

DETECCIÓN DE ANOMALÍAS EN SERIES DE TIEMPO FINANCIERAS USANDO APRENDIZAJE PROFUNDO

César Seijas¹, Egilda Pérez, Rafael Pacheco, Ricardo Villegas, Teodoro García, Sergio Villasana

Resumen

La detección de anomalías en series de tiempo es importante porque permite identificar patrones y tendencias que no son evidentes a simple vista y que pueden ser indicativos de problemas o riesgos, y que, en el caso financiero, puede representar serios desajustes económicos. La detección de anomalías se realiza usualmente mediante modelos estadísticos que capturan tendencias, estacionalidad y niveles en datos de series de tiempo. Cuando los datos nuevos difieren demasiado del modelo, se indica una anomalía o una falla del modelo. En este artículo, se describe un sistema automatizado para detección de anomalías en series de tiempo, usando algoritmos basados en aprendizaje profundo. Para el desarrollo de la investigación se usaron datos del comportamiento del índice de precios al consumidor, registrados mensualmente en nuestro país entre los años 1950 hasta agosto del 2022 por el Banco Central de Venezuela. Para la construcción del modelo de aprendizaje profundo se implementó una versión unidimensional del autocodificador U-Net. Se lograron resultados satisfactorios, alcanzando una precisión de 94,6 % en la detección de las anomalías.

Palabras clave: detección de anomalías, series de tiempo económicas o financieras, aprendizaje profundo.

DETECTION OF ANOMALIES IN FINANCIAL TIME SERIES USING DEEP LEARNING

Abstrac

The detection of anomalies in time series is important because it makes it possible to identify patterns and trends that are not obvious to the naked eye and that may be indicative of problems or risks, and which, in the financial case, may represent serious economic imbalances. Anomaly detection is usually done using statistical models that capture trends, seasonality, and levels in time series data. When the new data differs too much from the model, an anomaly or model failure is indicated. In this article, an automated system for detecting anomalies in time series is described, using algorithms based on deep learning. For the development of the research, data on the behavior of the consumer price index were used, registered monthly in our country between the years 1950 and August 2022 by the Central Bank of Venezuela. For the construction of the deep learning model, a one-dimensional version of the U-Net autoencoder was implemented. Satisfactory results were achieved, reaching an accuracy of 94.6% in the detection of anomalies.

Keywords: anomaly detection, economic or financial time series, deep learning.

¹ Docentes de la Universidad de Carabobo. cseijas@uc.edu.ve;
egiperez@uc.edu.ve

Introducción

Una serie de tiempo (ST) es un conjunto de observaciones de un proceso realizadas de modo secuencial en el tiempo. Su análisis es de gran importancia en una amplia gama de temas de investigación en ingeniería, medicina, finanzas y otros campos de la ciencia. El objetivo fundamental que se persigue con el análisis de ST es la predicción de valores futuros basados en mediciones observadas previamente.

Entre otros objetivos alcanzables se pueden mencionar la descripción, explicación, control o detección de anomalías. El caso específico de detección de anomalías se refiere a la búsqueda de patrones con comportamiento inusual, los cuales pueden ser interpretados como acciones no válidas o anómalas sobre los datos.

El análisis estadístico clásico de ST consiste en la estimación de métricas estadísticas significativas de patrones a partir de tres rasgos que típicamente exhiben: no estacionariedad, autocorrelación y estacionalidad. Estos rasgos se conjugan en el modelo por Box-Jenkins, ARIMA (Schaffer et al., 2021).

En el caso de detección de anomalías, el uso de técnicas estadísticas tradicionales ha producido resultados satisfactorios (Borges et al., 2022; Schaffer et al., 2021; Schmidl et al., 2021). Sin embargo, en los últimos años se han reportado resultados de mayor relevancia mediante el uso de Inteligencia Artificial (IA) y en particular de técnicas de aprendizaje profundo (*Deep Learning*, DL) (Bengio et al., 2016). Aplicada a sistemas financieros, la detección de anomalías con DL ha permitido la identificación y prevención de actividades maliciosas como fraude e intrusiones (Bakumenco et al., 2022), entre otras actividades irregulares.

Este artículo se enfoca en la detección de anomalías en series de tiempo económicas o financieras (STF) usando DL; para la implementación de un

sistema automático de detección de anomalías en IPC, se usó el algoritmo de DL conocido como autocodificador (Hinton, 2016).

