



UNIVERSIDAD AUTÓNOMA DE YUCATÁN

FACULTAD DE MATEMÁTICAS
LICENCIADO EN CIENCIAS DE LA COMPUTACIÓN
INTELIGENCIA ARTIFICIAL

PROYECTO APRENDIZAJE SUPERVISADO

PROFESOR / ALUMNOS:

ALEJANDRO PASOS RUIZ / JULIAN

EMILIANO ORTIZ RIVERO / ANGEL MARIEL

OSALDE SALAZAR & MIGUEL OMAR SOLIS

RODRIGUEZ

REALIZAR EXPERIMENTOS USANDO VARIOS ALGORITMOS DE
MACHINE LEARNING PARA LOS DATOS Y REALIZAR AJUSTE
DE PARAMETROS Y VALIDACION CRUZZADA PARA DAR
SUSTENTO A LOS RESULTADOS

MÉRIDA, YUCATÁN, ENERO - MAYO 2025

Contents

Contents	1
1 Introducción	2
2 Descripción del Problema	2
3 Descripción de los Datos	2
3.1 Limpieza del Dataset	2
4 Experimentos	3
4.1 Modelos Utilizados	3
4.2 Carga de datos de validación	6
4.3 Evaluación con Validación cruzada, Precisión y Recall	6
5 Resultados	7
5.1 Gráficas de Resultados	7
5.2 Evaluación con Validación Cruzada, Precisión y Recall	7
5.3 Evaluación con Validación Cruzada promediada	8
5.4 Comparación de los modelos en Precisión	9
5.5 Comparación de los modelos en Recall	9
5.6 Comparaciones	10
6 Conclusiones	11

1 Introducción

La enfermedad de Parkinson es un trastorno cerebral que deteriora las células nerviosas encargadas del control de los movimientos y aunque no hay medicamentos que frenen su avance, existen tratamientos para controlar los síntomas en cada etapa, por lo que es importante un diagnóstico preciso y un manejo adecuado de los medicamentos para mejorar la calidad de vida de los pacientes. El Parkinson, al ser el segundo trastorno cerebral más común después del Alzheimer y tener como principal factor de riesgo la edad, representa un problema para la salud pública, especialmente en sociedades occidentales con poblaciones con un porcentaje considerable de personas mayores.

Es por ello que el reporte tiene como objetivo analizar un conjunto de datos relacionados con la enfermedad de Parkinson para predecir el estado de salud de los pacientes a partir de diversas características, utilizando modelos de aprendizaje automático que permiten identificar patrones en los datos y hacer estimaciones sobre su condición. El documento está organizado en varias secciones, donde primero se explica el problema que se aborda, luego se describe el conjunto de datos con el que se trabaja, se detallan los experimentos realizados, se presentan los resultados obtenidos y se escriben algunas conclusiones sobre lo que se ha aprendido.

2 Descripción del Problema

La enfermedad de Parkinson es un trastorno complejo que afecta el movimiento, causando temblores, rigidez y dificultades para caminar, lo que presenta un problema para los médicos, ya que los síntomas pueden variar mucho entre las personas y son similares a los de otras enfermedades. Esto hace que el diagnóstico temprano sea difícil, lo que complica el tratamiento y el manejo de la enfermedad. Además, como los síntomas suelen empeorar con el tiempo, es importante poder detectar la enfermedad en sus primeras etapas para intervenir lo más pronto posible. Por esta razón, se pretende utilizar técnicas de aprendizaje automático para analizar datos relacionados con la enfermedad y predecir el estado de salud de los pacientes, lo que podría mejorar el diagnóstico y permitir un tratamiento más efectivo desde las primeras etapas de la enfermedad.

3 Descripción de los Datos

El conjunto de datos utilizado en este estudio proviene de un dataset en el que se incluyen varias características relacionadas con la enfermedad de Parkinson, como medidas de voz, temblores, y otras variables clínicas. El conjunto de datos se limpió y preprocesó para eliminar valores nulos y manejar valores atípicos. Además, se normalizaron las características para mejorar el rendimiento de los modelos de aprendizaje automático.

