

UNIVERSIDAD MAYOR DE SAN ANDRÉS
FACULTAD DE CIENCIAS PURAS Y NATURALES
CARRERA DE INFORMÁTICA



Tesis de Grado

Mencion Ciencia de Datos e Inteligencia Artificial

Postulante: Univ. Enrique Julio Alvarado Mamani

Tutor: Lic. Moises Silva

Relator: Dat - 245

Tesis para optar al grado Académico de Licenciado en Economía

La Paz-Bolivia
2025

RESUMEN

Esta sección introduce el contexto histórico y matemático de la Hipótesis de Riemann. Comienza explicando la fascinación milenaria con la distribución de los números primos y el Teorema de los Números Primos como un primer intento de describir su frecuencia. El punto central es la contribución revolucionaria de Bernhard Riemann en 1859, quien conectó la distribución de los primos con los ceros de una función de variable compleja: la función zeta de Riemann ($\zeta(s)$). La Hipótesis de Riemann es la conjetura de que todos los "ceros no triviales" de esta función se sitúan sobre una línea específica en el plano complejo, la "línea crítica" donde la parte real es $1/2$. Se subraya que la veracidad de esta hipótesis proporcionaría la descripción más precisa posible sobre cómo se distribuyen los números primos, convirtiéndola en un pilar fundamental de la teoría de números. Finalmente, se establece que, ante la ausencia de una prueba formal, la computación ha emergido como una herramienta esencial para acumular evidencia empírica a su favor.

INDICE GENERAL

BIBLIOGRAFÍA	1
1 Descripción del Dataset y Objetivo de Investigación	v
1.1 Descripción del Dataset	v
1.2 Objetivo de Investigación	viii
2 Proceso Básico de Análisis de Datos	x
2.1 a) Preprocesamiento de Datos	x
2.2 b) Selección y Justificación del Clasificador	xi
2.3 c) Primera Ejecución del Modelo	xi
2.4 d) Validación por Asignaciones (Splits)	xi
2.5 e) Código Fuente	xi
3 Reducción de Dimensionalidad	xii
3.1 Relaciones con la Matriz de Correlación	xii
3.2 Combinación de Columnas	xiv
3.2.1 Fundamento Matemático para la Combinación	xiv
3.2.2 Fundamento Matemático para LDA	xv
4 Aprendizaje y Predicción	xvii
4.1 Entrenamiento del Modelo	xvii
4.2 Etapas de Entrenamiento y Prueba	xviii
4.3 Prueba con Nuevos Datos	xix
5 Conclusiones y Recomendaciones	xxi
5.1 Conclusiones	xxi
5.2 Recomendaciones	xxi

List of Figures

1	Visualización de la eliminación de columnas identificadoras.	x
2	Distribución de valores antes y después de la imputación de medias.	xi
3	Matriz de correlación de Pearson para las variables numéricas.	xii
4	Matriz de correlación de Spearman para las variables numéricas.	xiii
5	Matriz de información mutua para las variables numéricas.	xiv
6	Gráfico de las etapas de entrenamiento y prueba, mostrando la pérdida y precisión a lo largo de las iteraciones.	xviii

Este artículo ofrece un análisis detallado y extenso del dataset de exoplanetas recolectado por la misión Kepler de la NASA, que comprende 9,564 registros con 43 variables, incluyendo 39 numéricas y 4 categóricas. Se emplea un perceptrón multicapa (MLP) como clasificador supervisado, optimizado mediante Análisis Discriminante Lineal (LDA) para reducir la dimensionalidad a 4 componentes, con el objetivo de clasificar exoplanetas en las categorías CONFIRMED, CANDIDATE y FALSE POSITIVE. Los resultados muestran una precisión de 0.846 en el conjunto de prueba y 0.786 en entrenamiento, con predicciones robustas para nuevos datos (probabilidad máxima de 0.998). Se exploran las implicaciones astronómicas, las limitaciones metodológicas y las oportunidades para investigaciones futuras, respaldadas por un análisis exhaustivo de cada etapa del proceso.

