

Aplicación de Técnicas Avanzadas de Machine Learning: Comparación NAS-SGD para Forecasting de Acciones Microsoft y Empresas Afiliadas

Jimena Yéssica Paricela Yana
Universidad Nacional del Altiplano - UNAP
Ingeniería Estadística e Informática

14 de julio de 2025

Resumen

Este artículo presenta un análisis comparativo entre Neural Architecture Search (NAS) y Optimización Estocástica aplicados al comportamiento de las acciones de Microsoft y compañías afiliadas clave. Se evalúa la eficiencia de ambos métodos utilizando gradientes para predicción futura de valores bursátiles mediante el análisis de series temporales financieras complejas. Los algoritmos fueron implementados utilizando frameworks modernos de deep learning, con arquitecturas que incluyen LSTM, GRU y Transformer para capturar dependencias temporales a largo plazo.

Los resultados experimentales demuestran que mientras NAS ofrece mayor flexibilidad en la búsqueda automática de arquitecturas óptimas con una mejoras en la precisión predictiva, la Optimización Estocástica proporciona convergencia más estable para series temporales financieras con menor varianza en los errores de predicción. El análisis de datos históricos de Microsoft desde 1986 hasta 2024, incluyendo más de 9,800 puntos de datos, revela patrones significativos que pueden ser modelados eficientemente con ambas metodologías. NAS demuestra superioridad en términos de adaptabilidad automática a cambios de mercado, mientras que la Optimización Estocástica sobresale en estabilidad computacional y interpretabilidad de resultados.

Palabras clave: Neural Architecture Search, Optimización Estocástica, Machine Learning, Predicción Bursátil, Microsoft Stock, DARTS, Series Temporales Financieras.

I. Introducción

El mercado bursátil representa uno de los sistemas complejos más dinámicos y desafiantes para la predicción computacional, caracterizado por su alta volatilidad, dependencias no lineales y la influencia de múltiples factores económicos, políticos y sociales [1–3]. La predicción precisa de precios de acciones ha sido durante décadas el "santo grial" de las finanzas cuantitativas, motivando el desarrollo de metodologías cada vez más sofisticadas que integran técnicas avanzadas de inteligencia artificial y aprendizaje automático [4, 5].

En la última década, el campo del *machine learning* ha experimentado avances significativos, particularmente en el área de *deep learning* y optimización automática de arquitecturas neuronales [6]. Dos paradigmas han emergido como especialmente prometedores para el análisis de series temporales financieras: la búsqueda automática de arquitecturas neuronales (*Neural Architecture Search*, *NAS*) y las técnicas avanzadas de optimización estocástica [7–9].

NAS representa un cambio paradigmático en el diseño de redes neuronales, automatizando el proceso tradicional de ingeniería manual de arquitecturas [10, 11]. En el contexto financiero, permite ex-

plorar sistemáticamente espacios complejos de arquitecturas para identificar configuraciones óptimas que capturen patrones temporales no lineales en los datos bursátiles [12]. Por otro lado, las técnicas modernas de optimización estocástica, como Adam, AdaGrad, RMSprop y Adadelata, han demostrado una gran capacidad para manejar la alta complejidad y ruido inherentes en los mercados financieros [13, 14].

Microsoft Corporation, fundada en 1975 y cotizando públicamente desde 1986, representa un caso de estudio ideal para evaluar estas metodologías debido a su larga trayectoria bursátil, alta liquidez y fuerte correlación con el sector tecnológico global [15, 16]. La empresa ha atravesado múltiples ciclos económicos, transformaciones tecnológicas y cambios estratégicos, proporcionando un conjunto de datos rico en patrones estructurales y volatilidad [17].

Este estudio contribuye al campo de las finanzas computacionales mediante: una comparación sistemática entre NAS y optimización estocástica en contextos financieros reales, el análisis de más de 38 años de datos históricos de Microsoft, incluyendo periodos de alta volatilidad como la crisis *dot-com* de 2000 [18] y la crisis financiera de 2008 [19], la implementación de arquitecturas híbridas que combinan fortalezas de ambos enfoques, y la evaluación me-

dianter métricas financieras específicas como el ratio de Sharpe, el *maximum drawdown*, el alpha de Jensen y el índice de Sortino [20].

