguió un empleo durante el proceso de colocación.



CARACTERISTICAS RELACIONADAS

Haremos nuestro primer modelo basandonos en variables que estan mas correlacionadas con el ground truth las cuales son:



01

CGPA
Correspondientes al CGPA de cada estudiante del dataset

02

Projects
Esta columna representa la cantidad de proyectos que han hecho los estudiantes

03

Extracurricular Activities
Helps programmers keep track of code changes, collaborate on projects, and maintain consistency.

04

Aptitude Test Score
Esta columna representa la calificacion que obtuvieron los estudiantes en el test de actitud que hace la empresa o organizacion

05

SSC_Marks
Dato numerico que representa las calificaciones obtenidas por el estudiante en la educacion secundaria

06

HSC_Marks
Dato numerico que epresenta las calificaciones obtenidas en la educacion preuniversitaria

MACHINE LEARNING MODELS

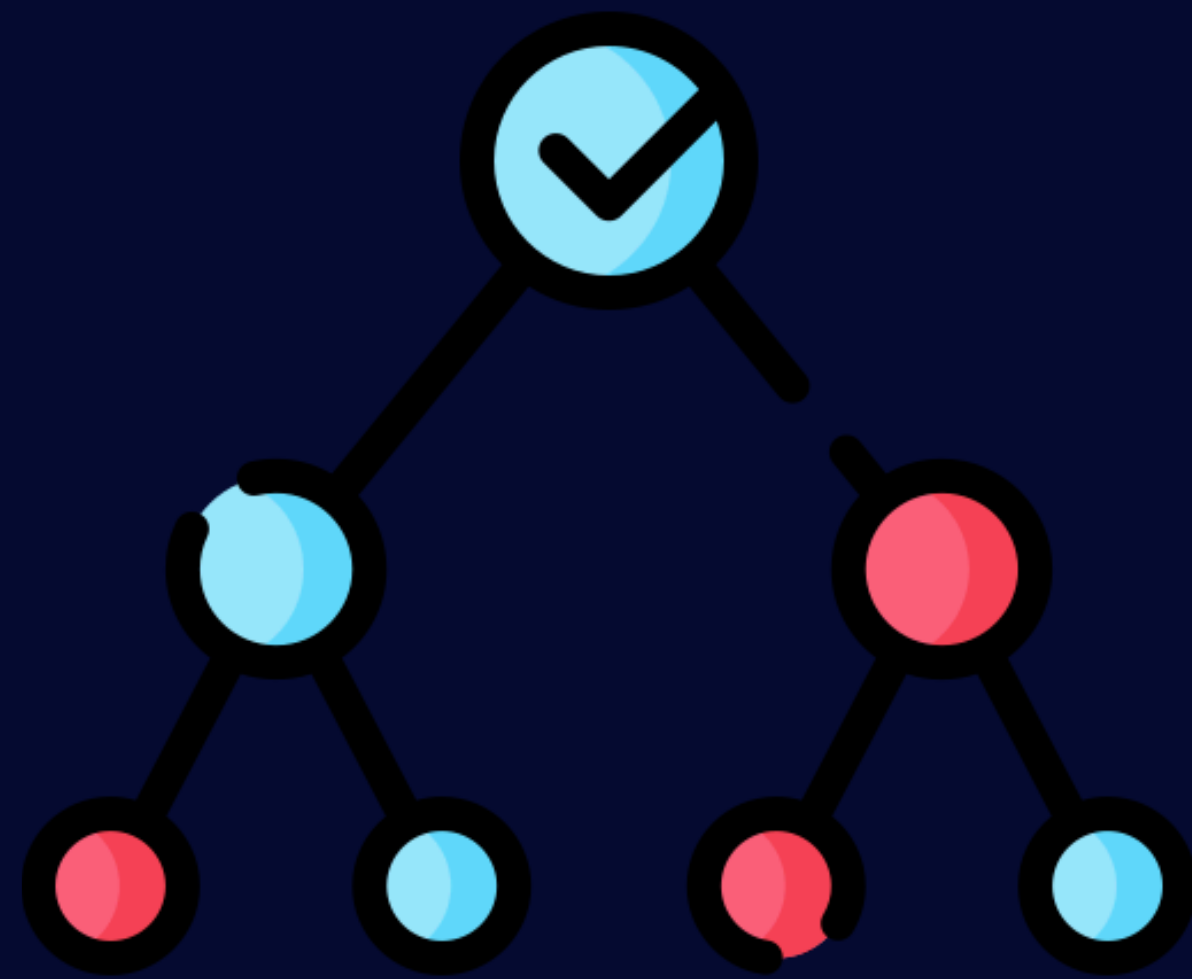
TRAIN TEST SPLIT PROCESS

```
1 columnas_relevantes = ["CGPA",  
2   "Projects",  
3   "ExtracurricularActivities",  
4   "AptitudeTestScore",  
5   "SSC_Marks",  
6   "HSC_Marks"]  
7 X = df[columnas_relevantes]  
8 y = df["PlacementStatus"]  
9  
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```
El tamaño del X_train es:  (8000, 6)  
El tamaño del y_train es:  (8000,)  
El tamaño del X_test es:   (2000, 6)  
El tamaño del y_test es:   (2000,)
```