Como STF de estudio, se usó el comportamiento del índice de precios al consumidor, IPC, registrado mensualmente en nuestro país entre los años 1950 hasta agosto del 2022 (Banco Central de Venezuela, 2022). El IPC es un indicador estadístico que cuantifica la evolución de los precios de una canasta representativa del consumo familiar durante un período determinado; para su cálculo, se selecciona un año de referencia o base (su valor se fija en 100). La canasta es un conjunto de bienes y servicios típicos que consume un hogar y la importancia relativa que tiene cada rubro en el gasto de consumo familiar establece la estructura de pesos del IPC.

En la figura 1 se presenta el gráfico de la STF a analizar: IPC 1950-2022. El gráfico, IPC 1950-2022, muestra el registro secuencial, expresado porcentualmente, de la variación del IPC mensual, acontecida en nuestro país, entre los meses de enero de 1950 hasta agosto del 2022, tomando como año de referencia (año base de IPC = 100), 2007. Obsérvese en dicho gráfico, la variación extrema del IPC a partir de aproximadamente el año 2017 cuando se alcanza un pico del orden del 200%.

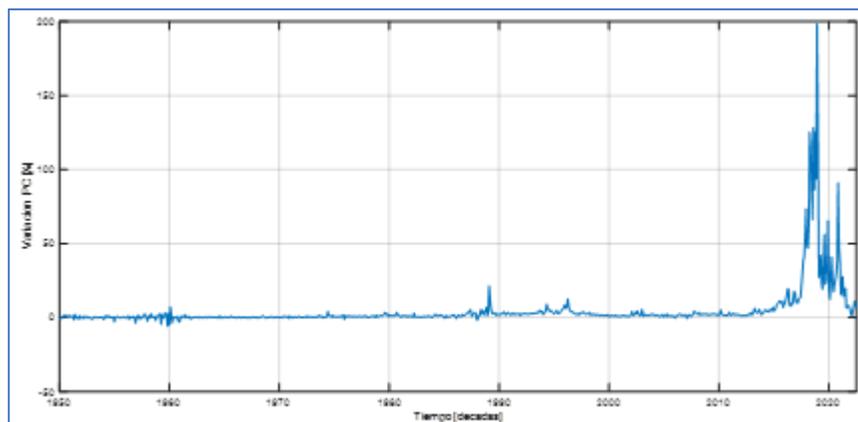


Figura 1. Variación del IPC mensual, expresado porcentualmente, acontecido en nuestro país entre los meses de enero de 1950 hasta agosto del 2022. Elaboración propia.

La estructura del presente artículo es la siguiente: esta primera sección, fue una explicación general del objetivo central del trabajo y la metodología usada; en la segunda sección se presentan investigaciones previas en el tema; la tercera sección, se dedica a fundamentación teórica del algoritmo; luego, la sección 4 detalla la implementación del sistema desarrollado, mientras que la sección 5 se ocupa de la parte experimental y análisis de resultados, finalmente presentamos las conclusiones y referencias bibliográficas.

Trabajos relacionados

El tema de la detección de anomalías en series de tiempo (ST) ha sido objeto de estudio de muchos investigadores. Crépey et al. (2022) desarrollaron un sistema detector de anomalías en series de tiempo económicas o financieras (STF), usando un modelo híbrido de redes neuronales artificiales (RNA) y la técnica de análisis de componentes principales (PCA); el uso de PCA como extractor de rasgos permite la reducción de la dimensionalidad. Zhou et al. (2022) usan un autocodificador contrastivo para detección de anomalías en series de tiempo multivariadas. Wei et al. (2022) usan un autocodificador basado en LSTM (siglas en inglés para la RNA conocida como: "*Long Short-Term Memory*") para detectar anomalías en series de tiempo de calidad del aire en ambientes interiores.

En el tema de pronósticos de STF, Korczak y Hernes (2017) desarrollaron métodos para pronóstico de STF usando DL en un sistema de negocio de acciones multi-agente, también conocido como "*A-trader*". Similarmente, Navon et al. (2017) usan una RNA de DL para predecir tendencias de las acciones en el mercado NASDAQ.

Por otra parte, el análisis del índice de precios al consumidor (IPC) como STF indicadora del status económico de un país, es un problema de interés en estudios académicos de las ciencias económicas. Tal como lo refleja el trabajo por López-Ávila et al. (2019), donde pronostican el IPC mexicano usando el

algoritmo FS-EPNet (siglas en inglés para “*Evolutionary Network Programing*”) y comparan sus resultados con modelos estadísticos tales como ARCH, ARIMA (Schaffer et al., 2021) y otros. Caicedo (2018) usa en su tesis de grado un modelo VAR para el pronóstico del IPC colombiano. Finalmente, se comenta un estudio por Quilis y Frutos (1999), donde los autores afirman que el análisis habitual de las condiciones inflacionarias de la economía española descansa, principalmente, en el IPC; en este trabajo, los autores usan un modelo VARMA (siglas en inglés para: “*Vectorial Autorregresive Moving Averaging*”) para la estimación del indicador.