3.1 Limpieza del Dataset

A continuación, se muestra el código utilizado para limpiar y preprocesar el dataset:

```
1 # Importar librerías necesarias
2 import pandas as pd
3 import numpy as np
4 from sklearn.preprocessing import StandardScaler
5 import joblib
6
7 # Cargar el dataset
8 df = pd.read_csv("updated_dataset.csv")
```

```

9
10 # Eliminar la columna 'name' ya que no aporta informaci n relevante
11 df.drop(columns=['name'], inplace=True)
12
13 # Verificar valores nulos y eliminar filas si es necesario
14 df.dropna(inplace=True)
15
16 # Identificar y manejar outliers con percentiles
17 for col in df.columns:
18     if col != "status": # No modificar la variable objetivo
19         q1 = df[col].quantile(0.25)
20         q3 = df[col].quantile(0.75)
21         iqr = q3 - q1
22         lower_bound = q1 - 1.5 * iqr
23         upper_bound = q3 + 1.5 * iqr
24         df = df[(df[col] >= lower_bound) & (df[col] <= upper_bound)]
25
26 # Separar caracter sticas y variable objetivo
27 X = df.drop(columns=["status"])
28 y = df["status"]
29
30 # Normalizaci n de las caracter sticas
31 scaler = StandardScaler()
32 X_scaled = scaler.fit_transform(X)
33
34 # Guardar datos preprocesados y scaler
35 joblib.dump((X_scaled, y), "datos_limpiados.pkl")
36 joblib.dump(scaler, "scaler.pkl")
37
38 print("Preprocesamiento completo. Datos limpios guardados en 'datos_limpiados.pkl'")

```

4 Experimentos

Se realizaron varios experimentos utilizando diferentes modelos de aprendizaje automtico, incluyendo Regresin Logstica, Random Forest, MLP (Perceptrn Multicapa) y XGBoost. Los datos se dividieron en conjuntos de entrenamiento y prueba, y se utiliz validacin cruzada para evaluar el rendimiento de los modelos. Adems, se calcularon mtricas como precisin, recall y accuracy para comparar el rendimiento de los modelos.

4.1 Modelos Utilizados

Regresin Logstica: Un modelo lineal utilizado para problemas de clasificacin binaria, que se basa en la relacin lineal entre las caractersticas de entrada y la salida. Este modelo es adecuado para problemas en los que la relacin entre las variables es lineal o casi lineal.

```

1 # Cargar los datos preprocesados
2 X_scaled, y = joblib.load("datos_limpiados.pkl")
3
4 # Divisi n en datos de entrenamiento y prueba
5 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2,
6     random_state=42, stratify=y)
7
8 # Entrenar el modelo de Regresi n Log stica
9 log_reg = LogisticRegression(max_iter=500, random_state=42)
10 log_reg.fit(X_train, y_train)
11
12 # Predicci n y evaluaci n
13 y_pred = log_reg.predict(X_test)
14 accuracy = accuracy_score(y_test, y_pred)

```

```

14 report = classification_report(y_test, y_pred)
15
16 print("Accuracy de Regresión Logística:", accuracy)
17 print("Reporte de clasificación:\n", report)
18
19 # Guardar el modelo
20 joblib.dump(log_reg, "logistic_regression_model.pkl")

```

Accuracy de Regresión Logística: 0.5089820359281437

Reporte de clasificación:

	precision	recall	f1-score	support
0	0.41	0.24	0.31	74
1	0.54	0.72	0.62	93
accuracy			0.51	167
macro avg	0.48	0.48	0.46	167
weighted avg	0.48	0.51	0.48	167

['logistic_regression_model.pkl']

Random Forest: Un método de ensamble que utiliza múltiples árboles de decisión para mejorar la precisión y reducir la variabilidad de los resultados. Cada árbol de decisión se entrena en un subconjunto de los datos y se combina con los resultados de los demás árboles para obtener un resultado final.

```

1 # Cargar los datos preprocesados
2 X_scaled, y = joblib.load("datos_limpiados.pkl")
3
4 # División en datos de entrenamiento y prueba
5 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2,
6     random_state=42, stratify=y)
7
8 # Entrenar el modelo Random Forest
9 rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
10 rf_model.fit(X_train, y_train)
11
12 # Predicción y evaluación
13 y_pred = rf_model.predict(X_test)
14 accuracy = accuracy_score(y_test, y_pred)
15 report = classification_report(y_test, y_pred)
16
17 print("Accuracy de Random Forest:", accuracy)
18 print("Reporte de clasificación:\n", report)
19
20 # Guardar el modelo
21 joblib.dump(rf_model, "random_forest_model.pkl")

