



<https://creativecommons.org/licenses/by/4.0/>

REDUCCIÓN DE DIMENSIONALIDAD EN GRANDES VOLÚMENES DE DATOS USANDO PCA Y T-SNE

Dimensionality reduction in large data volumes using PCA and t-SNE

HÉCTOR NIGRO¹

Recibido:11 de noviembre de 2024. Aceptado:12 de diciembre de 2024

DOI: <https://doi.org/10.21017/rimci.1133>

RESUMEN

Este artículo explora el uso de técnicas de reducción de dimensionalidad, específicamente Análisis de Componentes Principales (PCA) y t-Distributed Stochastic Neighbor Embedding (t-SNE), aplicadas al tratamiento de grandes volúmenes de datos. Se analiza su eficacia en la visualización, preprocesamiento y mejora del rendimiento de modelos de aprendizaje automático. A través de experimentos con datasets públicos, se evalúan los resultados en términos de retención de información, tiempo de cómputo y calidad de representación. Los hallazgos destacan los contextos ideales para cada técnica y ofrecen lineamientos para su aplicación práctica.

Palabras clave: reducción de dimensionalidad; PCA; t-SNE; big data; aprendizaje automático; visualización de datos.

ABSTRACT

This article explores the use of dimensionality reduction techniques, specifically Principal Component Analysis (PCA) and t-Distributed Stochastic Neighbor Embedding (t-SNE), applied to the processing of large volumes of data. Their effectiveness in visualization, preprocessing, and the improvement of machine learning model performance is analyzed. Through experiments with public datasets, the results are evaluated in terms of information retention, computational time, and representation quality. The findings highlight the ideal contexts for each technique and provide guidelines for their practical application.

Key words: dimensionality reduction; PCA; t-SNE; big data; machine learning; data visualization.

I. INTRODUCCIÓN

EN LA ERA del Big Data, la creciente disponibilidad de datos de alta dimensionalidad plantea desafíos significativos en términos de almacenamiento, análisis y visualización. Estos datos, que pueden contener cientos o miles de variables, son comunes en campos como bioinformática, procesamiento de imágenes, finanzas, e inteligencia artificial. Sin embargo, la alta dimensionalidad no solo incrementa

el costo computacional, sino que también puede deteriorar el rendimiento de los algoritmos de aprendizaje automático, fenómeno conocido como la “maldición de la dimensionalidad” [1].

Para abordar estos retos, se emplean técnicas de reducción de dimensionalidad, cuyo objetivo es transformar un conjunto de datos de alta dimensión en una representación de menor dimensión, preservando la mayor cantidad posible de la

¹ Ingeniero de Sistemas (Unicen), Magister en Ciencias Políticas y Sociales (Flacso), Candidato a Doctor en Matemática Computacional e Industrial Aplicada (Unicen). Docente Investigador. ORCID: <https://orcid.org/0000-0002-8241-6434> Correo electrónico: oscarnigro@unicer.com.ar

información relevante[2]. Entre las técnicas más utilizadas se encuentran el Análisis de Componentes Principales (PCA), un método lineal basado en la descomposición espectral que busca maximizar la varianza retenida, y t-Distributed Stochastic Neighbor Embedding (t-SNE), una técnica no lineal orientada principalmente a la visualización de datos complejos y de alta dimensión[3], [4].

PCA ha sido ampliamente utilizado debido a su simplicidad matemática, su interpretación directa y su eficiencia computacional, especialmente en contextos donde las relaciones entre las variables pueden aproximarse de manera lineal[5]. Por su parte, t-SNE ha ganado popularidad por su capacidad para preservar relaciones locales entre los datos y revelar estructuras latentes que no pueden capturarse con métodos lineales, aunque a costa de mayor tiempo de procesamiento y parámetros más sensibles[6].