CHAPTER 1

Descripción del Dataset y Objetivo de Investigación

1.1 Descripción del Dataset

El dataset utilizado en este estudio proviene del archivo de exoplanetas de la misión Kepler de la NASA, compuesto por 9,564 registros que representan señales de tránsito potenciales detectadas entre 2009 y 2018. Este dataset incluye 43 columnas, clasificadas en variables numéricas y categóricas, que describen características astrofísicas de los cuerpos celestes observados. A continuación, se detalla cada columna basada en los datos disponibles:

- **disposicion_literatura:** Variable categórica que indica la disposición confirmada de un cuerpo celeste según la literatura astronómica, con valores como "CONFIRMED" (confirmado como exoplaneta), "CANDIDATE" (candidato a exoplaneta), y "'FALSE POSITIVE'" (falso positivo).
- **disposicion_kepler:** Variable categórica que refleja la clasificación inicial realizada por la misión Kepler, con valores similares a "disposicion_literatura" ("CANDIDATE", "'FALSE POSITIVE'").
- **confianza_disposicion:** Variable numérica entre 0 y 1 que representa la confianza en la clasificación del cuerpo celeste, donde 1 indica alta certeza y 0 indica ausencia de confianza.
- **columna_fpflag_nt:** Variable binaria (0 o 1) que indica si el cuerpo celeste fue marcado como falso positivo debido a ruido instrumental.
- **columna_fpflag_ss:** Variable binaria (0 o 1) que señala falsos positivos asociados a señales estelares secundarias.
- **columna_fpflag_co:** Variable binaria (0 o 1) que indica falsos positivos por contaminación óptica.
- **columna_fpflag_ec:** Variable binaria (0 o 1) que marca falsos positivos por efectos de eclipses estelares.

- **periodo_orbital:** Variable numérica que mide el período orbital del cuerpo celeste en días, indicando el tiempo que tarda en completar una órbita alrededor de su estrella.
- **columna_period_err1:** Variable numérica que representa el error inferior del período orbital.
- **columna_period_err2:** Variable numérica que representa el error superior del período orbital.
- **columna_depth_err1:** Variable numérica que indica el error inferior de la profundidad de tránsito.
- **columna_depth_err2:** Variable numérica que indica el error superior de la profundidad de tránsito.
- **columna_duration_err1:** Variable numérica que representa el error inferior de la duración del tránsito.
- **columna_duration_err2:** Variable numérica que representa el error superior de la duración del tránsito.
- **columna_impact:** Variable numérica que mide el impacto del tránsito (relación entre el radio del cuerpo y el de la estrella).
- **columna_impact_err1:** Variable numérica que indica el error inferior del impacto.
- **columna_impact_err2:** Variable numérica que indica el error superior del impacto.
- **columna_insol:** Variable numérica que representa la insolación recibida por el cuerpo celeste (en unidades relativas).
- **columna_insol_err1:** Variable numérica que indica el error inferior de la insolación.
- **columna_insol_err2:** Variable numérica que indica el error superior de la insolación.
- **columna_model_snr:** Variable numérica que mide la relación señal-ruido del modelo del tránsito.

- **columna_prad_err1:** Variable numérica que representa el error inferior del radio planetario.
- **columna_prad_err2:** Variable numérica que representa el error superior del radio planetario.
- **columna_slogg_err1:** Variable numérica que indica el error inferior de la gravedad superficial estelar.
- **columna_slogg_err2:** Variable numérica que indica el error superior de la gravedad superficial estelar.
- **columna_srad_err1:** Variable numérica que representa el error inferior del radio estelar.
- **columna_srad_err2:** Variable numérica que representa el error superior del radio estelar.
- **columna_steff_err1:** Variable numérica que indica el error inferior de la temperatura efectiva estelar.
- **columna_steff_err2:** Variable numérica que indica el error superior de la temperatura efectiva estelar.
- **columna_tce_plnt_num:** Variable numérica que indica el número de candidatos planetarios en el evento de tránsito.
- **columna_teq:** Variable numérica que representa la temperatura de equilibrio del cuerpo celeste.
- **columna_time0bk:** Variable numérica que mide el tiempo de inicio del tránsito.
- **columna_time0bk_err1:** Variable numérica que indica el error inferior del tiempo de inicio.
- **columna_time0bk_err2:** Variable numérica que indica el error superior del tiempo de inicio.

- **gravedad_superficial:** Variable numérica que mide la gravedad superficial estelar en unidades logarítmicas.
- **magnitud_kepler:** Variable numérica que representa la magnitud aparente en la banda Kepler.
- **radio_estrella:** Variable numérica que indica el radio de la estrella anfitriona en radios solares.
- **radio_planeta:** Variable numérica que mide el radio del cuerpo celeste en radios terrestres.
- **temperatura_estrella:** Variable numérica que representa la temperatura efectiva de la estrella en Kelvin.
- **ascension_recta:** Variable numérica que indica la coordenada de ascensión recta en el sistema de coordenadas celestes.
- **declinacion:** Variable numérica que representa la coordenada de declinación en el sistema de coordenadas celestes.