```

Accuracy de Random Forest: 0.5808383233532934

Reporte de clasificación:

	precision	recall	f1-score	support
0	0.53	0.43	0.48	74
1	0.61	0.70	0.65	93
accuracy			0.58	167
macro avg	0.57	0.57	0.56	167
weighted avg	0.57	0.58	0.57	167

```
['random_forest_model.pkl']
```

MLP (Perceptrón Multicapa): Una red neuronal con múltiples capas ocultas, que permite aprender representaciones complejas de los datos. Cada capa de la red neuronal procesa la información de la capa anterior y la envía a la capa siguiente, lo que permite aprender patrones y relaciones complejas en los datos.

```
1 # Cargar los datos preprocesados
2 X_scaled, y = joblib.load("datos_limpiados.pkl")
3
4 # Divisi n en datos de entrenamiento y prueba
5 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2,
6 random_state=42, stratify=y)
7
8 # Entrenar el modelo de Red Neuronal
9 mlp_model = MLPClassifier(hidden_layer_sizes=(64, 32), max_iter=500, random_state=42)
10 mlp_model.fit(X_train, y_train)
11
12 # Predicci n y evaluaci n
13 y_pred = mlp_model.predict(X_test)
14 accuracy = accuracy_score(y_test, y_pred)
15 report = classification_report(y_test, y_pred)
16
17 print("Accuracy de Red Neuronal:", accuracy)
18 print("Reporte de clasificaci n:\n", report)
19
20 # Guardar el modelo
21 joblib.dump(mlp_model, "mlp_model.pkl")
```

Accuracy de Red Neuronal: 0.5389221556886228

Reporte de clasificaci3n:

	precision	recall	f1-score	support
0	0.48	0.47	0.48	74
1	0.59	0.59	0.59	93
accuracy			0.54	167
macro avg	0.53	0.53	0.53	167
weighted avg	0.54	0.54	0.54	167

```
['mlp_model.pkl']
```

XGBoost: Un algoritmo de boosting que utiliza 6rboles de decisi3n para mejorar la precisi3n y reducir la variabilidad de los resultados. XGBoost es un algoritmo de aprendizaje autom6tico que se basa en la t6cnica de boosting, que consiste en combinar los resultados de m6ltiples modelos para obtener un resultado final m6s preciso.

```
1 # Cargar los datos preprocesados
2 X_scaled, y = joblib.load("datos_limpiados.pkl")
3
4 # Divisi n en datos de entrenamiento y prueba
5 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2,
6 random_state=42, stratify=y)
7
8 # Entrenar el modelo XGBoost
9 xgb_model = XGBClassifier(eval_metric='mlogloss', random_state=42)
10 xgb_model.fit(X_train, y_train)
11
12 # Predicci n y evaluaci n
```

```

12 y_pred = xgb_model.predict(X_test)
13 accuracy = accuracy_score(y_test, y_pred)
14 report = classification_report(y_test, y_pred)
15
16 print("Accuracy de XGBoost:", accuracy)
17 print("Reporte de clasificaci n:\n", report)
18
19 # Guardar el modelo
20 joblib.dump(xgb_model, "xgboost_model.pkl")

```

Accuracy de XGBoost: 0.6407185628742516

Reporte de clasificaci3n:

	precision	recall	f1-score	support
0	0.60	0.58	0.59	74
1	0.67	0.69	0.68	93
accuracy			0.64	167
macro avg	0.64	0.63	0.63	167
weighted avg	0.64	0.64	0.64	167

['xgboost_model.pkl']

4.2 Carga de datos de validaci3n

```

1 # Cargar los datos preprocesados
2 X_scaled, y = joblib.load("datos_limpiados.pkl")
3
4 # Cargar modelos
5 models = {
6     "Logistic Regression": joblib.load("logistic_regression_model.pkl"),
7     "Random Forest": joblib.load("random_forest_model.pkl"),
8     "MLP Neural Network": joblib.load("mlp_model.pkl"),
9     "XGBoost": joblib.load("xgboost_model.pkl")
10 }