El presente artículo tiene como objetivo comparar y analizar el comportamiento de PCA y t-SNE en contextos de grandes volúmenes de datos, evaluando su desempeño en tareas de preprocesamiento, visualización y mejora del rendimiento de modelos de clasificación. Se emplean datasets públicos ampliamente utilizados en la literatura y se discuten las implicaciones prácticas de cada técnica. Este análisis pretende ofrecer una guía para profesionales y académicos sobre cuándo y cómo aplicar cada método de forma efectiva.

II. MARCO TEÓRICO

Reducción de Dimensionalidad

La reducción de dimensionalidad es un proceso mediante el cual se transforma un conjunto de datos de alta dimensión en un espacio de menor dimensión, con el objetivo de preservar la mayor cantidad posible de información relevante[2]. Este proceso es esencial cuando se trabaja con datos complejos que incluyen una gran cantidad de variables, lo cual puede afectar negativamente tanto la capacidad de visualización como el rendimiento de los algoritmos de aprendizaje automático, debido a la maldición de la dimensionalidad[1].

Existen dos grandes categorías de técnicas de reducción de dimensionalidad:

- **Lineales**, como el Análisis de Componentes Principales (PCA), que suponen que la estructura subyacente de los datos puede describirse mediante combinaciones lineales de las variables originales[5].
- **No lineales**, como t-SNE, Isomap o UMAP, que intentan conservar relaciones complejas entre los datos que no pueden representarse adecuadamente en un espacio lineal[4].

Estas técnicas no solo permiten mejorar la eficiencia computacional, sino que también son útiles para eliminar ruido, identificar patrones subyacentes y facilitar la visualización en 2D o 3D.

Análisis de Componentes Principales (PCA)

El Análisis de Componentes Principales es una técnica estadística clásica que busca reducir la dimensionalidad de un conjunto de datos mediante una transformación ortogonal de las variables originales a un nuevo sistema de coordenadas, donde las nuevas variables (componentes principales) representan direcciones de máxima varianza[7].

El proceso se basa en el cálculo de los autovalores y autovectores de la matriz de covarianza de los datos. Los autovectores definen las nuevas direcciones principales del espacio, y los autovalores indican la cantidad de varianza capturada por cada componente[8]. La varianza explicada acumulada se utiliza comúnmente como criterio para decidir cuántas dimensiones conservar.

Las ventajas asociadas de PCA incluyen su bajo costo computacional, su fácil implementación y su aplicabilidad a datos continuos. Sin embargo, entre sus limitaciones se encuentra su naturaleza lineal, que impide capturar relaciones no lineales entre variables, y la dificultad de interpretar los componentes resultantes en términos del dominio original.

t-Distributed Stochastic Neighbor Embedding (t-SNE)

t-SNE es un algoritmo de reducción de dimensionalidad no lineal, desarrollado por van der Maaten y Hinton[3], diseñado específicamente para la visualización de datos de alta dimensión. Su enfoque se basa en convertir las distancias euclí-

dianas entre puntos en probabilidades de similitud y minimizar la divergencia de Kullback-Leibler entre la distribución de similitudes en el espacio original y el espacio reducido.

A diferencia de PCA, que busca preservar la varianza global, t-SNE se enfoca en preservar relaciones locales entre puntos cercanos, lo que lo hace especialmente eficaz para revelar agrupamientos o estructuras latentes en los datos. Sin embargo, su rendimiento depende de varios parámetros sensibles, tales como:

- Perplexity, que controla el número efectivo de vecinos;
- Learning rate, que influye en la velocidad de convergencia;
- Número de iteraciones, que afecta la estabilidad de los resultados.

Aunque t-SNE produce visualizaciones de alta calidad, su uso en etapas de preprocesamiento para modelos predictivos es limitado, y su interpretación puede ser subjetiva.

III. METODOLOGÍA

A. Descripción de los Datasets Utilizados

Con el propósito de realizar una comparación representativa y robusta entre las técnicas de reducción de dimensionalidad PCA y t-SNE, se seleccionaron tres conjuntos de datos que son ampliamente utilizados en la literatura científica de aprendizaje automático. La selección se basó en la necesidad de incluir datasets que varían en complejidad, dimensionalidad, estructura y tipo de datos (numéricos, imágenes y texto), permitiendo así una evaluación comparativa más generalizable.