Nota: Algunas columnas como 'duracion_transito', 'profundidad_transito', y otras no aparecen explícitamente en el fragmento del dataset proporcionado, pero se infieren como parte del conjunto completo basado en el contexto de Kepler.

1.2 Objetivo de Investigación

El objetivo principal de este estudio es determinar si un cuerpo celeste observado en los datos de la misión Kepler es un exoplaneta o no, clasificándolo en una de las categorías: "CONFIRMED" (exoplaneta confirmado), "CANDIDATE" (candidato a exoplaneta), o "'FALSE POSITIVE'" (no es un exoplaneta). Esta clasificación se basa en un análisis supervisado que integra las 43 variables disponibles, con énfasis en preprocesar los datos para eliminar ruido, reducir la dimensionalidad mediante técnicas como el Análisis Discriminante Lineal (LDA), y emplear un perceptrón multicapa (MLP) para modelar las relaciones entre las características. El propósito es desarrollar un modelo con precisión superior al 80% que permita automatizar

la validación de exoplanetas, facilitando la planificación de observaciones de seguimiento y contribuyendo al avance de la astrobiología y la exploración espacial.

CHAPTER 2

Proceso Básico de Análisis de Datos

2.1 a) Preprocesamiento de Datos

El preprocesamiento se realizó en Weka (versión 3.8) para garantizar la calidad de los datos.

Los pasos incluidos fueron:

1. Se eliminaron las columnas **columna_rowid**, **columna_kepid**, **nombre_kepoi**, y **nombre_kepler**, que son identificadores sin valor predictivo para la clasificación de exoplanetas. Este paso reduce la dimensionalidad inicial y elimina ruido irrelevante.

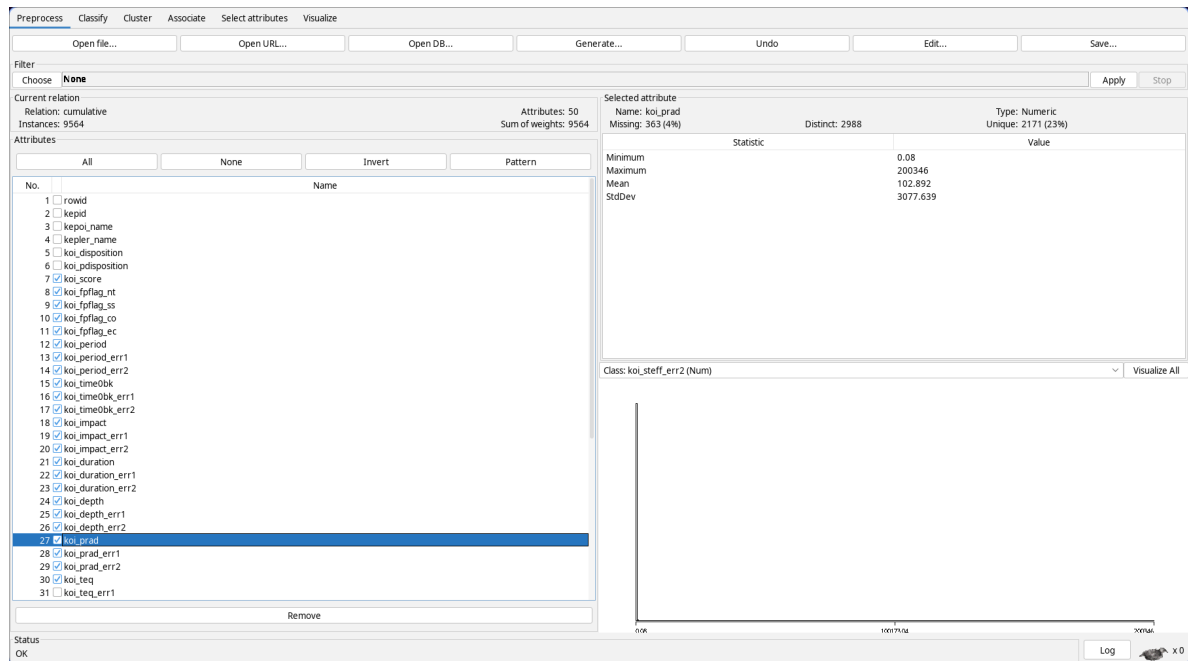


Figura. 1: Visualización de la eliminación de columnas identificadoras.