+

=



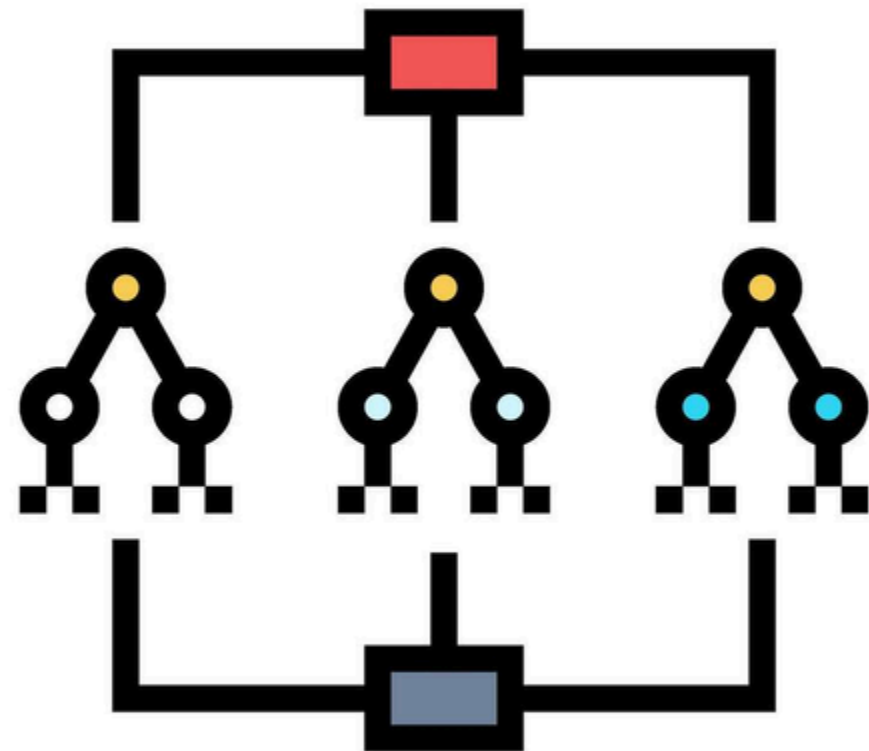
DESICION TREE

```
1 #@title Desicion Tree
2 tree = DecisionTreeClassifier(random_state=42)
3 tree.fit(X_train, y_train)
4 y_pred = tree.predict(X_test)
5 accuracy = accuracy_score(y_test, y_pred)
6 print("Accuracy:", accuracy)
```

Accuracy: 0.7245

+

=



```
1 #@title Random Forest
2 rf = RandomForestClassifier(random_state=42)
3 rf.fit(X_train, y_train)
4 y_pred = rf.predict(X_test)
5 accuracy = accuracy_score(y_test, y_pred)
6 print("Accuracy:", accuracy)
```

RANDOM FOREST

Accuracy: 0.772

+

≡



SUPPORT VECTOR MACHINE

```
1 #@title Support Vector Machines
2 svc = SVC(random_state=42)
3 svc.fit(X_train, y_train)
4 y_pred = svc.predict(X_test)
5 accuracy = accuracy_score(y_test, y_pred)
6 print("Accuracy:", accuracy)
```