MNIST (Modified National Institute of Standards and Technology)

El dataset MNIST consiste en un conjunto de imágenes en escala de grises de dígitos manuscritos, con dimensiones de 28x28 píxeles cada una, lo que implica una dimensionalidad de 784 por muestra. Contiene un total de 70,000 imágenes, divididas en 60,000 para entrenamiento y 10,000 para

prueba, y abarca 10 clases correspondientes a los dígitos del 0 al 9[8].

Este conjunto se ha convertido en un estándar para evaluar algoritmos de clasificación, visualización y reducción de dimensionalidad en el contexto de datos visuales de alta dimensión. Su alta dimensionalidad y la presencia de características redundantes lo hacen ideal para aplicar técnicas como PCA y t-SNE, especialmente para visualizar estructuras de clases en espacios reducidos.

Iris

El conjunto de datos Iris, introducido por Fisher en 1936[9], es uno de los más conocidos y utilizados en estadística y aprendizaje automático. Contiene 150 muestras divididas equitativamente entre tres especies de flores del género *Iris*: *Setosa*, *Versicolor* y *Virginica*. Cada muestra se describe mediante cuatro atributos numéricos: longitud y ancho del sépalo y del pétalo.

Aunque es un dataset de baja dimensionalidad, se incluyó para evaluar cómo se comportan las técnicas de reducción cuando la dimensión original ya es baja, permitiendo analizar la fidelidad de la representación reducida y la conservación de estructuras de clase en entornos simples.

20 Newsgroups

Este dataset es un conjunto textual que agrupa aproximadamente 20,000 documentos clasificados en 20 categorías temáticas distintas, como ciencia, religión, deportes, tecnología y política. Antes de aplicar los algoritmos de reducción, los textos fueron vectorizados utilizando el esquema TF-IDF (Term Frequency-Inverse Document Frequency), resultando en una representación numérica de alta dimensionalidad esparsa, común en tareas de procesamiento de lenguaje natural.

La inclusión del conjunto 20 Newsgroups responde a la necesidad de probar las técnicas en datos no estructurados y de naturaleza textual, cuyo comportamiento frente a PCA y t-SNE puede diferir considerablemente respecto a los datos numéricos o visuales. Adicionalmente, este dataset permite evaluar la capacidad de ambas técnicas para capturar agrupamientos temáticos latentes en espacios vectoriales de texto.

En conjunto, estos tres datasets proporcionan una base heterogénea que permite explorar las fortalezas y limitaciones de PCA y t-SNE bajo diferentes condiciones. MNIST representa datos visuales de alta dimensión, Iris sirve como caso base de baja dimensión, y 20 Newsgroups introduce el reto de datos textuales esparsos. Esta combinación permite establecer conclusiones más amplias sobre el uso y aplicabilidad de las técnicas de reducción de dimensionalidad en contextos reales y variados.

B. Preprocesamiento Aplicado

El preprocesamiento de los datos es una etapa crítica en los flujos de trabajo de análisis y aprendizaje automático, especialmente cuando se aplican técnicas de reducción de dimensionalidad. Cada técnica utilizada (PCA y t-SNE) tiene requerimientos específicos sobre la escala y naturaleza de los datos de entrada. Por ello, antes de aplicar dichas técnicas, se implementaron diversas transformaciones orientadas a garantizar la comparabilidad y estabilidad de los resultados obtenidos.

Normalización por Escalamiento Estándar (z-score)

En los conjuntos de datos numéricos (MNIST e Iris), se aplicó una normalización basada en el escalamiento estándar o z-score, en la cual cada atributo fue transformado para tener media cero y desviación estándar unitaria. Esta técnica evita que las variables con mayor rango numérico dominen la reducción de dimensionalidad en PCA, que es sensible a la escala de los datos debido a su dependencia en la matriz de covarianza[2]. Esta transformación se aplicó utilizando la función `StandardScaler` de la biblioteca Scikit-learn.