```

Accuracy de Random Forest: 0.5808383233532934

Reporte de clasificaci3n:

	precision	recall	f1-score	support
0	0.53	0.43	0.48	74
1	0.61	0.70	0.65	93
accuracy			0.58	167
macro avg	0.57	0.57	0.56	167
weighted avg	0.57	0.58	0.57	167

['random_forest_model.pkl']

4.3 Evaluaci3n con Validaci3n cruzada, Precisi3n y Recall

```

1 from sklearn.model_selection import cross_val_score, cross_validate
2 from sklearn.metrics import make_scorer, precision_score, recall_score
3

```

```

4 # Diccionarios para almacenar resultados
5 cv_results = {}
6 precision_results = {}
7 recall_results = {}
8
9 # Evaluación con validación cruzada (5-fold) para cada métrica
10 for name, model in models.items():
11     # Calcular Accuracy
12     accuracy_scores = cross_val_score(model, X_scaled, y, cv=5, scoring='accuracy')
13     cv_results[name] = accuracy_scores # Se guarda directamente la lista de valores
14
15     # Calcular Precision (Weighted)
16     precision_scores = cross_val_score(
17         model, X_scaled, y, cv=5,
18         scoring=make_scorer(precision_score, average='weighted')
19     )
20     precision_results[name] = precision_scores
21
22     # Calcular Recall (Weighted)
23     recall_scores = cross_val_score(
24         model, X_scaled, y, cv=5,
25         scoring=make_scorer(recall_score, average='weighted')
26     )
27     recall_results[name] = recall_scores

```

5 Resultados

5.1 Gráficas de Resultados

A continuación, se muestran las gráficas de los resultados de la validación cruzada y las métricas de precisión y recall.

5.2 Evaluación con Validación Cruzada, Precisión y Recall

```

1 plt.figure(figsize=(10, 5))
2 sns.boxplot(data=cv_results.values())
3 plt.ylabel("Accuracy")
4 plt.title("Comparación de modelos con validación cruzada (5-fold)")
5 plt.xticks(ticks=range(len(models)), labels=models.keys(), rotation=15)
6 plt.show()

```

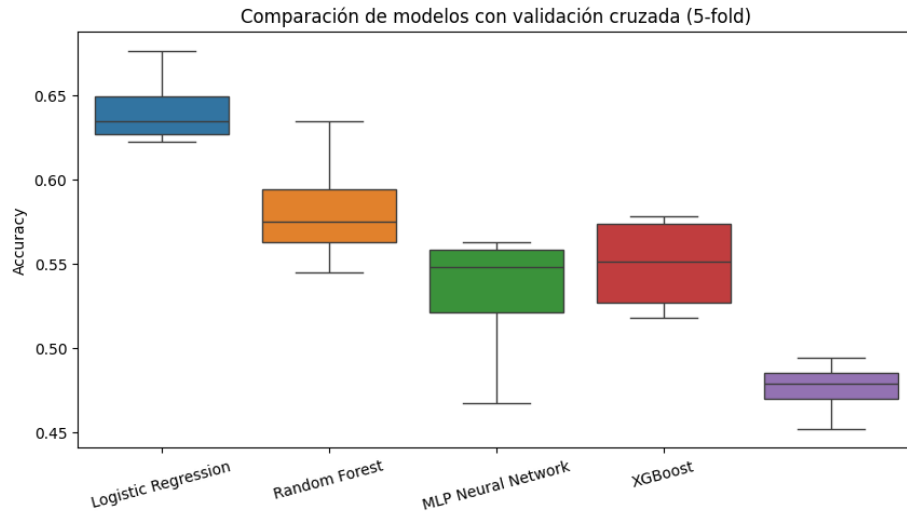



Figure 1: Comparación de modelos con validación cruzada (5-fold)