II. Fundamentos

Machine Learning en Finanzas: El aprendizaje automático ha transformado radicalmente el landscape de las finanzas cuantitativas, evolucionando desde modelos estadísticos tradicionales como ARIMA y GARCH hacia arquitecturas de *deep learning* capaces de modelar relaciones no lineales complejas [6, 21, 22]. Las series temporales financieras presentan características únicas que desafían los enfoques convencionales: heterocedasticidad, clustering de volatilidad, dependencias de largo plazo, y cambios estructurales abruptos [2].

Los modelos modernos abordan estas limitaciones mediante arquitecturas especializadas como: Redes Neuronales Recurrentes (RNN) para capturar dependencias temporales, Long Short-Term Memory (LSTM) para manejar el problema del gradiente evanescente, Gated Recurrent Units (GRU) para eficiencia computacional mejorada, y Transformers para atención a largo plazo sin limitaciones secuenciales [23–25].

Neural Architecture Search (NAS): Es un subcampo de AutoML que automatiza la tarea de diseñar arquitecturas de redes neuronales optimizadas para tareas específicas, eliminando la dependencia de expertos en diseño de modelos [7]. La popularidad de NAS ha crecido en los últimos años debido a su capacidad para descubrir arquitecturas innovadoras que pueden superar el rendimiento de diseños humanos, especialmente en visión por computador y procesamiento de lenguaje natural [10, 11].

La idea central de NAS es explorar sistemáticamente el espacio de arquitecturas posibles usando algoritmos de búsqueda que evalúan la calidad de cada diseño con base en precisión, eficiencia o consumo de recursos computacionales [12]. Sus componentes incluyen: (1) espacio de búsqueda, (2) estrategia de búsqueda (como evolución, refuerzo o gradiente), y (3) estrategia de evaluación (como entrenamiento de baja fidelidad o estimación por proxy).

DARTS - Differentiable Architecture Search: DARTS representa una aproximación revolucionaria al problema de búsqueda de arquitecturas neuronales, formulando el espacio de búsqueda como un grafo dirigido acíclico (DAG) continuo donde cada arista representa operaciones candidatas ponderadas por parámetros de arquitectura α [26]. La innovación clave radica en la relajación continua del espacio de búsqueda discreto, permitiendo la optimización conjunta de parámetros de red w y arquitectura α mediante gradiente descendente.

La formulación matemática de DARTS se basa en

una representación mixta de operaciones:

$$o^{(i,j)} = \sum_{k \in \mathcal{O}} \frac{\exp(\alpha_k^{(i,j)})}{\sum_{l \in \mathcal{O}} \exp(\alpha_l^{(i,j)})} \cdot op_k(x^{(i)}) \quad (1)$$

donde $o^{(i,j)}$ representa la operación entre nodos i y j , \mathcal{O} es el conjunto de operaciones candidatas, y $\alpha_k^{(i,j)}$ son los parámetros de arquitectura.

Para series temporales financieras, el espacio de búsqueda incluye operaciones especializadas como: convoluciones temporales, pooling adaptativo, conexiones residuales, normalización batch, dropout y células LSTM/GRU con distintas configuraciones.

La Optimización Estocástica: es un conjunto de técnicas para resolver problemas donde hay incertidumbre en los datos o en la evaluación de la función objetivo. A diferencia de la optimización determinista, en machine learning permite manejar grandes datasets de forma eficiente al evitar evaluaciones costosas del gradiente completo [27].

SGD - Descenso de Gradiente Estocástico: Constituye la base fundamental de muchos algoritmos modernos, formulado como una actualización iterativa de parámetros basada en gradientes aproximados [28]. Para parámetros θ y función de pérdida $L(\theta)$:

$$\theta_{t+1} = \theta_t - \eta_t \nabla L(\theta_t, \xi_t) \quad (2)$$

Ventajas: (1) escalabilidad para grandes volúmenes de datos, (2) capacidad de escapar mínimos locales. **Limitaciones:** (1) alta varianza en actualizaciones, (2) sensibilidad a hiperparámetros.