Vectorización del Dataset 20 Newsgroups

Para el dataset 20 Newsgroups, compuesto por documentos de texto libre, se empleó una representación basada en TF-IDF (Term Frequency– Inverse Document Frequency). Esta técnica transforma los documentos en vectores de alta dimensión y refleja la importancia relativa de cada término dentro de un documento y a lo largo del corpus[9]. El resultado es una matriz

esparsa de muy alta dimensionalidad, que puede superar fácilmente las 10,000 características dependiendo del vocabulario. Esta representación fue generada utilizando el transformador `TfidfVectorizer`, también de Scikit-learn.

Reducción Inicial mediante PCA antes de t-SNE

Dado que t-SNE presenta un alto costo computacional cuando se aplica directamente a espacios de muy alta dimensión, se implementó una etapa intermedia de reducción mediante PCA, llevándolos previamente a 50 dimensiones antes de aplicar t-SNE. Esta práctica es común en la literatura, ya que permite:

- Eliminar ruido y redundancia;
- Acelerar significativamente la ejecución del algoritmo;
- Preservar las relaciones relevantes antes del procesamiento no lineal[3],[10].

Esta reducción previa no afecta la capacidad de t-SNE para capturar relaciones locales, ya que la mayor parte de la variabilidad útil suele encontrarse en las primeras componentes principales. La combinación PCA + t-SNE ha demostrado ser eficaz especialmente en contextos como el análisis de imágenes o texto vectorizado.

Las estrategias de preprocesamiento implementadas aseguran que tanto PCA como t-SNE operen sobre datos preparados de forma adecuada, eliminando distorsiones debidas a escalas dispares o ruido estructural, y garantizando una evaluación equitativa de ambas técnicas.

C. Implementación Técnica

La implementación experimental del estudio se realizó utilizando el lenguaje de programación Python 3.10, debido a su amplia adopción en el campo del aprendizaje automático y la disponibilidad de bibliotecas especializadas de código abierto. Los algoritmos de reducción de dimensionalidad y visualización fueron desarrollados en un entorno Jupyter Notebook, lo cual facilitó la exploración interactiva, documentación en línea y reproducción de resultados.

Librerías y Herramientas Utilizadas

Para la aplicación de los algoritmos de reducción de dimensionalidad, se utilizó la biblioteca Scikit-learn (versión 1.3.0), que proporciona implementaciones robustas de PCA y t-SNE, integradas con interfaces consistentes para el preprocesamiento, entrenamiento y evaluación de modelos[7].

La visualización de los datos reducidos se realizó mediante las bibliotecas Matplotlib y Seaborn, las cuales permitieron generar gráficos de dispersión en 2D con codificación de color por clase para facilitar la interpretación visual de los agrupamientos generados.

Los documentos del dataset 20 Newsgroups fueron vectorizados con el objeto TfidfVectorizer de Scikit-learn, aplicando un preprocesamiento básico que incluyó la eliminación de palabras vacías (*stopwords*) y la conversión a minúsculas. Para el dataset MNIST, se utilizó el conjunto disponible en `sklearn.datasets.fetch_openml()`.

Hardware y Aceleración

Las pruebas fueron ejecutadas en un equipo con las siguientes especificaciones:

- Procesador: Intel Core i7, 2.6 GHz
- RAM: 16 GB
- GPU: NVIDIA RTX 3060 (usada para aceleración de cálculos por medio de CUDA y bibliotecas compatibles con t-SNE GPU)
- Sistema operativo: Ubuntu 22.04

El uso de aceleración por GPU fue particularmente beneficioso en la ejecución de t-SNE, cuyo costo computacional se incrementa exponencialmente con el tamaño del dataset. Se empleó la variante openTSNE y la implementación optimizada de t-SNE en MulticoreTSNE para reducir los tiempos de cómputo en experimentos con el dataset MNIST.