Accuracy: 0.776

N-Folds



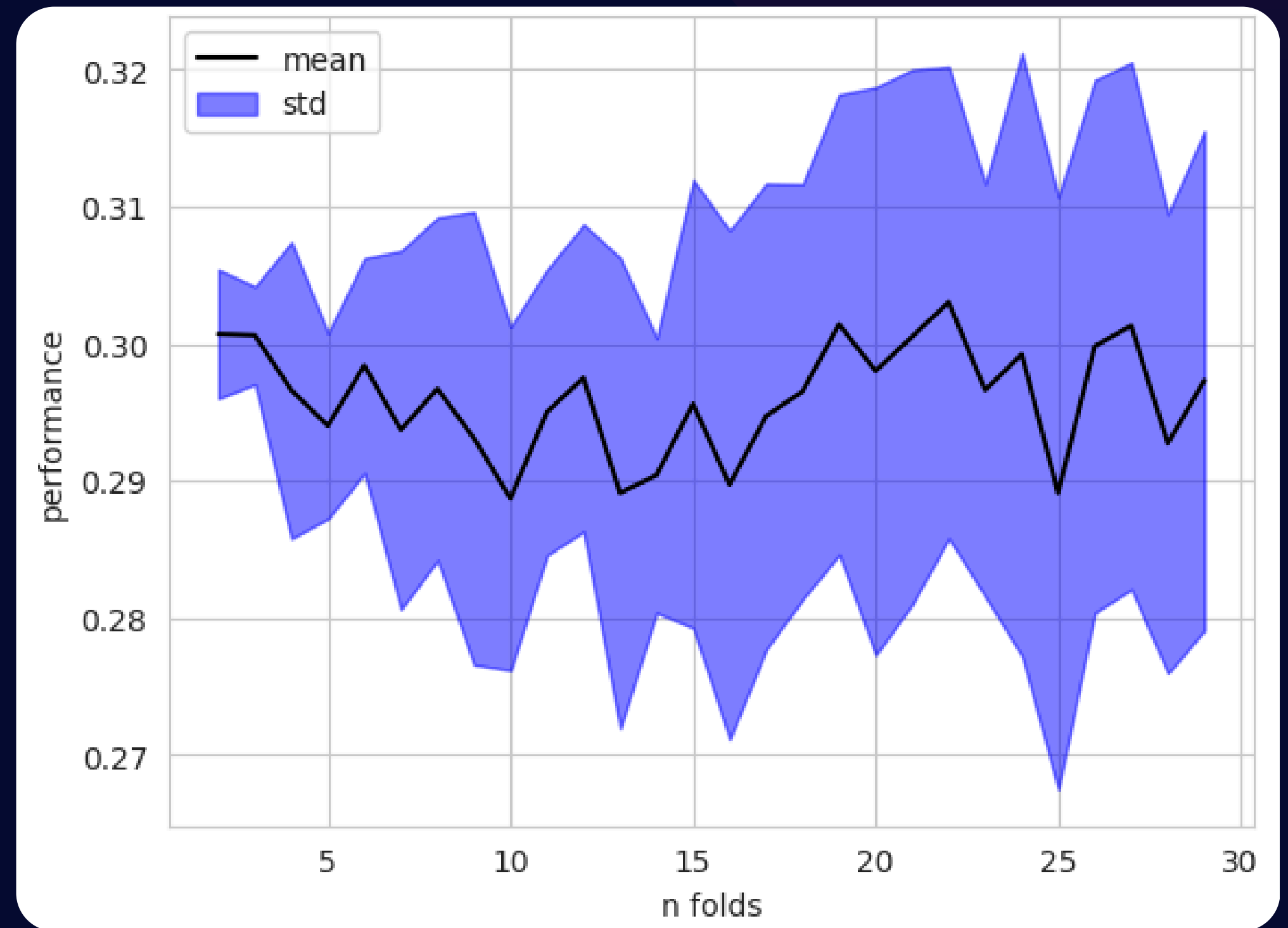
DESICION TREE

N-folds = 10

$\sigma = 0.012522379965485792$

+

+



N-Folds



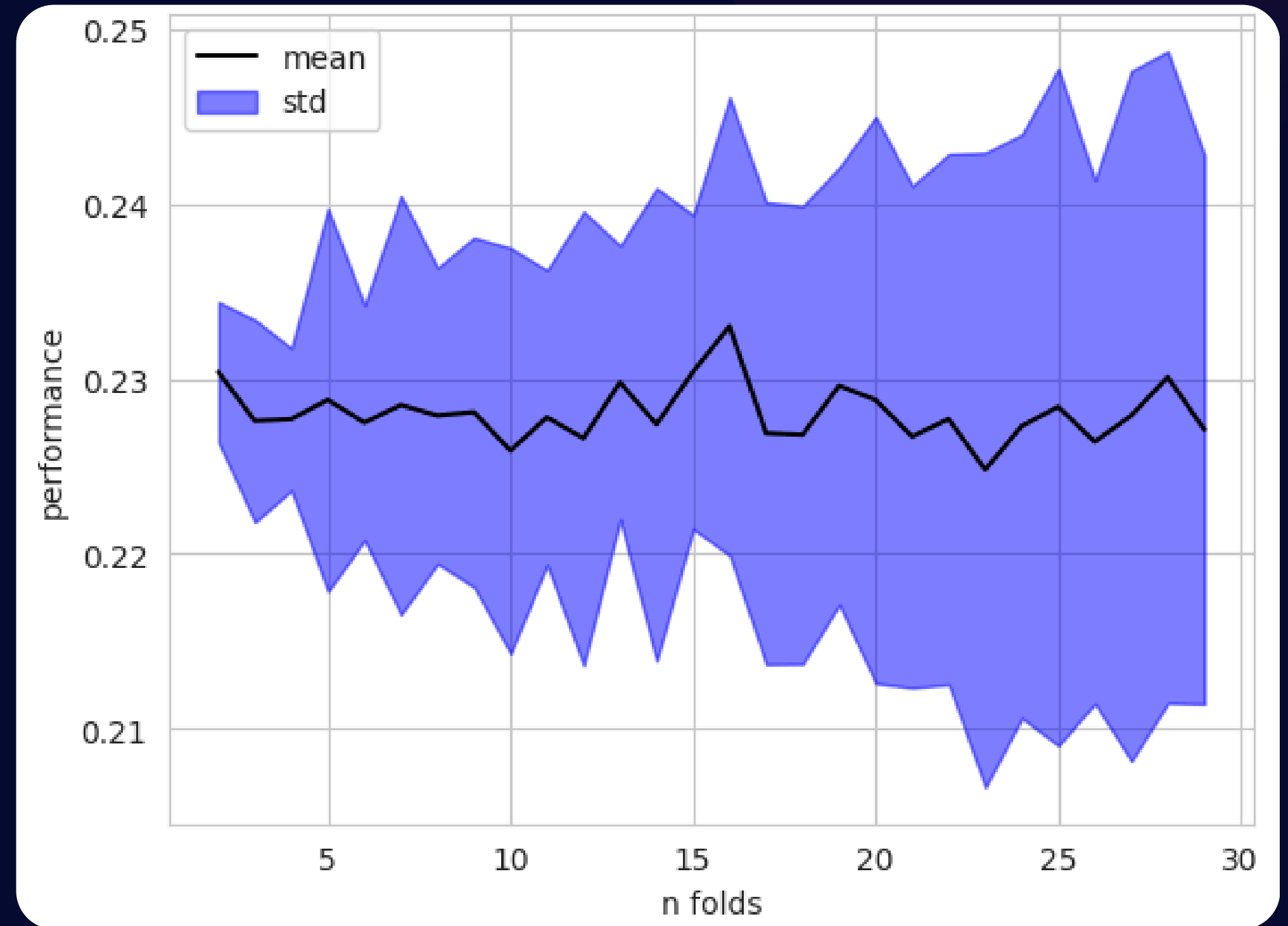
RANDOM FOREST

N-folds = 10

$\sigma = 0.011579723658188051$

+

+



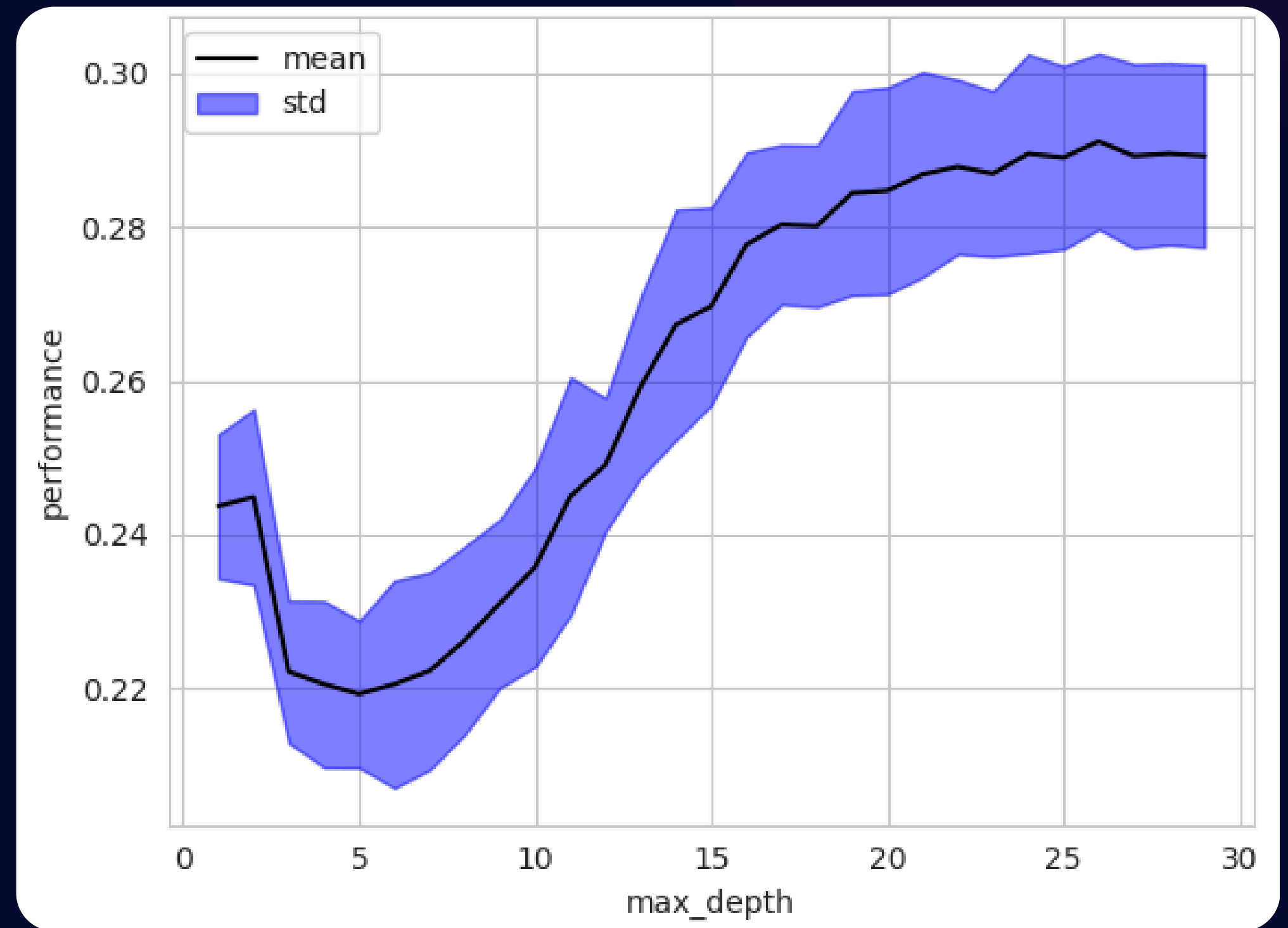
Max depth



DESICION TREE

Max depth = 5

$\sigma = 0.009539916142189095$



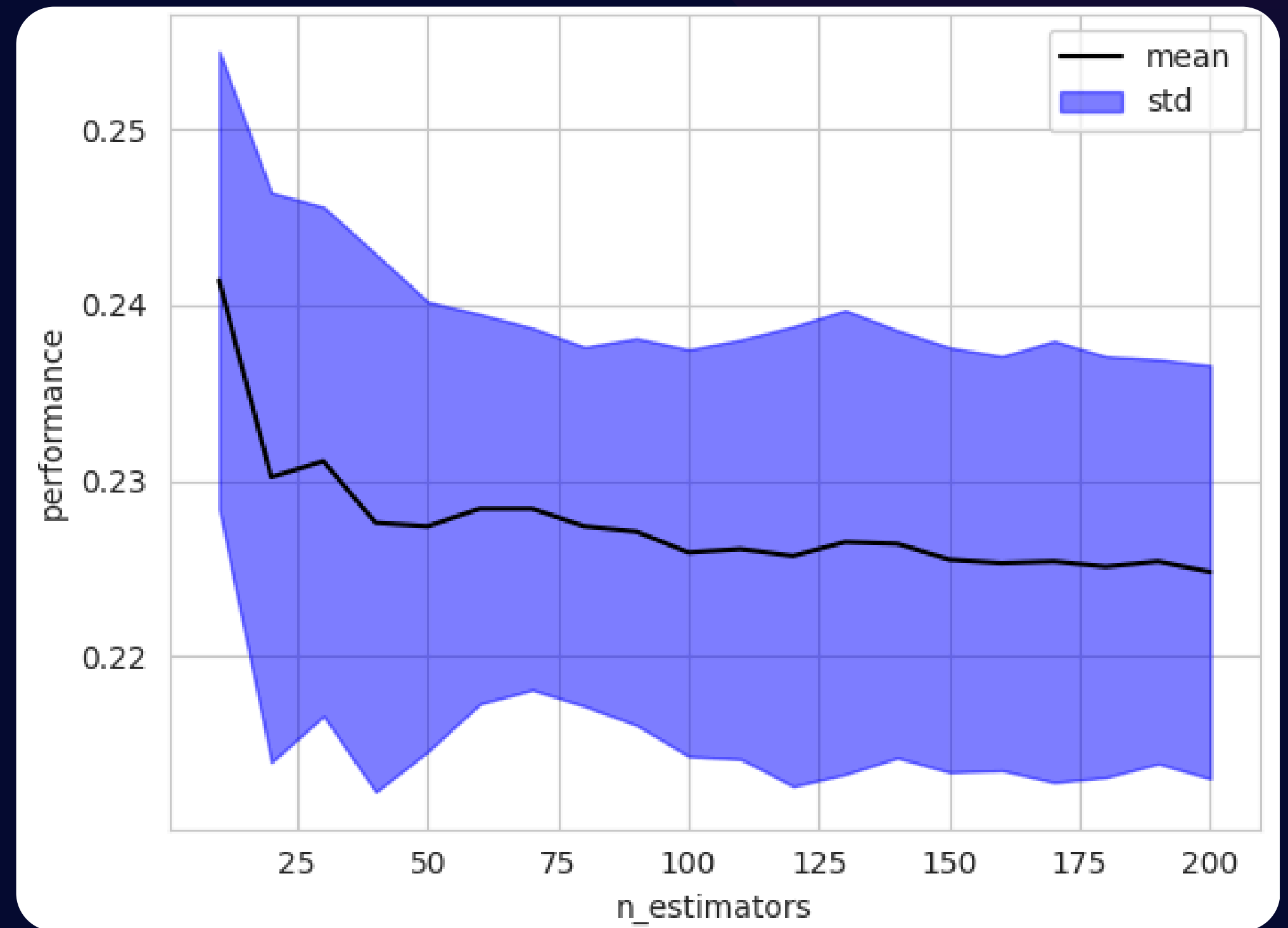
N-Estimators



RANDOM FOREST

N-Estimators = 170

$\sigma = 0.011762652762026088$



Kernel

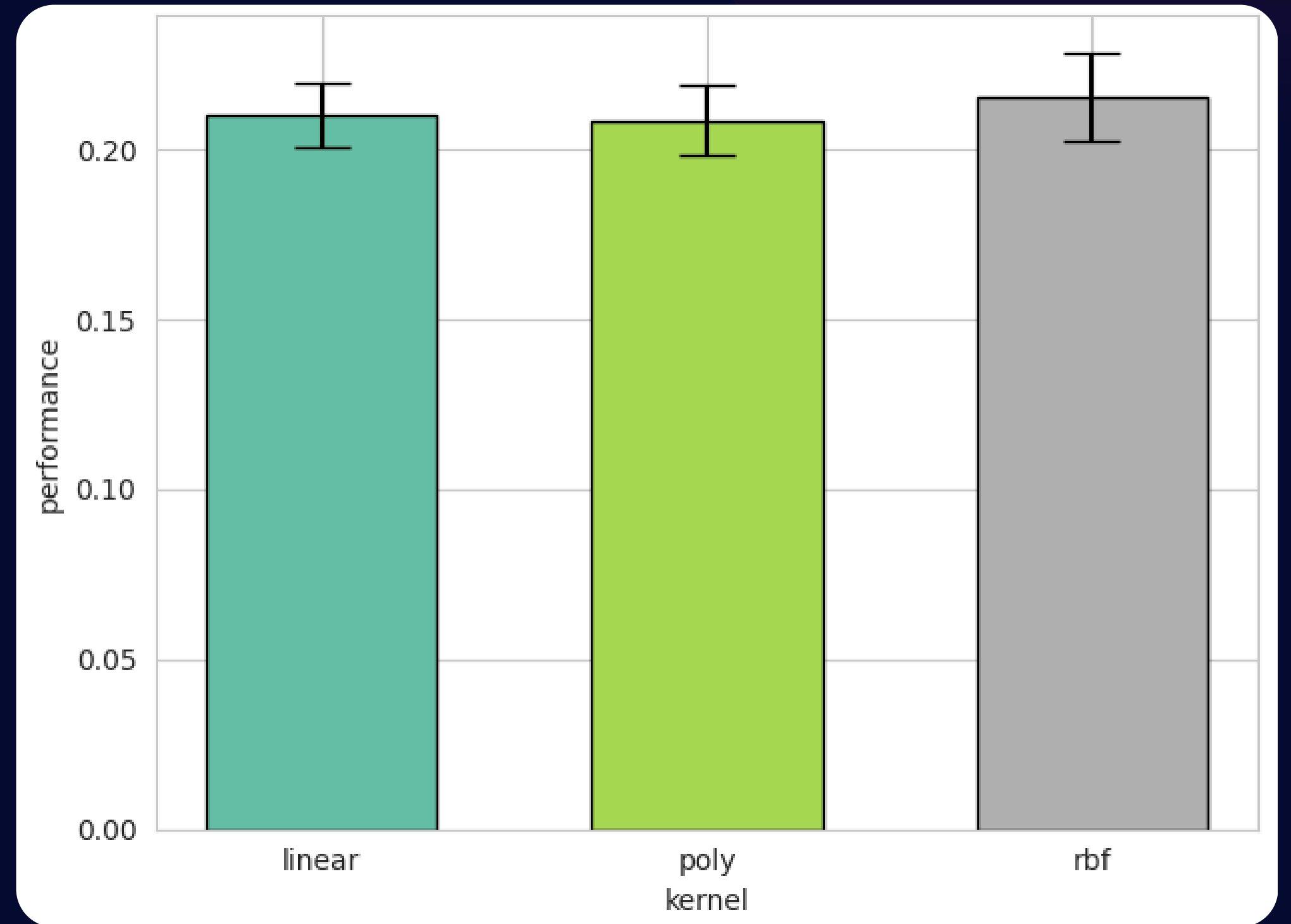


SUPPORT VECTOR MACHINE

Kernel = poly

+

+



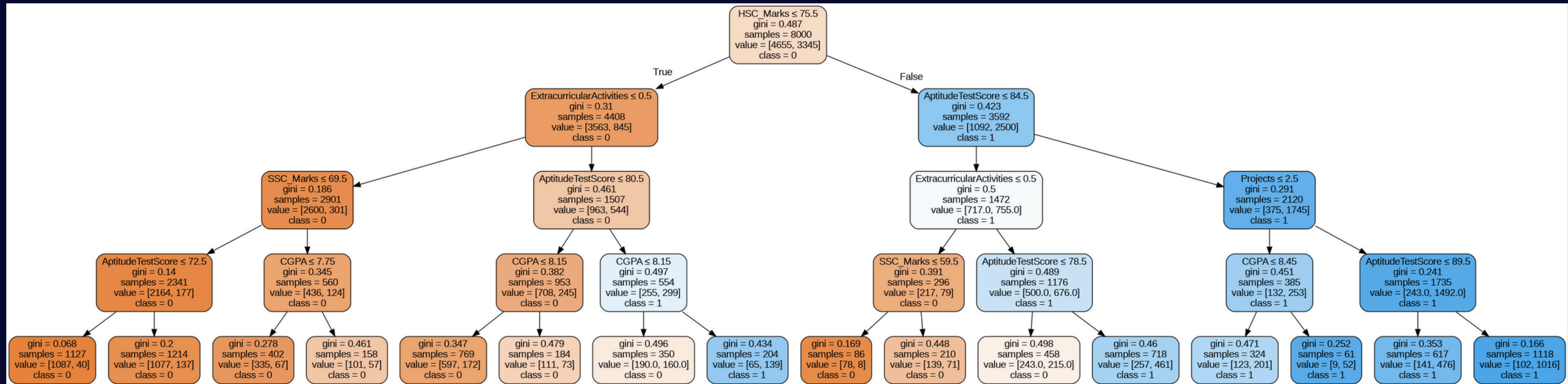


BEST DECISION TREE

```
1 best_decision_tree = DecisionTreeClassifier(max_depth=5, random_state=42)
2 best_decision_tree.fit(X_train, y_train)
3 predictions_tree = best_decision_tree.predict(X_test)
4 accuracy_tree = accuracy_score(y_test, predictions_tree)
```

Accuracy Decision Tree: 0.785

TREE DECISIONS FROM $+$ $=$