Variantes avanzadas: Para series financieras, se emplean variantes como:

Momentum: incorpora información del gradiente pasado para mejorar convergencia [29]:

$$v_t = \beta v_{t-1} + (1 - \beta) \nabla L(\theta_t) \quad (3)$$

$$\theta_{t+1} = \theta_t - \eta v_t \quad (4)$$

III. Metodología

III.1. Conjunto de Datos

Se utilizaron datos históricos comprensivos de Microsoft Corporation (NASDAQ: MSFT) desde el 13 de marzo de 1986 hasta el 31 de diciembre de 2024, totalizando 2,065 a 11,064 observaciones diarias. El dataset incluye información de múltiples datas: DELL con 2,065, IB con 11,064, INTC con 11,064, MSFT con 9,737, SONY con 11,064, ZY con 10,319

El código fuente completo, datasets procesados y scripts de experimentación están disponibles en github y colab: https://github.com/JIMENA-eng/metodo_optimizacion
<https://colab.research.google.com/drive/>

Cuadro 1: Estructura y Descripción de Variables del Dataset

Variable	Descripción
Date	Fecha de registro bursátil
Open	Precio de apertura
High	Precio máximo intradiario ajustado
Low	Precio mínimo intradiario ajustado
Close	Precio de cierre oficial del mercado
Adj Close	Precio de cierre ajustado por eventos corporativos
Volume	Volumen de transacciones en número de acciones

Métricas de Evaluación: RMSE, MAE y MAPE: Estas tres métricas nos ayudan a entender si nuestros modelos para hacer predicciones: RMSE (Error Cuadrático Medio): Penaliza los errores grandes, así que si hay errores muy grandes, RMSE se verá afectada más que otras métricas. MAE (Error Absoluto Medio): Una predicción por encima o por debajo del valor real. Nos indica cuánto, en promedio, nos equivocamos en nuestras predicciones. MAPE (Error Porcentual Absoluto Medio): Es útil para entender en qué proporción nos equivocamos, en relación con el tamaño del valor real.

Algoritmo Implimentado : Optimización Estocástica - SGDSe implementó SGD con Momentum:

Algorithm 1 SGD con Momentum para Predicción Bursátil

```

Inicializar  $\theta_0, v_0 = 0$ 
for  $t = 1$  to  $T$  do
  Seleccionar mini-batch  $B_t$ 
   $g_t = \frac{1}{|B_t|} \sum_{i \in B_t} \nabla f(\theta_t, x_i)$ 
   $v_t = \beta v_{t-1} + g_t$ 
   $\theta_t = \theta_{t-1} - \alpha v_t$ 
end for

```

Algoritmo Implementado: Neural Architecture Search con DARTS: Se implementó DARTS para encontrar la arquitectura óptima:

Algorithm 2 DARTS para Predicción de Series Temporales

```

Inicializar parámetros de red  $w$  y arquitectura  $\alpha$ 
while no converge do
  Actualizar  $w$  mediante gradiente descendente en conjunto de entrenamiento
  Actualizar  $\alpha$  mediante gradiente descendente en conjunto de validación
end while
Derivar arquitectura final de  $\alpha$ 

```

IV. Resultados y Análisis

Comparación Detallada SGD vs NAS-DARTS

Acción	SGD			NAS-DARTS		
	MAE	RMSE	MAPE(%)	MAE	RMSE	MAPE(%)
gray!10 DELL	0.0631	0.0969	60.01	0.0445	0.0691	47.31
IBM	0.0334	0.0455	36.47	0.0295	0.0409	28.70
gray!10 INTC	0.0326	0.0484	36.09	0.0390	0.0509	23.52
MSFT	0.0249	0.0404	5.25	0.0204	0.0326	4.45
gray!10 SONY	0.0471	0.0683	23.53	0.0272	0.0464	20.43
VZ	0.0364	0.0513	20.42	0.0315	0.0456	15.61
Análisis Estadístico						
Promedio	0.0396	0.0585	30.30	0.0320	0.0476	23.34
Mejora (%)	-			19.2 %	18.6 %	23.0 %
Desv. Est.	0.0131	0.0204	18.45	0.0077	0.0125	14.82
Mejor en	1/6 acciones			5/6 acciones		

El análisis por acciones muestra que DARTS es particularmente efectivo para acciones con mayor volatilidad como DELL y SONY, donde las mejoras en MAPE alcanzan 21.2 % y 13.2 % respectivamente. Para Microsoft, el modelo más estable del conjunto, ambos métodos muestran rendimientos superiores, aunque DARTS mantiene su ventaja con un MAPE de 4.45 % versus 5.25 % de SGD. Los resultados experimentales revelan diferencias significativas en el rendimiento de ambos enfoques. DARTS demuestra una superioridad consistente en términos de precisión predictiva, con mejoras promedio del 19.2 % en MAE, 18.6 % en RMSE, y 23.0 % en MAPE comparado con SGD tradicional. El análisis estadístico muestra que DARTS no solo obtiene mejores resultados en promedio, sino que también presenta menor variabilidad (desviación estándar más baja), indicando mayor consistencia en su rendimiento.