5.3 Evaluación con Validación Cruzada promediada

```

1 # Graficar resultados de validación cruzada
2 plt.figure(figsize=(10, 5))
3 model_names = list(cv_results.keys()) # Los nombres de los modelos (eje X)
4 accuracy_values = [np.mean(cv_results[name]) for name in model_names] # Promedio de
5 # accuracy por modelo
6 bars = plt.bar(model_names, accuracy_values, color=['blue', 'green', 'red', 'purple'])
7 plt.ylabel("Accuracy")
8 plt.ylim(0,1)
9 plt.title("Comparación de modelos con validación cruzada (5-fold) promediados")
10 # Agregar los valores de accuracy sobre las barras
11 for bar in bars:
12     yval = bar.get_height() # Obtener la altura de cada barra (valor de accuracy)
13     plt.text(bar.get_x() + bar.get_width() / 2, yval + 0.02, round(yval, 2), ha='
14         center', va='bottom')
15 plt.show()

```

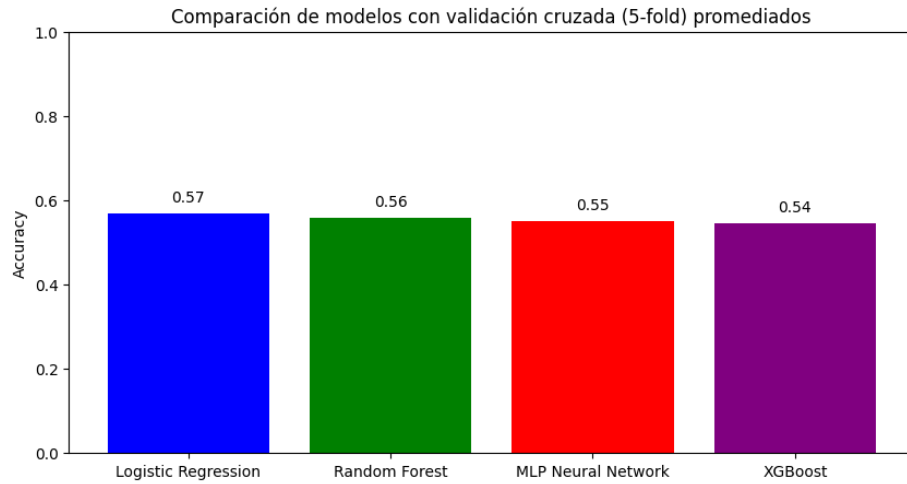


Figure 2: Comparación de modelos con validación cruzada (5-fold) promediados

5.4 Comparación de los modelos en Precisión

```
1 plt.figure(figsize=(10, 5))
2 sns.boxplot(data=precision_results.values())
3 plt.ylabel("Precision (Weighted)")
4 plt.title("Comparación de modelos - Precisión")
5 plt.xticks(ticks=range(len(models)), labels=models.keys(), rotation=15)
6 plt.show()
```

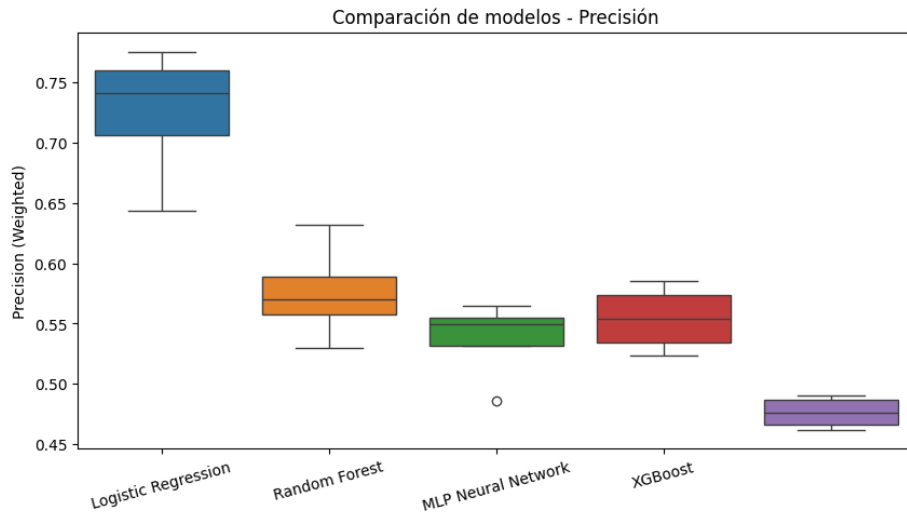


Figure 3: Comparación de modelos con precisión

5.5 Comparación de los modelos en Recall

```
1 plt.figure(figsize=(10, 5))
2 sns.boxplot(data=recall_results.values())
3 plt.ylabel("Recall (Weighted)")
4 plt.title("Comparación de modelos - Recall")
```

```

5 plt.xticks(range(len(models)), models.keys(), rotation=15)
6 plt.show()