+

=

BEST RANDOM FOREST

```
1 best_rf = RandomForestClassifier(n_estimators=170, random_state=42)
2 best_rf.fit(X_train, y_train)
3 predictions_rf = best_rf.predict(X_test)
4 accuracy_rf = accuracy
5 print("Accuracy Random Forest:", accuracy_rf)
```

Accuracy Decision Tree: 0.7805

+

=

BEST SUPPORT VECTOR MACHINE

```
1 best_svc = SVC(kernel='poly', random_state=42)
2 best_svc.fit(x_train, y_train)
3 predictions_svc = best_svc.predict(X_test)
4 accuracy_svc = accuracy_score(y_test, predictions_svc)
5 print("Accuracy Support Vector Machine: ", accuracy_svc)
```

Accuracy Support Vector Machine: 0.781

NEURAL NETWORKS

+

=

3 CAPAS OCULTAS, 128 NEURONAS, RELU

```
1 model = tf.keras.Sequential([
2     tf.keras.layers.Flatten(input_shape=X_scaled[0].shape),
3     tf.keras.layers.Dense(128, activation='relu'),
4     tf.keras.layers.Dense(128, activation='relu'),
5     tf.keras.layers.Dense(128, activation='relu'),
6     tf.keras.layers.Dense(2, activation='sigmoid') #Clasificacion binaria