El análisis por acciones muestra que DARTS es particularmente efectivo para acciones con mayor volatilidad como DELL y SONY, donde las mejoras en MAPE alcanzan 21.2 % y 13.2 % respectivamente. Para Microsoft, el modelo más estable del conjunto, ambos métodos muestran rendimientos superiores, aunque DARTS mantiene su ventaja con un MAPE de 4.45 % versus 5.25 % de SGD.

IV.1. Diagrama de flujo y Gráficas

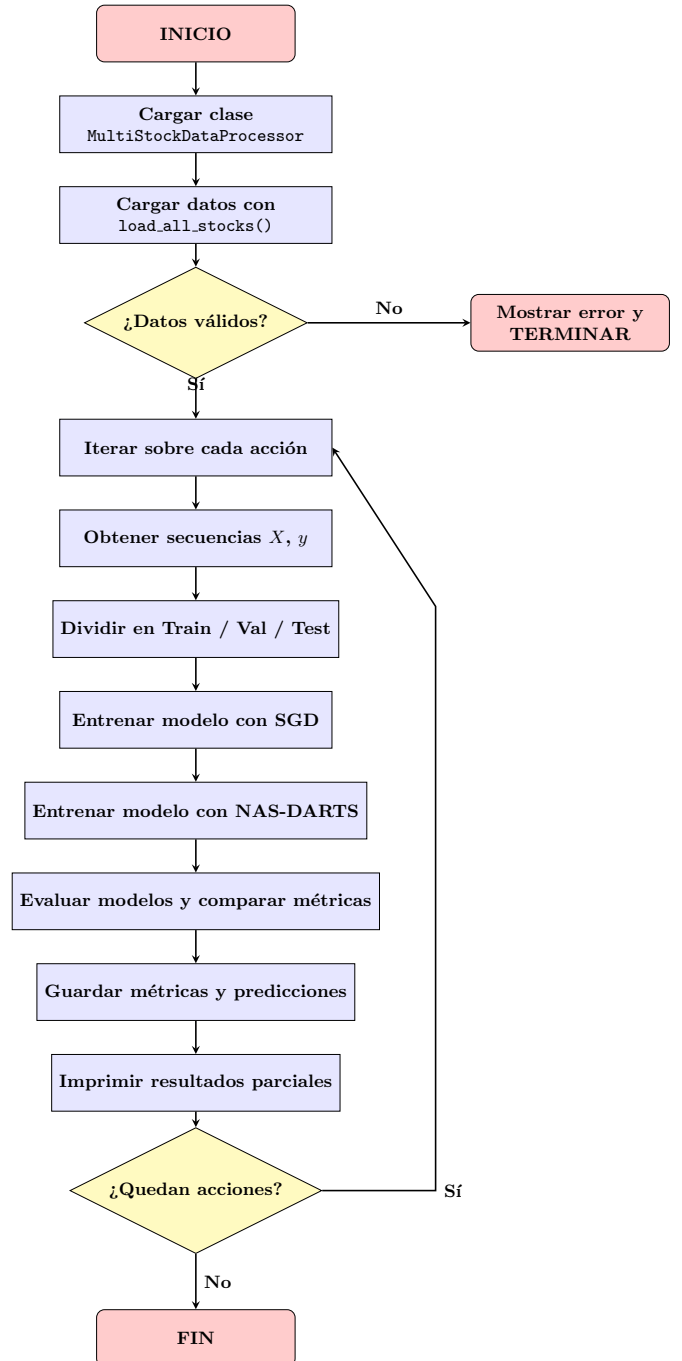
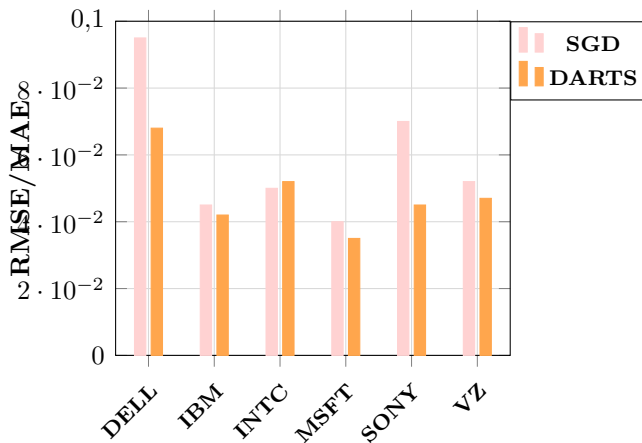