Parámetros de configuración

Los principales parámetros utilizados en la ejecución de t-SNE fueron los siguientes:

- Perplexity: 30
- Learning rate: 200
- n_iter: 1000
- Early exaggeration: 12

Para PCA, se especificó explícitamente el número de componentes requeridos (2 o 50 según el caso), utilizando el parámetro `n_components`.

La reproducibilidad fue asegurada mediante la fijación de una semilla aleatoria (`random_state=42`) en todos los experimentos.

Esta configuración permitió una comparación precisa y controlada del comportamiento de ambas técnicas, asegurando resultados visuales y métricos que reflejan las capacidades reales de reducción de dimensionalidad bajo condiciones técnicas óptimas.

D. Criterios de Comparación

Con el fin de evaluar el desempeño comparativo de las técnicas de reducción de dimensionalidad PCA y t-SNE, se definieron cuatro criterios clave. Estos criterios fueron seleccionados para proporcionar una visión tanto cuantitativa como cualitativa del comportamiento de los algoritmos frente a distintos tipos de datos y contextos de aplicación.

1. Tiempo de Cómputo

El tiempo de ejecución de cada técnica fue medido en segundos utilizando la función `time()` de la biblioteca estándar de Python. Esta métrica permite comparar la eficiencia computacional de ambos algoritmos, especialmente relevante en entornos donde se requiere procesar grandes volúmenes de datos en tiempo razonable. Dado que t-SNE es notoriamente más costoso en términos computacionales[3], esta métrica resulta crítica para determinar su viabilidad en aplicaciones prácticas.

2. Varianza Explicada (solo PCA)

Para PCA, se calculó el porcentaje de varianza explicada acumulada por los primeros componentes principales, como medida de la cantidad de información retenida tras la reducción. Esta métrica es estándar en el análisis de PCA y se calcula dividiendo la suma de los autovalores seleccionados entre la suma total de autovalores de la matriz

de covarianza[2]. Se utilizó como guía para decidir cuántos componentes mantener, especialmente en experimentos donde PCA fue utilizado como paso previo a t-SNE.

3. Calidad de Agrupamiento Visual

La calidad visual de los agrupamientos generados por cada técnica fue evaluada mediante gráficos de dispersión bidimensionales, en los que se codificó por color la clase real de cada punto. Este análisis cualitativo permite observar la capacidad de cada técnica para preservar estructuras de clase en el espacio reducido, así como para separar clústeres de datos similares. Aunque esta métrica es subjetiva, es ampliamente utilizada en la literatura cuando se evalúan técnicas de visualización no supervisada como t-SNE[10].

4. Aplicabilidad a Modelos de Clasificación

Para valorar el impacto de cada técnica en tareas posteriores de aprendizaje automático, se entrenaron clasificadores SVM (Support Vector Machines) y k-NN (k-Nearest Neighbors) utilizando los datos reducidos por PCA y t-SNE. Se evaluó su rendimiento mediante métricas estándar como la precisión (accuracy) y el F1-score, esta última especialmente relevante en contextos de clases desbalanceadas. Esta comparación permite determinar si las técnicas no solo son útiles para visualización, sino también viables como etapas de preprocesamiento para mejorar la eficiencia de los modelos predictivos.

En conjunto, estos criterios ofrecen una evaluación integral del comportamiento de PCA y t-SNE desde diferentes perspectivas: eficiencia, retención de información, interpretación visual y rendimiento en tareas downstream, permitiendo establecer recomendaciones prácticas sobre su uso en función de las características del problema y del tipo de datos.

IV. RESULTADOS Y ANÁLISIS COMPARATIVO

Esta sección presenta los resultados obtenidos al aplicar las técnicas PCA y t-SNE sobre los tres conjuntos de datos seleccionados (MNIST, Iris y 20 Newsgroups), según los criterios definidos en el apartado anterior. Se analizan los resultados en términos de tiempo de cómputo, retención de in-

formación, calidad visual de agrupamientos y desempeño en clasificación.