Fundamentos teóricos

El sistema de detección de anomalías utiliza DL y consta de cuatro bloques en cascada. El primer bloque se encarga del preprocesamiento de la STF, dividiendo la secuencia de entrada en una subsecuencia de entrenamiento (libre de anomalías) y una subsecuencia de prueba (con tramos anómalos). De cada subsecuencia (entrenamiento y prueba) se extraen subsecuencias menores, pero con longitud suficiente para garantizar su reconstrucción con modelos autorregresivos (AR(p)). Los subconjuntos de entrenamiento y prueba son normalizados para el entrenamiento óptimo del autocodificador. El siguiente bloque es el autocodificador, que extrae rasgos relevantes de cada secuencia de entrada.

Autocodificador

Un autocodificador es una red neuronal artificial no recurrente (Cholet, 2018), entrenada de modo no supervisado para replicar a la salida el mismo vector de entrada (Yan et al., 2020). Fue concebida inicialmente para reducción de dimensionalidad (Hinton, 2006), pero internamente tiene una capa escondida intermedia que describe un código de dimensión reducida respecto al vector de entrada utilizado para representar la entrada (representación

latente). La arquitectura de un autocodificador consta de una ruta de contracción para capturar contexto (rasgos) y una ruta de expansión simétrica que permite una localización precisa (figura 2).

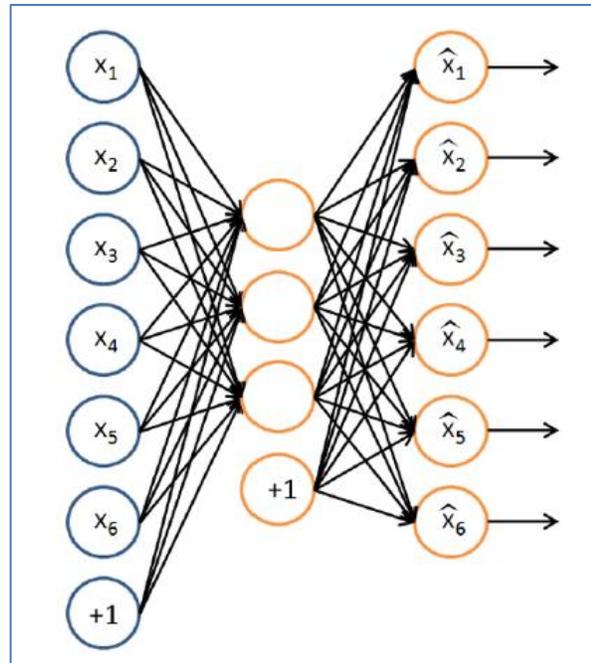


Figura 2. Arquitectura de un Autocodificador (Ronneberger et al. 2015).

Tal como se puede inferir de la arquitectura mostrada, se puede considerar que el autocodificador consta de dos etapas: una etapa codificadora $h = \sigma(x)$ y una decodificadora, que produce una reconstrucción $x' = \sigma'(z)$. En el caso de procesar secuencias temporales, el vector a reconstruir corresponde a esta secuencia. Matemáticamente, esto es:

La etapa codificadora toma la entrada $x \in \mathcal{R}^d = X$ y la proyecta a la capa escondida $h \in \mathcal{R}^p = F$, de tal manera que:

$$h = \sigma(W \cdot x + b) \quad (\text{Ec. 1})$$

Donde:

\mathbf{h} representa la variable latente,

σ es una función de activación

\mathbf{W} corresponde a la matriz de pesos y

\mathbf{b} es el vector de polarización.

La etapa decodificadora toma la representación latente para reconstruir la entrada, esto es:

$$\mathbf{x}' = \sigma'(\mathbf{W}' \cdot \mathbf{h} + \mathbf{b}') \quad (\text{Ec. 2})$$

La réplica de la entrada, o vector reconstruido, \mathbf{x}' se obtiene minimizando una función de pérdidas de reconstrucción, $L(\mathbf{x}, \mathbf{x}')$ dada por:

$$L(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|^2 = \|\mathbf{x} - \sigma'(\mathbf{W}' \cdot \sigma(\mathbf{W} \cdot \mathbf{x} + \mathbf{b}) + \mathbf{b}')\|^2 \quad (\text{Ec. 3})$$

La función global de error es:

$$J_{AE}(\Theta) = \sum_x L(\mathbf{x}, \mathbf{x}') \quad (\text{Ec. 4})$$

donde:

$\Theta = (\mathbf{W}, \mathbf{W}', \mathbf{b}, \mathbf{b}')^\tau$ es el vector de parámetros entrenables.