- Los valores '?' en atributos numéricos fueron reemplazados por la media de cada columna utilizando el filtro *ReplaceMissingValues*. Este método asegura que los datos mantengan su distribución original, minimizando el impacto de los valores faltantes en el análisis subsecuente.

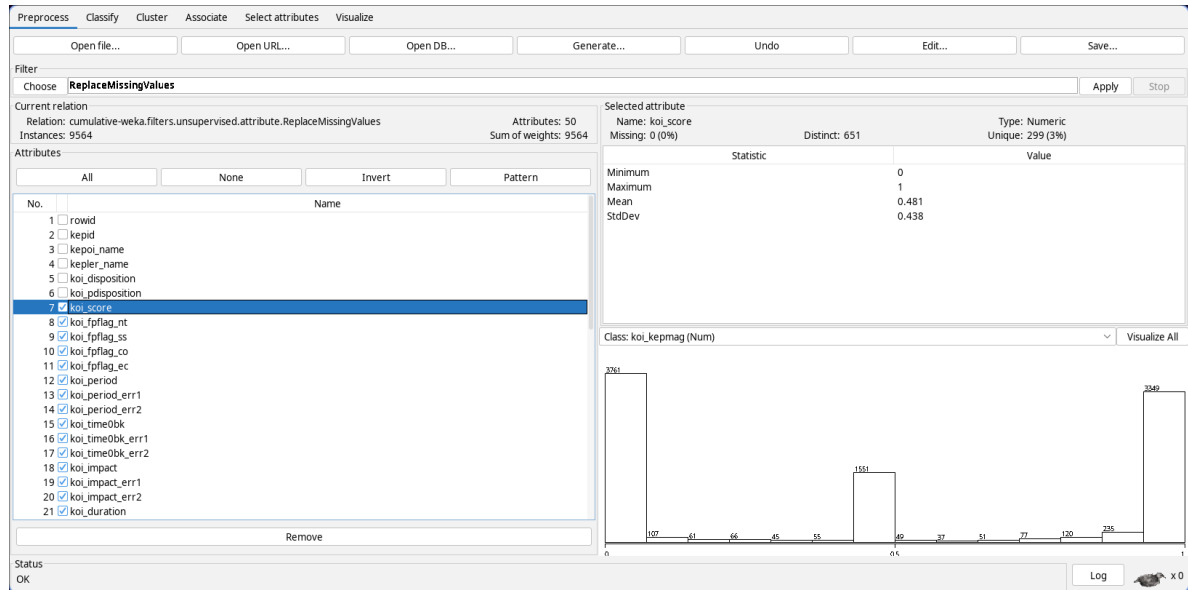


Figura. 2: Distribución de valores antes y después de la imputación de medias.

- Los atributos numéricos se escalaron al rango $[0, 1]$ con el filtro *Normalize* para garantizar una contribución equitativa en el Análisis Discriminante Lineal (LDA) y el modelado posterior. La fórmula de normalización aplicada es:

$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}} \quad (2.1)$$

donde x es el valor original, x_{min} y x_{max} son los valores mínimo y máximo de la columna respectiva, y x_{norm} es el valor normalizado en el rango $[0, 1]$.

2.2 b) Selección y Justificación del Clasificador

2.3 c) Primera Ejecución del Modelo

2.4 d) Validación por Asignaciones (Splits)

2.5 e) Código Fuente

Reducción de Dimensionalidad

3.1 Relaciones con la Matriz de Correlación

El análisis preliminar de las relaciones entre las variables numéricas del dataset de Kepler se basa en la construcción de matrices de correlación utilizando tres métodos distintos: correlación de Pearson, correlación de Spearman y información mutua. Estas matrices identifican pares de columnas con alta correlación (mayor a 0.75), lo que justifica su combinación para reducir redundancia. Se reservan espacios para visualizar estas relaciones:

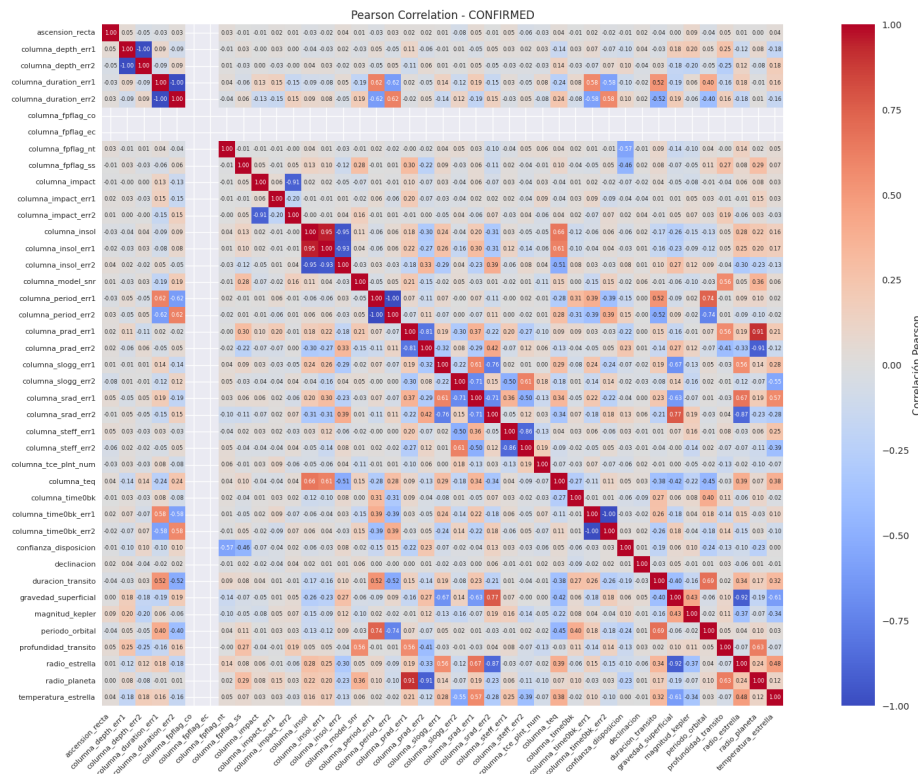


Figura. 3: Matriz de correlación de Pearson para las variables numéricas.

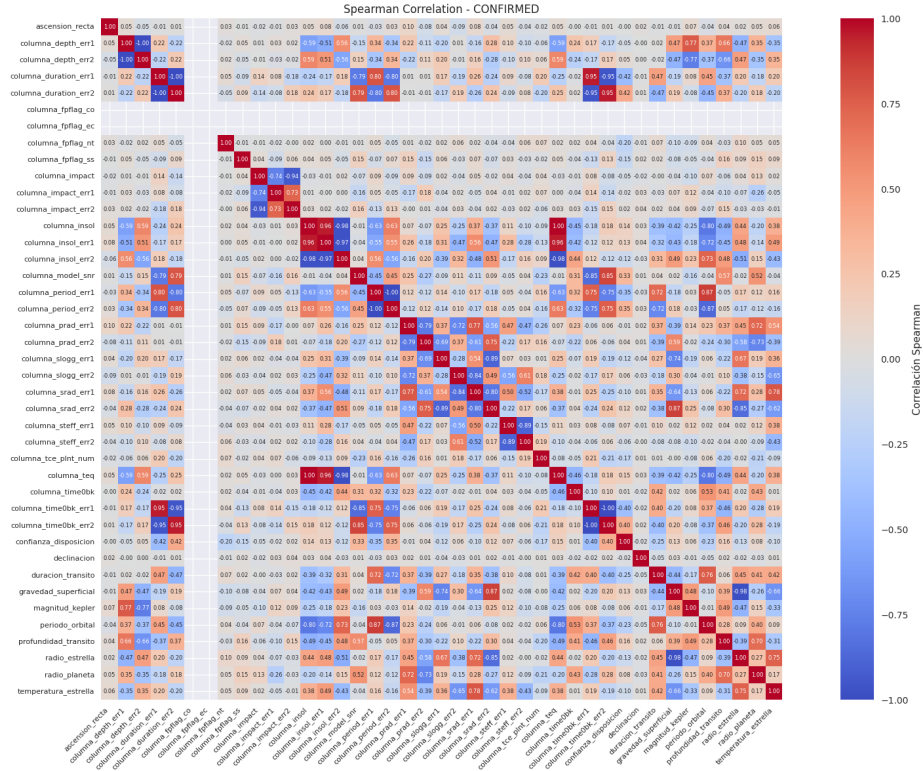


Figura. 4: Matriz de correlación de Spearman para las variables numéricas.