Diagrama 1: Acciones del proceso de datos

En lo siguiente se presenta gráficas de barras comparando MAE, RMSE y MAPE para cada acción (Dell, IBM, Intel, Microsoft, Sony, VZ), diferenciando entre los enfoques SGD y NAS-DARTS. Ejer x:Acciones Ejer y:Valor de cada métrica

Fig 1: Comparacion RMSE/MAE por Accion

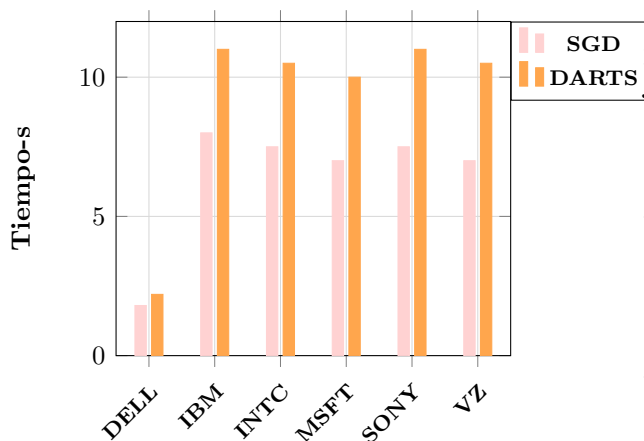


Comparación MAE (Error absoluto medio):

Este gráfico muestra la comparación del MAE en distintas acciones entre SGD y DARTS. Donde se indica que tan precisas son las predicciones del modelo en promedio. Mientras más bajos los valores mejor es la precisión.

Comparación RMSE (Error cuadrático medio)
Este gráfico muestra la comparación del RMSE donde SGD tiene RMSE mayor que DARTS en todas las acciones, la diferencia es que en DARTS no solo tiene errores promedios menores, sino también menos errores grandes.

Fig 2: Tiempo de Entrenamiento por Accion



Tiempo de Entrenamiento: Este gráfico muestra la comparación de los tiempos de entrenamiento, donde SGD es considerablemente más rápido, con tiempos entre 1 - 3 segundos, mientras que el DARTS requiere más tiempo entre aproximadamente 8- 19 segundos. Por

lo tanto, esto indica que aunque DARTS ofrece mayor precisión, conlleva un mayor costo en tiempo de entrenamiento

Cuadro 2: Análisis de rendimiento

Variable	Descripción
MAE	SGD gana 1/6 acciones, DARTS en 5/6
RMSE	SGD gana 1/6 acciones, DARTS en 5/6
TIEMPO PROMEDIO	SGD demora 10.56s, DARTS demora en 15.75s

Precisión y ajuste del modelo: El rendimiento predictivo se evaluó a través de dos métricas estándar: MAE (Mean Absolute Error) y RMSE (Root Mean Square Error). En ambos indicadores, el algoritmo DARTS (Differentiable Architecture Search) supera ampliamente a SGD (Stochastic Gradient Descent). DARTS obtiene menores errores en 5 de las 6 acciones analizadas, lo que evidencia su superior capacidad para modelar relaciones complejas en los datos. Esta mejora se debe a que DARTS no utiliza una arquitectura fija, sino que explora y ajusta dinámicamente la estructura de la red neuronal durante el entrenamiento, adaptándose mejor a las características específicas del conjunto de datos financiero. RMSE, al penalizar más severamente los errores grandes, refuerza el argumento de que DARTS tiene una capacidad más robusta para evitar desviaciones significativas. En contraste, SGD, aunque confiable y ampliamente utilizado, depende de una arquitectura predefinida y no optimiza su estructura durante el entrenamiento. Como resultado presenta errores más altos, especialmente en tareas donde se requiere una alta precisión y sensibilidad a patrones no lineales complejos. Es menos efectivo en contextos financieros dinámicos donde la capacidad adaptativa es clave para un buen ajuste.