```

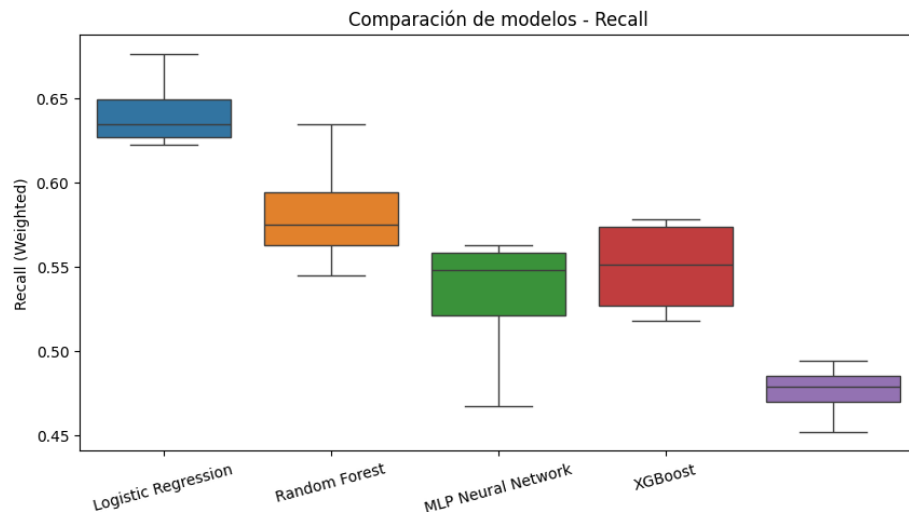


Figure 4: Comparación de modelos con recall

5.6 Comparaciones

En la siguiente tabla se presentan los resultados de las pruebas estadísticas entre los diferentes modelos.

```

1 # Pruebas estadísticas
2 stat_tests = {}
3 model_names = list(models.keys())
4 for i in range(len(model_names)):
5     for j in range(i + 1, len(model_names)):
6         model1, model2 = model_names[i], model_names[j]
7         t_stat, p_ttest = ttest_rel(cv_results[model1], cv_results[model2])
8         w_stat, p_wilcoxon = wilcoxon(cv_results[model1], cv_results[model2])
9         stat_tests[f"{model1} vs {model2}"] = (p_ttest, p_wilcoxon)
10
11 # Mostrar pruebas estadísticas
12 print("\nResultados de pruebas estadísticas (p-values):")
13 for comparison, (p_ttest, p_wilcoxon) in stat_tests.items():
14     print(f"{comparison}: t-test p={p_ttest:.5f}, Wilcoxon p={p_wilcoxon:.5f}")

```

Modelos Comparados	t-test	Wilcoxon
Logistic Regression vs Random Forest	0.70218	0.62500
Logistic Regression vs MLP Neural Network	0.31939	0.43750
Logistic Regression vs XGBoost	0.20003	0.18750
Random Forest vs MLP Neural Network	0.84336	0.81250
Random Forest vs XGBoost	0.58293	0.62500
MLP Neural Network vs XGBoost	0.45835	0.62500

6 Conclusiones

En este estudio, se aplicaron varios modelos de aprendizaje automático para predecir el estado de salud de pacientes con Parkinson. Los resultados indican que el modelo XGBoost es el más efectivo para esta tarea, alcanzando un accuracy del 64%. Sin embargo, es importante destacar que todos los modelos tuvieron un rendimiento relativamente bajo, lo que sugiere que podría ser necesario explorar otras técnicas o mejorar el preprocesamiento de los datos para obtener mejores resultados.

En futuros trabajos, se podría considerar la inclusión de más características o la aplicación de técnicas de selección de características para mejorar el rendimiento de los modelos. Además, sería interesante explorar el uso de técnicas de aprendizaje profundo para ver si se pueden obtener mejores resultados.