La ecuación de entrenamiento (actualización) del vector de parámetros Θ es:

$$\Theta_{i+1} := \Theta_i - \alpha \cdot \frac{\partial L_{AE}(\Theta_i)}{\partial \Theta_i} \quad (\text{Ec. 5})$$

donde $\alpha \geq 0$ es la constante de aprendizaje (Chollet, 2018; Yen, 2020).

A continuación, en la siguiente figura 3 se puede observar la representación gráfica de la Arquitectura de U-Net:

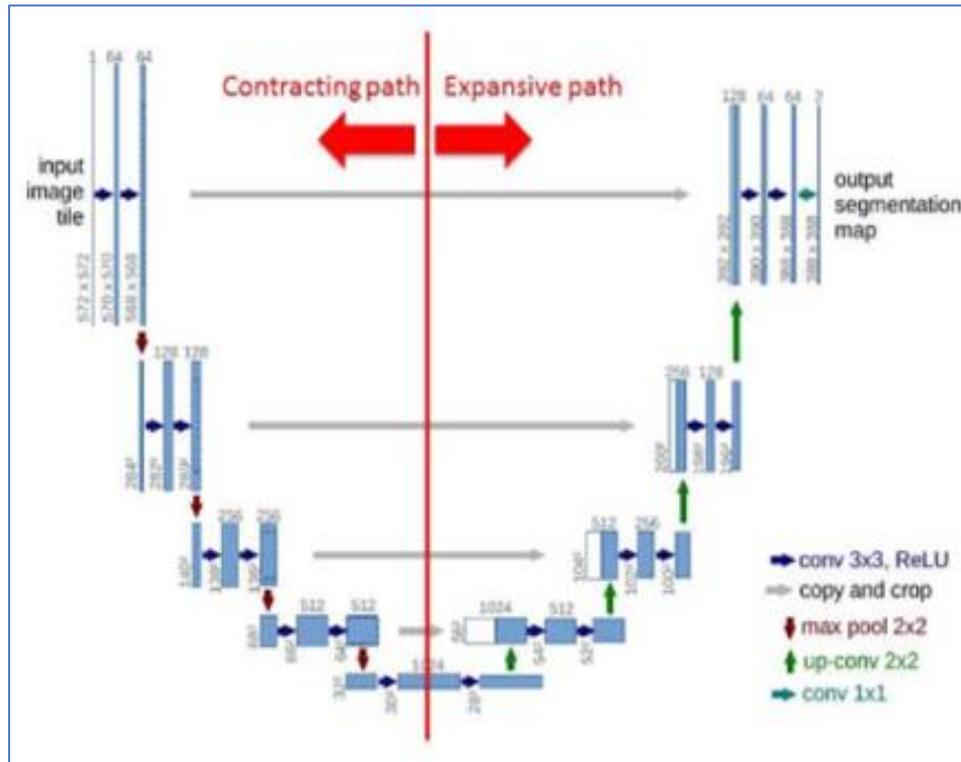


Figura 3. Arquitectura de U-Net. (Ronnerberger et al., 2015)

Metodología

En esta sección se describen los materiales y métodos usados en el experimento de detección de anomalías en la base de datos, STF: IPC 1950-2022. Esto incluye: pre-procesamiento, plataformas, librerías y códigos empleados en la implementación del modelo de detección.

Base de datos

El conjunto de datos corresponde a la STF: IPC 1950-2022, la cual contiene el registro secuencial de la variación del IPC porcentual mensual, nacional, desde enero de 1950 hasta agosto del 2022, tomando como año de referencia (año base de IPC = 100), 2007. El total de registros representa un conjunto de 864 muestras, equiespaciadas mensualmente, univariadas; del

cual se seleccionó, aleatoriamente, un subconjunto de 750 muestras para entrenamiento y se reservó el resto para evaluación del modelo, una relación 90%-10%, para entrenamiento/pruebas.

De los conjuntos de entrenamiento y prueba, se extrajeron 726 secuencias univariadas de 288 muestras cada una, repartidas en 463 secuencias de entrenamiento y 163 para prueba. En términos tensoriales (Cholet, 2018), el conjunto de entrenamiento es un tensor de dimensiones [463, 288, 1] y el conjunto de pruebas uno de dimensiones [163, 288, 1]. Es oportuno indicar que la región de entrenamiento corresponde al segmento que visualmente mostró mayor estacionariedad (años previo a 2010). En las figuras 4 y 5 se representan segmentos de los conjuntos de entrenamiento y prueba.