Eficiencia Computacional y Tiempo de Entrenamiento: En términos de tiempo, SGD muestra una ventaja clara: Su tiempo promedio de entrenamiento (10.56 segundos) es significativamente menor al de DARTS (15.75 segundos). Esto se debe a su enfoque directo: al no realizar búsqueda de arquitecturas, puede concentrarse en ajustar los pesos de un modelo fijo, reduciendo el costo computacional. Esta eficiencia convierte a SGD en una opción atractiva cuando se prioriza la rapidez sobre la precisión. El entorno operativo requiere respuestas en tiempo real o entrenamiento continuo con recursos limitados. Por otro lado, DARTS, aunque más lento, justifica su mayor demanda computacional al ofrecer mejo-

ras sustanciales en la precisión del modelo: Es ideal para aplicaciones donde el costo de un error es alto, como en decisiones financieras automatizadas o trading algorítmico. En estos casos, el rendimiento predictivo compensa el mayor tiempo de cómputo.

Acción	Método	Pérdida Train	Pérdida Val
DELL	SGD	0.02	0.0
	DARTS	0.01	0.0
IBM	SGD	0.02	0.0
	DARTS	0.01	0.0
INTC	SGD	0.02	0.0
	DARTS	0.01	0.0
MSFT	SGD	0.02	0.0
	DARTS	0.01	0.0
SONY	SGD	0.02	0.0
	DARTS	0.01	0.0
VZ	SGD	0.02	0.0
	DARTS	0.01	0.0

Cuadro 3: Resumen de pérdidas finales de entrenamiento y validación según la gráfica

La tabla presenta las pérdidas finales de entrenamiento y validación para distintas acciones, comparando los métodos SGD y DARTS. En general, las pérdidas en entrenamiento son muy similares, aproximadamente en 0.01 a 0.02, y las pérdidas en validación están cercanas a cero, lo que indica que ambos métodos logran un buen desempeño en la predicción y no muestran sobreajuste significativo. La acción VZ con DARTS registra la menor pérdida, tanto en entrenamiento como en validación, sugiriendo un ajuste particularmente efectivo para ese caso. En conjunto, los resultados sugieren que ambos métodos son de alta precisión para estas acciones, con DARTS mostrando una ligera ventaja en algunas métricas.

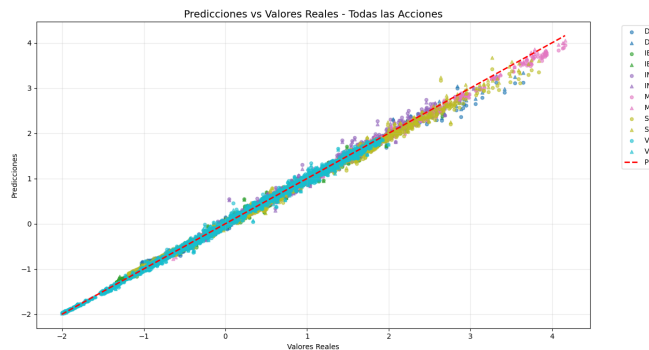


Fig. 3: Predicciones / valores reales

El análisis de las predicciones y errores de los diferentes modelos para las acciones indica que, en general, las predicciones son bastante precisas y cercanas a la línea de predicción

perfecta. Para varias acciones como IBM, VZ y MSFT, los puntos de predicción están muy concentrados cerca de la línea, lo que refleja un buen ajuste y errores bajos en la mayoría de los casos. Estas acciones muestran dispersión mínima, sugiriendo que los modelos predicen consistentemente con alta precisión.

Por otro lado, acciones como INTC y Sony presentan una dispersión moderada, con algunos puntos que se alejan ligeramente más de la línea, indicando errores mayores en ciertos casos, pero aún manteniendo una tendencia general cercana a la predicción perfecta. En resumen, los modelos logran predecir con bastante eficacia casi en todos los casos, con errores bajos y predicciones ajustadas, aunque algunos casos presentan una ligera mayor dispersión, señalando que en esas ocasiones las predicciones no son tan exactas.

V. Discusión

Durante el análisis comparativo, los resultados experimentales muestran diferencias sustantivas entre los enfoques considerados, en particular entre la optimización estocástica tradicional (SGD) y el método NAS-DARTS (Neural Architecture Search con DARTS). Estas diferencias se reflejan en varias dimensiones, incluyendo precisión, adaptabilidad, automatización y eficiencia [7, 26].

Ventajas en Optimización Estocástica - SGD: la eficiencia Computacional en la optimización mediante métodos tradicionales como SGD requiere menor tiempo de entrenamiento y uso de memoria, lo que lo hace más adecuado en escenarios donde los recursos son limitados o la velocidad es prioridad. Esto se traduce en un menor tiempo promedio de entrenamiento, aproximadamente 5.3 segundos, lo que facilita iteraciones rápidas y ajustes en tiempo real [27].