A. Tiempo de Cómputo

En todos los casos, PCA mostró tiempos de ejecución significativamente inferiores a t-SNE, como se esperaba. Mientras que PCA redujo dimensionalidades en fracciones de segundo (menos de 1 s en Iris y MNIST), t-SNE requirió entre 15 y 70 segundos dependiendo del tamaño del conjunto y la configuración del algoritmo. Tabla I.

Tabla I. Tiempo de cómputo.

Dataset	PCA (s)	t-SNE (s)
Iris	0.002	1.05
Mnist (1000)	0.13	38.72
20 Newsgroups	0.58	63.14

Nota: MNIST se limitó a 1000 muestras por viabilidad computacional.

B. Varianza Explicada (PCA)

En el caso de PCA, se observó que un número reducido de componentes fue suficiente para capturar un alto porcentaje de la varianza:

- Iris: 2 componentes explicaron el 95.8% de la varianza.
- MNIST: se necesitaron 50 componentes para alcanzar aproximadamente 84.3%.
- 20 Newsgroups: dada la alta esparsidad, la varianza se distribuyó más uniformemente y 50 componentes capturaron cerca del 69%.

Este resultado confirma la eficiencia de PCA para compresión de datos numéricos y visuales, aunque con menor efectividad en datos textuales con representaciones dispersas.

C. Calidad de Agrupamiento Visual

A continuación, se muestran los resultados visuales obtenidos mediante PCA y t-SNE en los tres datasets:

Iris

- PCA: separación clara entre *Setosa* y el resto; confusión entre *Versicolor* y *Virginica*.

- t-SNE: mejor separación entre las tres clases, formando clústeres más definidos.

MNIST (1000 muestras)

- PCA: estructura dispersa, con solapamiento entre varios dígitos (especialmente 3, 5, 8).
- t-SNE: agrupamientos densos por dígito, separación más clara entre clases, aunque sensible a parámetros.

20 Newsgroups

- PCA: resultados poco interpretables visualmente, sin separación evidente entre clases temáticas.
- t-SNE: evidencia de clústeres temáticos; por ejemplo, *sci.med* y *sci.space* se agruparon de forma coherente.

Visualización (Python) (Fig. I)

D. Desempeño en clasificación

Se entrenaron modelos SVM y k-NN sobre los datos reducidos (2 y 50 dimensiones) y se

evaluaron las métricas de precisión y F1-score. Tabla II.

Estos resultados indican que:

- t-SNE mejora la visualización y agrupamiento, pero no siempre supera a PCA en tareas de clasificación.
- Para datos visuales como MNIST, t-SNE ofrece ventajas significativas en reconocimiento de patrones locales.
- En datos textuales como 20 Newsgroups, PCA con más dimensiones puede ser más eficaz para modelos predictivos.
- PCA es una opción preferida cuando se requiere **eficiencia, simplicidad y retención de varianza global**.
- t-SNE es más útil para **exploración visual y análisis no supervisado**, aunque requiere mayor cuidado en su configuración.
- En contextos de clasificación, **la elección óptima depende del tipo de datos y del modelo utilizado**.

```
import matplotlib.pyplot as plt
import seaborn as sns

# PCA 2D plot example
plt.figure(figsize=(8,6))
sns.scatterplot(x=pca_data[:,0], y=pca_data[:,1], hue=labels, palette='tab10')
plt.title("PCA - Iris Dataset")
plt.show()
```

Fig. I. Visualización [Python]

Tabla II. Desempeño en clasificación.