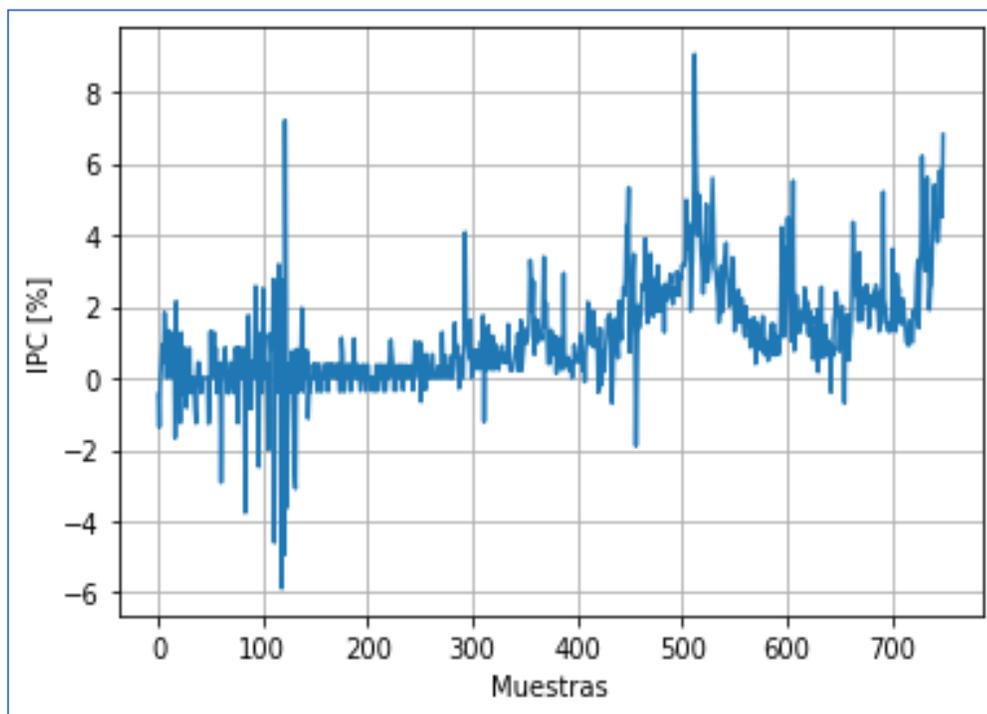


Figura 4. Segmento de IPC 1950-2022 de entrenamiento. Elaboración propia.

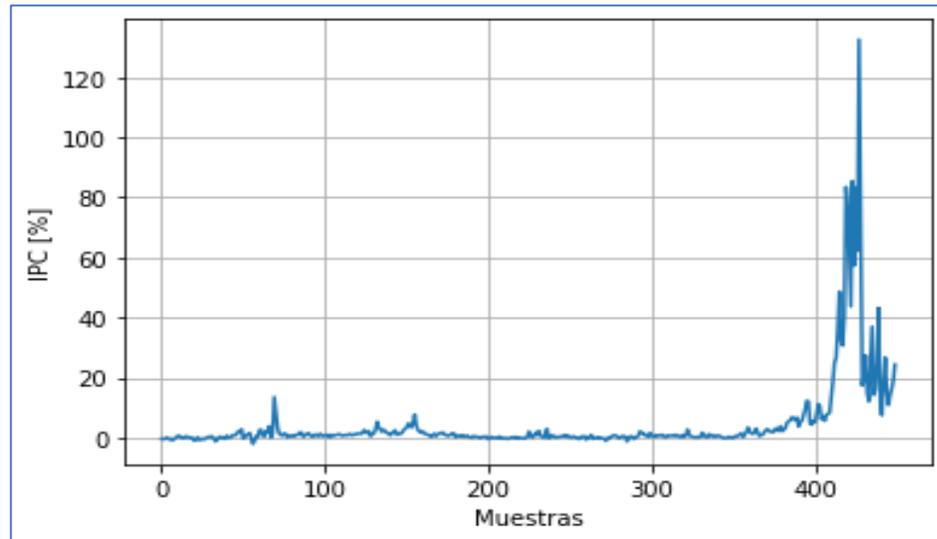


Figura 5. Segmento de IPC 1950-2022 de prueba. Elaboración propia.

Plataforma de codificación y estructura y pre-procesamiento de datos

El modelo autocodificador se implementó en las plataformas *Tensorflow* y *Keras*. Como autocodificador se usó una red *U-Net* modificada (Ronneberger et al, 2015) para procesamiento univariado, porque la original procesa imágenes (datos bidimensionales o 2D). La arquitectura de la *U-Net* implementada sustituye las capas convolucionales 2D (*Conv2D*), por unidimensionales *Conv1D* en la etapa de compresión y *Conv1DTranspose* en la etapa de expansión, con capas *Dropout* (Chollet, 2018) entre capas convolucionales para propósitos de regularización.