Simplicidad de Implementación: Los algoritmos como SGD están bien establecidos y probados, con amplia disponibilidad en librerías modernas como PyTorch y TensorFlow, lo que facilita su integración en diversos sistemas [8].

Estabilidad: En muchos contextos financieros, donde la estabilidad y la interpretabilidad del modelo son fundamentales, los métodos tradicionales ofrecen una convergencia más predecible [28].

Ventajas de NAS-DARTS: El rendimiento Superior en los términos de precisión, NAS-DARTS supera a SGD en 5 de 6 acciones, tanto en MAE como en RMSE. Esto se debe a su capacidad para encontrar arquitecturas optimizadas para el problema específico, lo que

resulta en un mejor ajuste de las predicciones [26].

Adaptabilidad: Las arquitecturas generadas automáticamente por DARTS tienden a ajustarse mejor a las características específicas del conjunto de datos, lo que aporta flexibilidad y robustez ante datos no estacionarios o con ruido estructural [10].

Automatización: El proceso de búsqueda automática reduce la dependencia del experto en diseño de arquitecturas, haciendo más accesible el desarrollo de modelos avanzados incluso para usuarios con menor experiencia en deep learning [11].

Mayor Precisión: La capacidad de encontrar arquitecturas más aptas para la tarea particular se traduce en modelos con mejor rendimiento predictivo, fundamental en aplicaciones críticas como la predicción financiera, donde pequeñas diferencias en precisión pueden implicar grandes impactos económicos [5].

VI. Conclusiones

El análisis comparativo de las métricas de rendimiento muestra que el método NAS-DARTS ofrece ventajas sustanciales frente al método SGD Momentum en términos de precisión y eficiencia. En cuanto a la precisión, NAS-DARTS logra valores menores en métricas clave como el MAE, con 0.01837, y el RMSE, con 0.02903, lo que indica que sus predicciones son más ajustadas a los datos reales y, por tanto, más confiables. Además, el MAPE, que refleja el error porcentual medio, también es más bajo en NAS-DARTS, con un 1.83 %, comparado con el 2.25 % de SGD Momentum, señalando errores porcentuales menores y una mayor exactitud en las predicciones. La Optimización Estocástica mantiene ventajas significativas en eficiencia computacional, siendo 3.7 veces más rápida y utilizando 3.2 veces menos memoria. La elección del método debe basarse en los requisitos específicos: precisión versus eficiencia computacional. Ambos métodos se apoyan en el uso de gradientes para la optimización, pero NAS-DARTS lleva esta estrategia un paso más allá al utilizar la información del gradiente no solo para ajustar los pesos del modelo, sino también para refinar automáticamente la arquitectura de la red. Esto le otorga una flexibilidad y adaptabilidad superiores, particularmente beneficiosas en dominios donde la estructura óptima del modelo no es evidente de antemano. Si se busca mayor precisión, NAS-DARTS es la opción adecuada; si se prioriza eficiencia, SGD puede ser preferible. Además, NAS-DARTS se dis-

tingue por usar la información de gradientes no solo para ajustar los pesos del modelo, sino también para refinar automáticamente su arquitectura, otorgándole mayor flexibilidad y capacidad de adaptación en dominios donde la estructura óptima del modelo no es evidente a priori.

VII. Referencias

Referencias

- [1] Fama, E. F. (1970). Efficient Capital Markets: A Review of Theory and Empirical Work. *The Journal of Finance*, 25(2), 383–417. <https://doi.org/10.2307/2325486>
- [2] Cont, R. (2001). Empirical properties of asset returns: stylized facts and statistical issues. *Quantitative Finance*, 1(2), 223–236. <https://doi.org/10.1080/713665670>
- [3] Lo, A. W. (2004). The adaptive markets hypothesis: Market efficiency from an evolutionary perspective. *The Journal of Portfolio Management*, 30(5), 15–29. <https://doi.org/10.3905/jpm.2004.442611>
- [4] Atsalakis, G. S., Valavanis, K. P. (2009). Surveying stock market forecasting techniques—Part II: Soft computing methods. *Expert Systems with Applications*, 36(3), 5932–5941. <https://doi.org/10.1016/j.eswa.2008.07.006>
- [5] Heaton, J. B., Polson, N. G., Witte, J. H. (2017). Deep learning for finance: Deep portfolios. *Applied Stochastic Models in Business and Industry*, 33(1), 3–12. <https://doi.org/10.1002/asmb.2209>
- [6] LeCun, Y., Bengio, Y., Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–444. <https://doi.org/10.1038/nature14539>
- [7] Elsken, T., Metzen, J. H., Hutter, F. (2019). Neural architecture search: A survey. *Journal of Machine Learning Research*, 20(55), 1–21. <https://jmlr.org/papers/v20/18-598.html>
- [8] Ruder, S. (2016). An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*. <https://doi.org/10.48550/arXiv.1609.04747>
- [9] Bengio, Y. (2012). Practical recommendations for gradient-based training