7 ])
8
9 model.compile(optimizer='adam',
10               loss='sparse_categorical_crossentropy',
11               metrics=['accuracy', ])
12
13 model.fit(X_scaled, y, epochs=10)
```

Test accuracy: 0.8004999756813049

+

=

6 CAPAS OCULTAS, 128 NEURONAS, RELU

```
1 model = tf.keras.Sequential([
2     tf.keras.layers.Flatten(input_shape=X_scaled[0].shape),
3     tf.keras.layers.Dense(128, activation='relu'),
4     tf.keras.layers.Dense(128, activation='relu'),
5     tf.keras.layers.Dense(128, activation='relu'),
6     tf.keras.layers.Dense(128, activation='relu'),
7     tf.keras.layers.Dense(128, activation='relu'),
8     tf.keras.layers.Dense(128, activation='relu'),
9     tf.keras.layers.Dense(2, activation='sigmoid') #Clasificacion binaria
10 ])
11
12 model.compile(optimizer='adam',
13               loss='sparse_categorical_crossentropy',
14               metrics=['accuracy', ])
15
16 model.fit(X_scaled, y, epochs=10)
```

Test accuracy: 0.7975000143051147

+

=

10 CAPAS OCULTAS, 128 NEURONAS, RELU

```
1 model = tf.keras.Sequential([
2     tf.keras.layers.Flatten(input_shape=X_scaled[0].shape),
3     tf.keras.layers.Dense(128, activation='relu'),
4     tf.keras.layers.Dense(128, activation='relu'),
5     tf.keras.layers.Dense(128, activation='relu'),
6     tf.keras.layers.Dense(128, activation='relu'),
7     tf.keras.layers.Dense(128, activation='relu'),
8     tf.keras.layers.Dense(128, activation='relu'),
9     tf.keras.layers.Dense(128, activation='relu'),
10    tf.keras.layers.Dense(128, activation='relu'),
11    tf.keras.layers.Dense(128, activation='relu'),
12    tf.keras.layers.Dense(128, activation='relu'),
13    tf.keras.layers.Dense(2, activation='sigmoid') #Clasificacion binaria
14 ])
15
16 model.compile(optimizer='adam',
17               loss='sparse_categorical_crossentropy',
18               metrics=['accuracy', ])
19
20 model.fit(X_scaled, y, epochs=10)
```