El procedimiento empleado para la detección de anomalías en la serie de tiempo es el propuesto por Pavithrasv (2020), el cual sigue los siguientes pasos:

1. Determine error absoluto medio (*Mean absolute error*, MAE) en el conjunto de entrenamiento.
2. Encuentre el valor máximo del MAE, este valor se establece como

umbral para detección de anomalías.

3. Si la pérdida de reconstrucción de una muestra es mayor que el valor umbral, la muestra es etiquetada como anomalía.

Resultados

En esta sección se presentan los hallazgos obtenidos durante el entrenamiento, validación y prueba del modelo. Los hiperparámetros de entrenamiento seleccionados en este caso de estudio fueron:

- a) Número de épocas: 50.
- b) Tamaño del lote de muestras por paso de entrenamiento (*batch_size*): 128 muestras.
- c) Tamaño del conjunto de validación: 10% del conjunto de datos total.

La figura 6, muestra la distribución del MAE en el segmento de entrenamiento, allí puede observarse el umbral del error de reconstrucción de aproximadamente 13%.

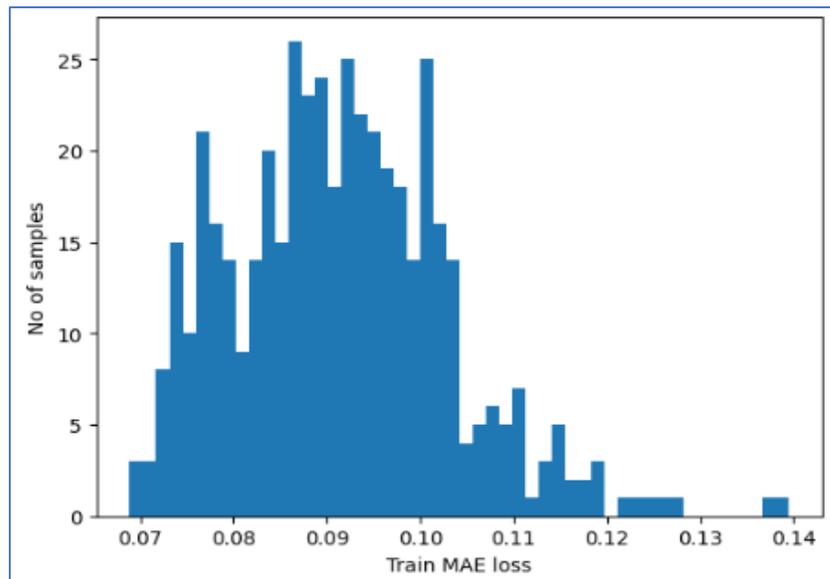


Figura 6. Distribución del MAE en el segmento de entrenamiento. Elaboración propia.

Por su parte, la figura 7 permite comparar la similitud entre la señal de entrenamiento y la reconstruida en el segmento de entrenamiento, mientras que la figura 8 muestra la distribución del MAE para el segmento de prueba.

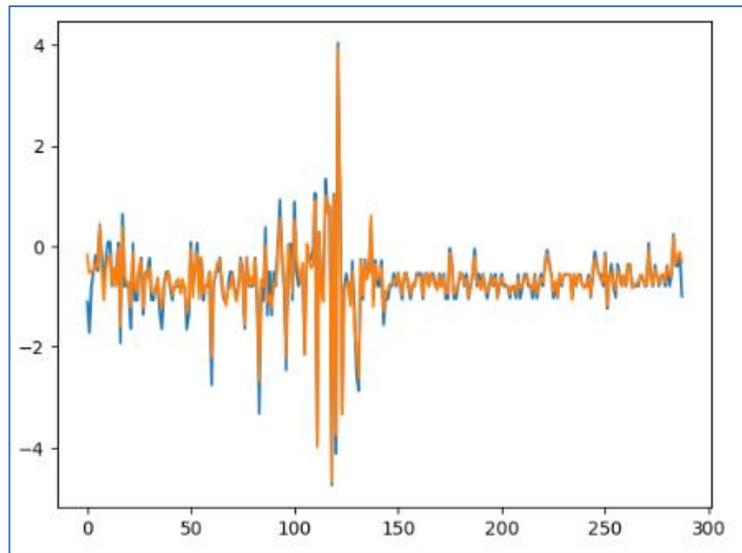


Figura 7. Señal original (azul) y reconstruida (naranja) en el segmento de entrenamiento. Elaboración propia.