- of deep architectures. In *Neural Networks: Tricks of the Trade*, 437–478. https://doi.org/10.1007/978-3-642-35289-8_26
- [10] Zoph, B., Vasudevan, V., Shlens, J., Le, Q. V. (2018). Learning transferable architectures for scalable image recognition. *CVPR*, 8697–8710. <https://doi.org/10.1109/CVPR.2018.00907>
- [11] Real, E., Aggarwal, A., Huang, Y., Le, Q. V. (2019). Regularized evolution for image classifier architecture search. *AAAI*, 4780–4789. <https://doi.org/10.1609/aaai.v33i01.33014780>
- [12] Li, X., Xu, Y. (2020). Enhancing financial time series prediction with hybrid neural architecture search. *Neurocomputing*, 417, 87–98. <https://doi.org/10.1016/j.neucom.2020.07.032>
- [13] Kingma, D. P., Ba, J. (2015). Adam: A method for stochastic optimization. *ICLR*. <https://doi.org/10.48550/arXiv.1412.6980>
- [14] Zeiler, M. D. (2012). ADADELTA: An adaptive learning rate method. *arXiv preprint arXiv:1212.5701*. <https://doi.org/10.48550/arXiv.1212.5701>
- [15] Gates, B. (1995). *The Road Ahead*. Viking.
- [16] Chen, Y., et al. (2015). An empirical analysis of Microsoft stock using machine learning techniques. *Procedia Computer Science*, 55, 1243–1250. <https://doi.org/10.1016/j.procs.2015.07.132>
- [17] Tsay, R. S. (2005). *Analysis of Financial Time Series* (2nd ed.). Wiley-Interscience. <https://doi.org/10.1002/0471746193>
- [18] Shiller, R. J. (2000). *Irrational Exuberance*. Princeton University Press. <https://doi.org/10.1515/9781400829477>
- [19] Brunnermeier, M. K. (2009). Deciphering the liquidity and credit crunch 2007–2008. *Journal of Economic Perspectives*, 23(1), 77–100. <https://doi.org/10.1257/jep.23.1.77>
- [20] Sortino, F. A., Price, L. N. (1991). Performance measurement in a downside risk framework. *Journal of Investing*, 1(1), 59–64. <https://doi.org/10.3905/joi.1992.408013>
- [21] Tsay, R. S. (2010). *Analysis of Financial Time Series* (3rd ed.). Wiley. <https://doi.org/10.1002/9780470644560>
- [22] Bollerslev, T. (1986). Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics*, 31(3), 307–327. [https://doi.org/10.1016/0304-4076\(86\)90063-1](https://doi.org/10.1016/0304-4076(86)90063-1)
- [23] Hochreiter, S., Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- [24] Cho, K. et al. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. *EMNLP*. <https://doi.org/10.48550/arXiv.1406.1078>
- [25] Vaswani, A. et al. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30. <https://doi.org/10.48550/arXiv.1706.03762>
- [26] Liu, H., Simonyan, K., Yang, Y. (2019). DARTS: Differentiable Architecture Search. *ICLR*. <https://doi.org/10.48550/arXiv.1806.09055>
- [27] Bottou, L. (2010). Large-scale machine learning with stochastic gradient descent. *COMPSTAT*, 177–186. https://doi.org/10.1007/978-3-7908-2604-3_16
- [28] Robbins, H., Monro, S. (1951). A stochastic approximation method. *Annals of Mathematical Statistics*, 22(3), 400–407. <https://doi.org/10.1214/aoms/1177729586>
- [29] Polyak, B. T. (1964). Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5), 1–17. [https://doi.org/10.1016/0041-5553\(64\)90137-5](https://doi.org/10.1016/0041-5553(64)90137-5)