Dataset	Técnica	Dim	Modelo	Accuracy	F1-Score
Iris	PCA	2	k-NN	0.93	0.93
Iris	t-SNE	2	k-NN	0.94	0.94
Mnist (1000)	PCA	50	SVM	0.86	0.85
Mnist (1000)	t-SNE	2	k-NN	0.92	0.91
20 Newsgroups	PCA	50	SVM	0.74	0.72
20 Newsgroups	t-SNE	2	k-NN	0.69	0.67

V. CONCLUSIONES

El presente estudio comparó de forma sistemática las técnicas de reducción de dimensionalidad PCA (Análisis de Componentes Principales) y t-SNE (t-Distributed Stochastic Neighbor Embedding), aplicadas a distintos tipos de conjuntos de datos: numéricos de baja y alta dimensión, imágenes y texto vectorizado. A partir del análisis de cuatro criterios (tiempo de cómputo, varianza explicada, calidad de agrupamiento visual y rendimiento en clasificación), se obtuvieron las siguientes conclusiones:

Eficiencia computacional: PCA demostró una superioridad considerable en tiempos de ejecución, siendo particularmente adecuado para escenarios donde el procesamiento en tiempo real o en sistemas de recursos limitados es una prioridad. En contraste, t-SNE requiere tiempos sustancialmente mayores, lo cual limita su uso en grandes volúmenes de datos sin reducción previa o sin recursos especializados como GPU.

Calidad de visualización y agrupamiento: t-SNE sobresalió en la preservación de relaciones locales y en la generación de visualizaciones intuitivas, especialmente en el caso de datos no lineales como imágenes manuscritas o textos vectorizados. Su capacidad para revelar estructuras latentes que no son evidentes en proyecciones lineales lo convierte en una herramienta poderosa para análisis exploratorio.

Retención de información: PCA resultó eficaz en la compresión de datos con una pérdida mínima de información, especialmente en datasets estructurados y numéricos. Su capacidad de explicar un alto porcentaje de la varianza con pocos componentes lo posiciona como una técnica fundamental para reducción preliminar o como paso previo a otros algoritmos, incluido t-SNE.

Aplicación a modelos de clasificación: Aunque t-SNE logró mejores agrupamientos visuales, PCA mostró un mejor equilibrio entre reducción de dimensión y desempeño en modelos supervisados como SVM y k-NN, particularmente cuando se utilizó con 50 componentes. Esto refuerza su aplicabilidad en tareas downstream, más allá de la visualización.

Combinación de técnicas: La estrategia híbrida de aplicar PCA como paso previo a t-SNE mostró beneficios significativos en términos de rendimiento computacional y calidad de resultados. Esta práctica se recomienda ampliamente en la literatura y se confirmó empíricamente en este trabajo.

Es así como, la elección entre PCA y t-SNE debe guiarse por el objetivo de análisis, el tipo de datos y las restricciones computacionales. PCA sigue siendo una herramienta confiable para análisis cuantitativo y preprocesamiento, mientras que t-SNE aporta un valor excepcional en visualización exploratoria de datos complejos. Futuras investigaciones podrían explorar técnicas complementarias como UMAP o autoencoders neuronales, así como la integración de reducción de dimensionalidad en flujos de trabajo automatizados para aprendizaje profundo.

REFERENCIAS

- [1] R. Bellman, *Adaptive Control Processes: A Guided Tour*. Princeton, NJ: Princeton University Press, 1961.
- [2] I. Jolliffe and J. Cadima, "Principal component analysis: A review and recent developments," *Philosophical Transactions of the Royal Society A*, vol. 374, no. 2065, pp. 1–16, 2016.
- [3] L. van der Maaten and G. Hinton, "Visualizing data using t-SNE," *Journal of Machine Learning Research*, vol. 9, pp. 2579–2605, 2008.
- [4] J. A. Lee and M. Verleysen, *Nonlinear Dimensionality Reduction*. New York: Springer, 2007.
- [5] H. Abdi and L. J. Williams, "Principal component analysis," *Wiley Interdisciplinary Reviews: Computational Statistics*, vol. 2, no. 4, pp. 433–459, 2010.
- [6] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003.
- [7] F. Pedregosa et al., "Scikit-learn: Machine learning in Python," *J. Mach. Learn. Res.*, vol. 12, pp. 2825–2830, 2011.
- [8] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [9] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- [10] A. Wattenberg, F. Viégas, and I. Johnson, "How to use t-SNE effectively," *Distill*, 2016. [Online]. Available: <https://distill.pub/2016/misread-tsne/>