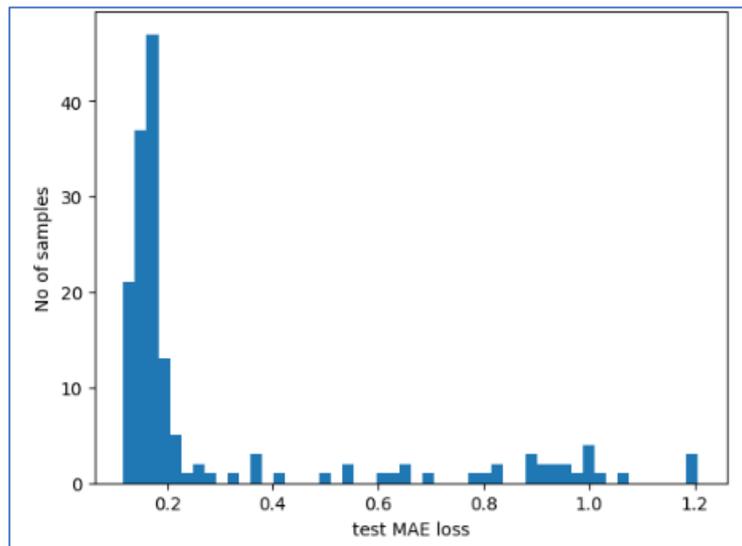


Figura 8. Distribución de MAE en el segmento de prueba. Elaboración propia.

En la anterior figura 8, puede observarse que aproximadamente 38 secuencias de las 163 secuencias de prueba resultaron anómalas ($MAE > \text{umbral de error de reconstrucción}$). Estas secuencias, en la STF inicial corresponden a los años posteriores a 2010, lo que es consistente con la dinámica económica nacional, en nuestra historia contemporánea.

Finalmente, en la figura 9 se pueden visualizar las anomalías. El balance de error de detección de anomalías en el conjunto de pruebas determino que 2 de las 38, fueron etiquetadas incorrectamente, lo que representa un 94,7% de acierto en la detección.

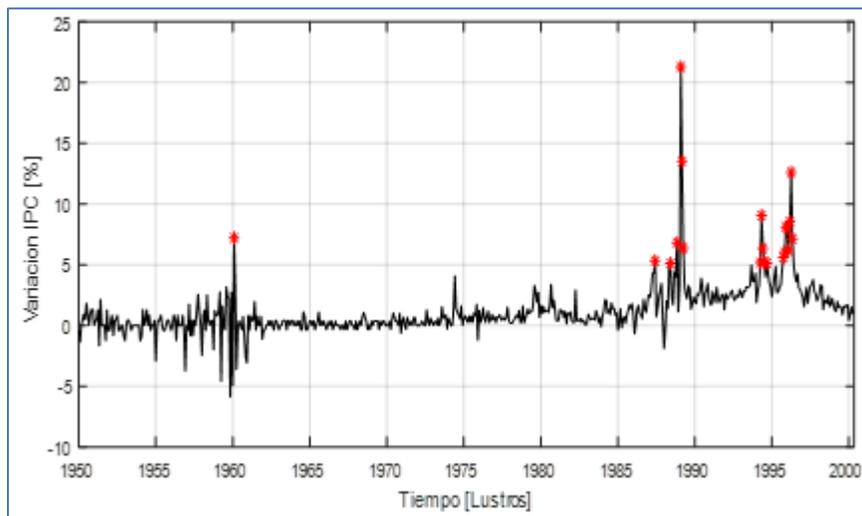


Figura 9. Ubicación de anomalías en la serie de tiempo financiera en estudio: IPC 1950-2022. Elaboración propia.

Conclusiones

En este trabajo, se describe un sistema automatizado para detección de anomalías en series de tiempo económicas o financieras (STF) usando DL. Para la construcción del modelo de DL, se empleó el algoritmo conocido como autocodificador, específicamente se implementó una versión unidimensional del autocodificador U-Net. Como STF de estudio, se usó el comportamiento del índice de precios al consumidor, IPC (Banco Central de Venezuela, 2022).

Como puede constatarse, en la sección de experimentos de este artículo, se lograron resultados satisfactorios, alcanzando una precisión de 94,6 % en el objetivo.