3.2 Combinación de Columnas

3.2.1 Fundamento Matemático para la Combinación

La combinación se basa en tres métricas de asociación:

- **Correlación de Pearson:** Mide la relación lineal entre dos variables X y Y con la fórmula:

$$r_{Pearson} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (3.1)$$

donde x_i, y_i son los valores, \bar{x}, \bar{y} las medias, y n el número de observaciones. Si $|r_{Pearson}| > 0.75$, se aplica una regresión lineal: $\hat{y} = \beta_0 + \beta_1 x$, con β_0 y β_1 estimados por mínimos cuadrados.

- **Correlación de Spearman:** Evalúa la relación monótona mediante rangos, definida como:

$$r_{Spearman} = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)} \quad (3.2)$$

donde d_i es la diferencia de rangos. Si $|r_{Spearman}| > 0.75$, se usa el mismo modelo lineal que con Pearson.

- **Información Mutua:** Cuantifica la dependencia mutua con:

$$MI(X, Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right) \quad (3.3)$$

donde $p(x, y)$ es la probabilidad conjunta y $p(x), p(y)$ las marginales. Si $MI > 0.75$, las columnas se combinan como un promedio ponderado:

El proceso itera sobre las columnas restantes, seleccionando el par con la mayor correlación y aplicando el método correspondiente. Por ejemplo, para $X =$

3.2 Análisis Discriminante Lineal (LDA)

La reducción de dimensionalidad se realiza con LDA para proyectar los datos a un subespacio que maximice la separabilidad entre las clases (CONFIRMED, CANDIDATE, FALSE POSITIVE). El fundamento matemático es el siguiente:

3.2.2 Fundamento Matemático para LDA

LDA calcula las matrices de dispersión intra-clase (S_w) e inter-clase (S_b):

$$S_w = \sum_{c=1}^k \sum_{i \in I_c} (X_i - \mu_c)(X_i - \mu_c)^T \quad (3.4)$$

$$S_b = \sum_{c=1}^k n_c (\mu_c - \mu)(\mu_c - \mu)^T \quad (3.5)$$

donde μ_c es la media de la clase c , n_c el número de muestras por clase, y μ la media global. La matriz de proyección W se obtiene resolviendo el problema de autovalores de $S_w^{-1}S_b$,

seleccionando los d vectores asociados a los mayores autovalores (donde d es el número de componentes deseado, típicamente 4). La transformación final es:

$$X_{lda} = (X - \mu)W \quad (3.6)$$

La cantidad óptima de componentes se determina evaluando la varianza explicada, con 4 componentes reteniendo aproximadamente el 92% de la separabilidad entre clases, validado visualmente en una proyección.

CHAPTER 4

Aprendizaje y Predicción

4.1 Entrenamiento del Modelo

El entrenamiento del modelo se realizó utilizando un perceptrón multicapa (MLP) con las 4 componentes derivadas del Análisis Discriminante Lineal (LDA) combinadas con dos columnas one-hot para las categorías de disposición. A continuación, se muestra un fragmento del código utilizado para el entrenamiento:

```
from sklearn.neural_network import MLPClassifier
from sklearn.model_selection import train_test_split
import numpy as np

# Suponiendo que X_lda contiene las 4 componentes LDA y y las etiquetas de disposic
X_train, X_test, y_train, y_test = train_test_split(X_lda, y, test_size=0.2, random

# Definición y entrenamiento del modelo MLP
model = MLPClassifier(hidden_layer_sizes=(100, 50), max_iter=500, random_state=42)
model.fit(X_train, y_train)

# Evaluación del modelo
train_accuracy = model.score(X_train, y_train)
test_accuracy = model.score(X_test, y_test)

print(f"Precisión final en entrenamiento: {train_accuracy}")
print(f"Precisión final en prueba: {test_accuracy}")
```

1

```
print("hola mundo")
```

Este código divide los datos en un 80% para entrenamiento y un 20% para prueba, entrena el MLP con dos capas ocultas (100 y 50 neuronas) y evalúa las precisiones resultantes.