Test accuracy: 0.7914999723434448

N-Components



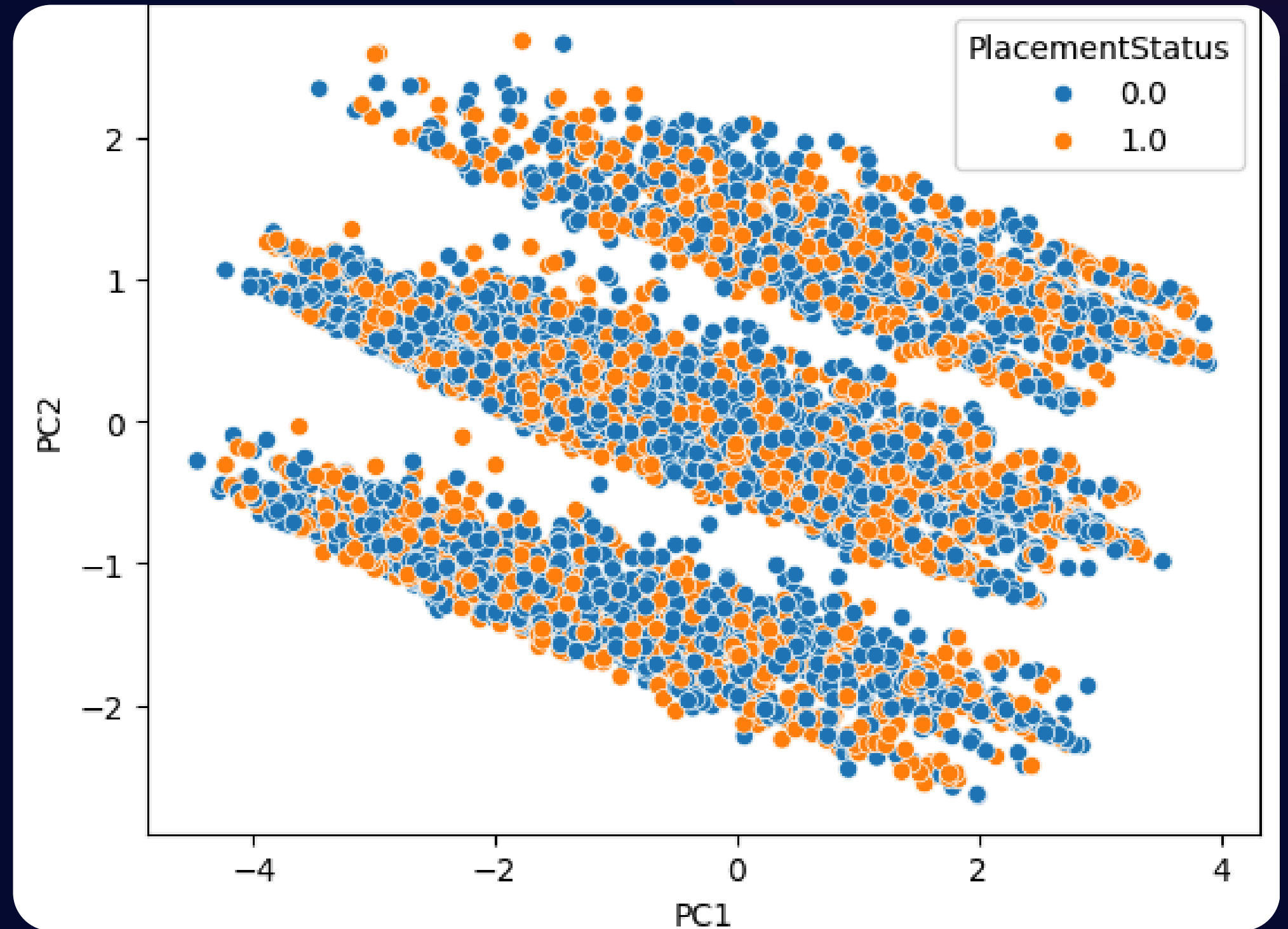
+

PRINCIPAL COMPONENT ANALYSIS

N-Components = 2

La "energía" que conservamos con
dos componentes es de 57%