Referencias

- Bakumenko, A. y Elragal, A. (2022). **Detecting Anomalies in Financial Data Using Machine Learning Algorithms, Algorithms**. *Systems* 2022, 10, 130. Disponible en: <https://doi.org/10.3390/systems10050130>
- Banco Central de Venezuela (01 de agosto de 2022). **Índice General de Precios al Consumidor, Área Metropolitana de Caracas, Serie desde 1950, (Base: diciembre 2007=100)**. Disponible en: <https://www.bcv.org.ve/estadisticas/consumidor>
- Bengio Yoshua, I. Goodfellow, R. Y Courville, A. (2016). **Deep Learning**. MIT Press.
- Borges H., Reza A. y Maseglia, F. (2021). **Anomaly Detection in Time Series. Transactions on Large-Scale Data-and Knowledge-Centered Systems L, LNCS. TLDKS-12930**, pp.46-62, *Lecture Notes in Computer Science*. Transactions on Large-Scale Data- and Knowledge-Centered Systems, 978-3-662-64553-6. Doi: 10.1007/978-3-662-64553-6_3. lirmm-03359500.
- Caicedo Palacio, M. (2018). **Análisis de las Series de Tiempo del IPC en la Predicción de la Inflación en Colombia**. Tesis de Grado. Universidad del Valle Facultad de Ciencias Naturales y Exactas.
- Chollet, F. (2018). **Deep Learning with Python**. NY, USA: Manning Publications Co, Shelter Island.
- Crépey S., N. Lehdili, N., Madhar, Y. y Thomas, M. (2022). **Anomaly Detection on Financial Time Series by Principal Component Analysis and Neural Networks**. arXiv:2209.11686v1 [q-fin.ST] 22 Sep 2022.
- Hinton G. y Salakhutdinov, R. (2006). **Reducing the Dimensionality of Data with Neural Networks**. *Science, New Series*, 313(5786), pp. 504-507.
- Korczak, J. y Marcin, R. (2017). **Deep Learning for Financial Time Series Forecasting in A-Trader System**. Proceedings of the Federated Conference on Computer Science and Information Systems, pp. 905–912, DOI: 10.15439/2017F449, ISSN 2300-5963 ACSIS.
- León Anaya, L., Landassuri Moreno, V., Orozco Aguirre, H. y Quintana López, M. (2018). **Predicción del IPC mexicano combinando modelos**

económicos e inteligencia artificial. *Rev. Mex. Econ. Finanz.*
Disponible en: <https://doi.org/10.21919/remef.v13i4.342>.

López-Avila L., Acosta-Mendoza, N. y Gago-Alonso, A., (2019). Detección de anomalías basada en aprendizaje profundo: Revisión. *Revista Cubana de Ciencias Informáticas*. Disponible en: <http://rcci.uci.cu>

Navon, A. y Keller, Y. (2017). **Financial Time Series Prediction using Deep Learning**. arXiv:1711.04174v1 [eess.SP]

Pavithrasv, A. (2020). **Timeseries Anomaly Detection Using an Autoencoder, Keras**. Disponible en: https://keras.io/examples/timeseries/timeseries_anomaly_detection/.

Quilis, E. y Frutos Vivar, R. (1999). **Características inflacionarias de la economía española. Un análisis ARMA vectorial**. Instituto de Estudios Fiscales, Papel de Trabajo N. 9/99.

Ronneberger, O., Fischer, P. and Brox, T.. (2015). **U-Net: Convolutional Networks for Biomedical Image Segmentation**. *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. Disponible en: <https://lmb.informatik.uni-freiburg.de/people/ronneber/u-net/>

Schaffer A., Timothy, A., y Pearson, S. (2021). **Interrupted time series analysis using autoregressive integrated moving average (ARIMA)**. *BMC Medical Research Methodology*. Disponible en: <https://doi.org/10.1186/s12874-021-01235-8>

Schmidl S., Wenig, P. y Papenbrock, T. (2022). **Anomaly Detection in Time Series: A Comprehensive Evaluation**. *Proceedings of the VLDB Endowment*, 15 (9), ISSN 2150-8097. doi:10.14778/3538598.3538602.

Wei Y., Jang-Jaccard, J. Wen Xu, F., Camtepe, S. y Boulic, M. (2022). **LSTM-Autoencoder based Anomaly Detection for Indoor Air Quality Time Series Data**, arXiv:2204.06701v1 [cs.LG]. Disponible en: <https://doi.org/10.48550/arXiv.2204.06701>.

Yan W., (2020). **Computational Methods for Deep Learning**. *Springer, Texts in Computer Science*. Disponible en: https://doi.org/10.1007/978-3-030-61081-4_4.

Zhou H., Xuan Zhang, K., Wu, G. y Yazidi, A. (2022). **Contrastive autoencoder for anomaly detection in multivariate time series**. *Elsevier, Information Sciences*, 610, pp. 266-280.