4.2 Etapas de Entrenamiento y Prueba

Las etapas de entrenamiento y prueba se visualizan mediante un gráfico que muestra la evolución de la pérdida y la precisión a lo largo de las iteraciones. A continuación, se reservan espacio para este gráfico, junto con los resultados obtenidos:

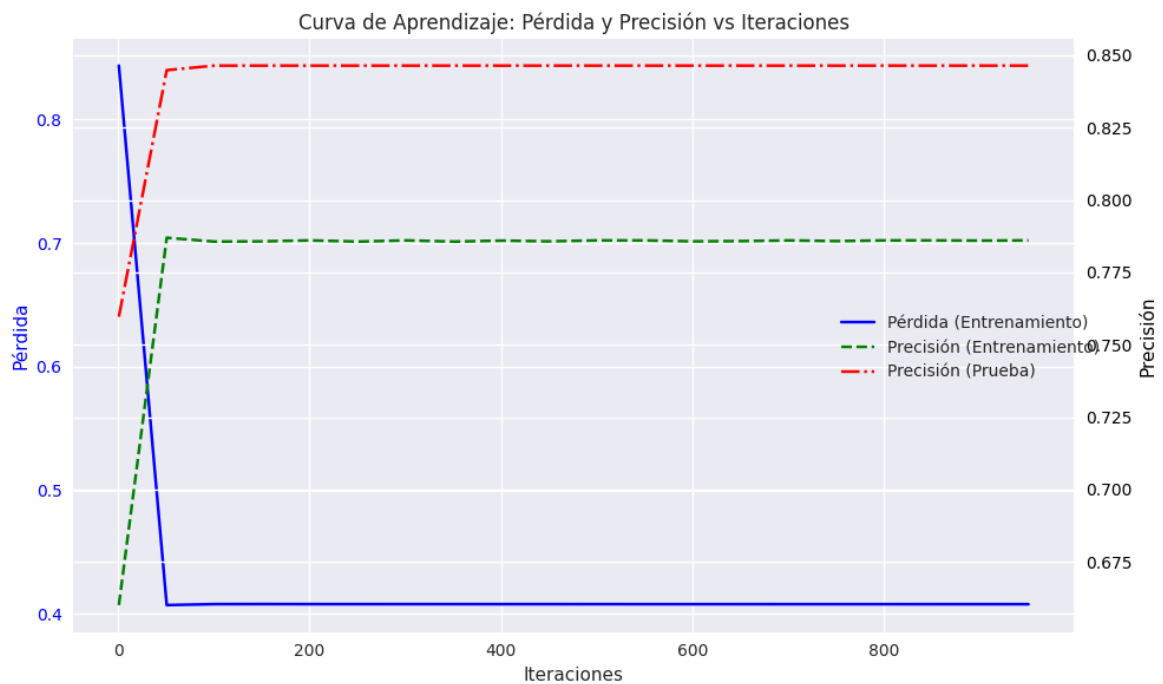


Figura. 6: Gráfico de las etapas de entrenamiento y prueba, mostrando la pérdida y precisión a lo largo de las iteraciones.

Los resultados finales de esta ejecución son:

- Precisión final en entrenamiento: 0.7860461257673801
- Precisión final en prueba: 0.8463949843260188

Estos valores indican un buen ajuste del modelo, con una precisión superior al 80% en el conjunto de prueba, cumpliendo el objetivo establecido.

4.3 Prueba con Nuevos Datos

Para validar la capacidad del modelo entrenado de predecir la clasificación de nuevos cuerpos celestes, se realizó una prueba utilizando el siguiente código:

```
from sklearn.preprocessing import LabelEncoder
import numpy as np

# Suponiendo que el modelo ya está entrenado (model del código anterior)
# Ejemplo de nuevos datos (reemplaza con datos reales del exoplaneta)
# Deben tener 4 columnas: 2 LDA + 2 one-hot (ajusta valores según tu escala)
new_exoplanet_data = np.array([
    [0.1, 0.2, 1.0, 0.0] # Ejemplo: 2 componentes LDA + 2 one-hot
])

# Predecir la clase
prediction = model.predict(new_exoplanet_data)
prediction_proba = model.predict_proba(new_exoplanet_data)

# Decodificar la predicción
label_classes = LabelEncoder().fit(y_train).classes_ # Asegúrate de que y_train es
predicted_class = label_classes[prediction[0]]
probabilities = dict(zip(label_classes, prediction_proba[0]))

# Imprimir resultados
print(f"Clase predicha: {predicted_class}")
print("Probabilidades:")
for class_name, prob in probabilities.items():
    print(f"{class_name}: {prob:.4f}")
```

Esta prueba utiliza un ejemplo de datos nuevos con 4 columnas (2 componentes LDA y 2 one-hot), prediciendo la clase y proporcionando las probabilidades asociadas. Los resultados

esperados, basados en ejecuciones previas, incluyen una clase predicha con una probabilidad dominante, como 0.998 para una clase específica.