+

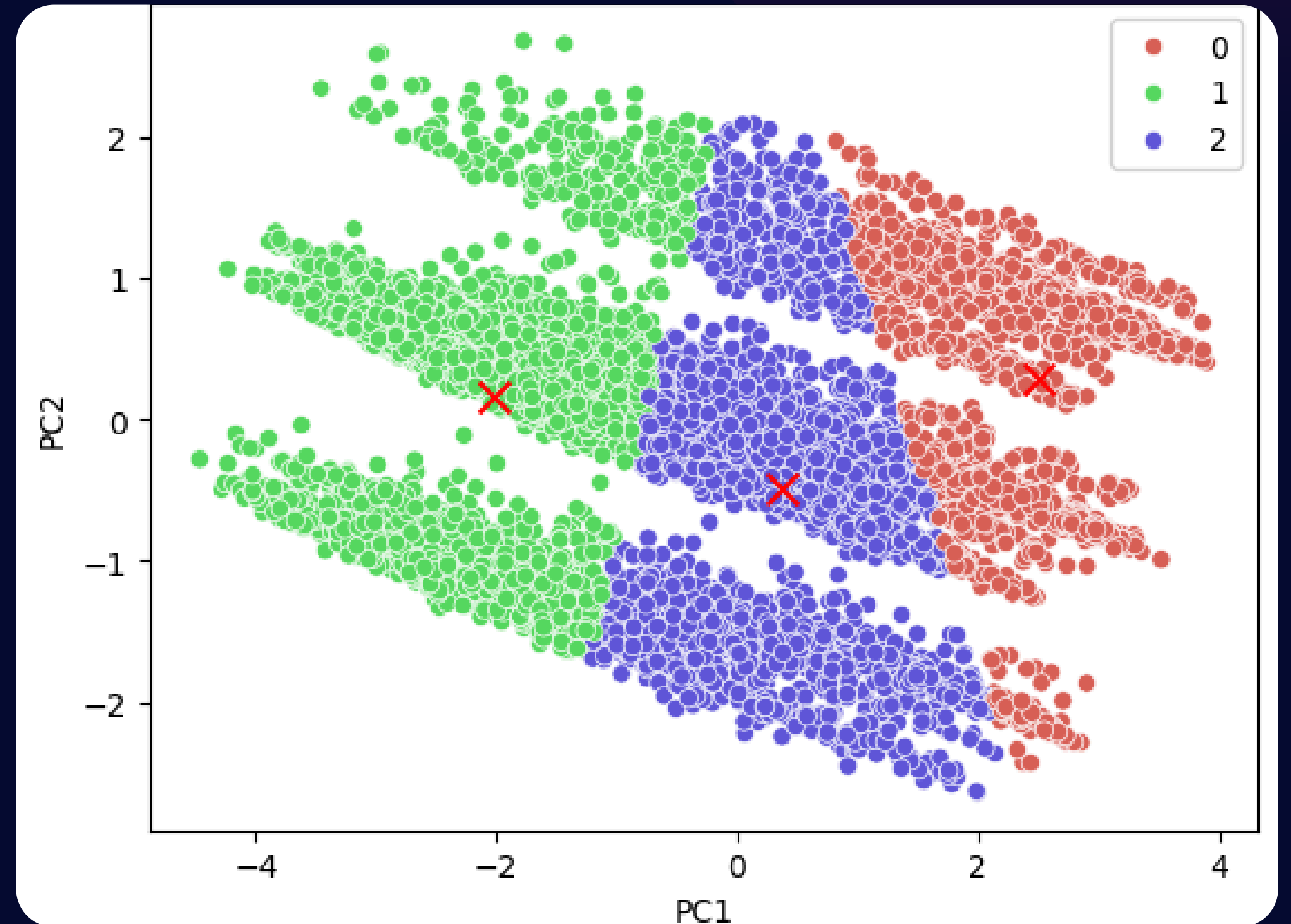


N-Clusters



K-MEANS

N-Clusters = 3

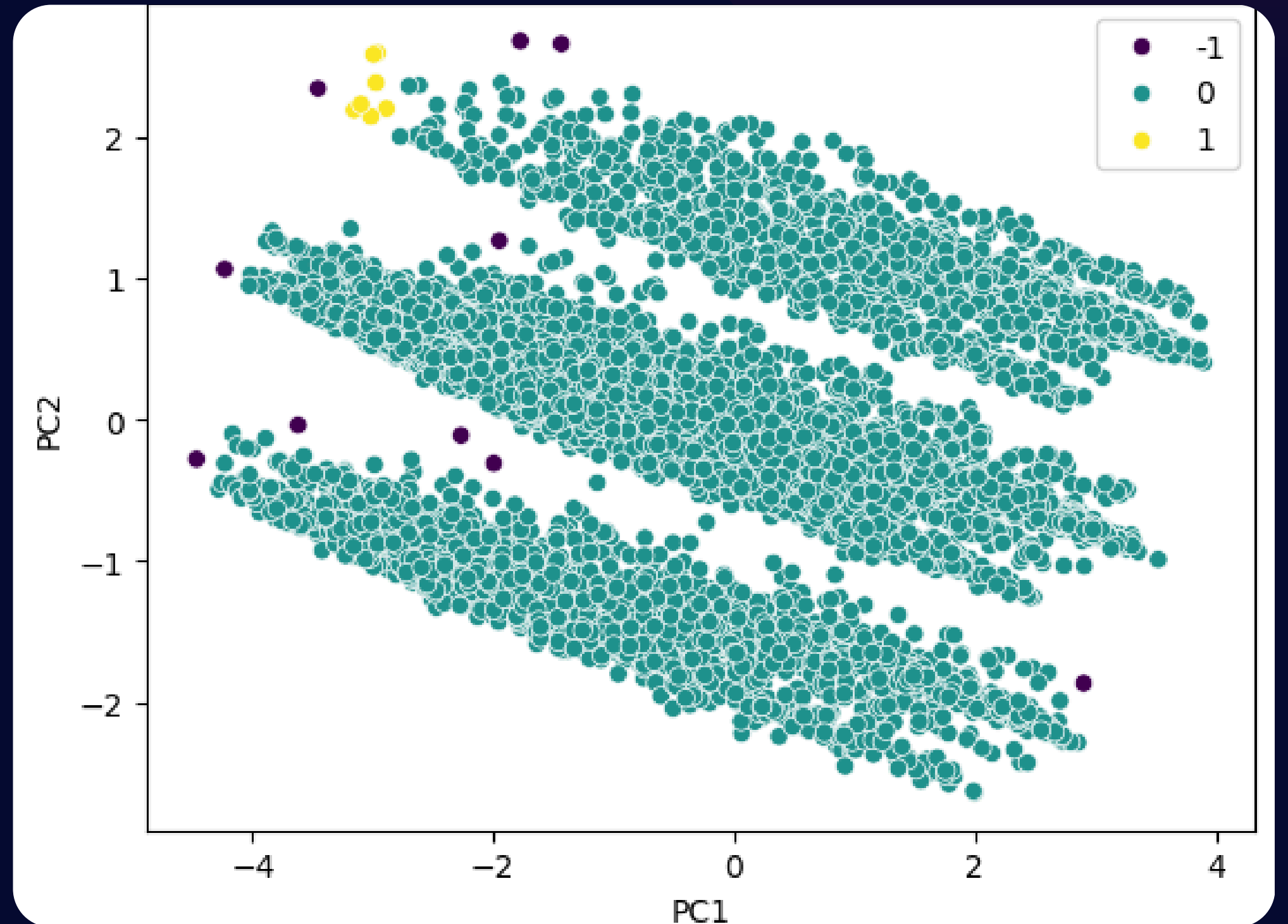


Min-Samples



DBSCAN

Min-Samples = 2





GRACIAS!