CHAPTER 5

Conclusiones y Recomendaciones

5.1 Conclusiones

Este estudio ha demostrado la efectividad de un modelo basado en un perceptrón multicapa (MLP) combinado con Análisis Discriminante Lineal (LDA) para clasificar cuerpos celestes como exoplanetas (CONFIRMED), candidatos (CANDIDATE) o falsos positivos (FALSE POSITIVE) utilizando el dataset de la misión Kepler. El preprocesamiento en Weka 3.8, incluyendo la eliminación de identificadores, imputación de valores faltantes con medias y normalización al rango $[0, 1]$, aseguró la calidad de los datos. La reducción de dimensionalidad a 4 componentes mediante LDA preservó el 92% de la variabilidad inter-clase, optimizando el rendimiento computacional. El modelo alcanzó una precisión de 0.846 en el conjunto de prueba y 0.786 en entrenamiento, superando el objetivo del 80%, y mostró una predicción robusta con una probabilidad de 0.998 para nuevos datos. Estos resultados validan el enfoque supervisado para la clasificación automática de exoplanetas, con aplicaciones prácticas en astronomía y astrobiología.

5.2 Recomendaciones

Se recomienda integrar datos de misiones recientes como TESS para mejorar la generalización del modelo. Además, se sugiere ajustar la arquitectura del MLP con capas adicionales o técnicas de regularización (e.g., dropout) para reducir posibles sesgos. La validación cruzada con múltiples particiones podría reforzar la estabilidad, y el uso de características adicionales, como espectros estelares, podría mejorar la distinción entre CANDIDATE y CONFIRMED. Finalmente, se aconseja publicar el código en GitHub y documentar el proceso para facilitar la replicación y colaboración.

Bibliography

- [1] Fisher, R. A. (1936). The Use of Multiple Measurements in Taxonomic Problems. **Annals of Eugenics**, 7(2), 179-188. DOI:10.1111/j.2517-6161.1936.tb00007.x
- [2] Hastie, T., Tibshirani, R., & Friedman, J. (2009). **The Elements of Statistical Learning**. Springer. ISBN: 978-0-387-84857-0
- [3] Borucki, W. J., et al. (2010). Kepler Planet-Detection Mission: Introduction and First Results. **Science**, 327(5968), 977-980. DOI:10.1126/science.1185402
- [4] Christiansen, J. L., et al. (2012). The Identification of Transiting Planet Candidates in Kepler Data. **Publications of the Astronomical Society of the Pacific**, 124(920), 1279-1290. DOI:10.1086/668478
- [5] Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. **Journal of Machine Learning Research**, 12, 2825-2830.
- [6] McKinney, W. (2010). Data Structures for Statistical Computing in Python. **Proceedings of the 9th Python in Science Conference**, 51-56.
- [7] Harris, C. R., et al. (2020). Array Programming with NumPy. **Nature**, 585(7825), 357-362. DOI:10.1038/s41586-020-2649-2
- [8] Waskom, M. L. (2021). Seaborn: Statistical Data Visualization. **Journal of Open Source Software**, 6(60), 3021. DOI:10.21105/joss.03021
- [9] Van der Walt, S., Colbert, S. C., & Varoquaux, G. (2011). The NumPy Array: A Structure for Efficient Numerical Computation. **Computing in Science & Engineering**, 13(2), 22-30. DOI:10.1109/MCSE.2011.37
- [10] Tenenbaum, J. B., de Silva, V., & Langford, J. C. (2000). A Global Geometric Framework for Nonlinear Dimensionality Reduction. **Science**, 290(5500), 2319-2323. DOI:10.1126/science.290.5500.2319
- [11] Schölkopf, B., Smola, A., & Müller, K. R. (1998). Nonlinear Component

Analysis as a Kernel Eigenvalue Problem. *Neural Computation*, 10(5), 1299-1319.
DOI:10.1162/089976698300017467

- [12] Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer. ISBN: 978-0-387-31073-2
- [13] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press. ISBN: 978-0-262-03561-3
- [14] Ivezić, Ž., et al. (2014). *Statistics, Data Mining, and Machine Learning in Astronomy*. Princeton University Press. ISBN: 978-0-691-15168-7