

**Predicción de los niveles de NO_2 en el aire de Seúl
mediante el modelo Temporal Fusion Transformer (TFT)**

Carlos López Pérez^a (carlosdl@uninorte.edu.co),
Miguel Herrera Rocha^a (mherrerar@uninorte.edu.co)

^a Msc. en Estadística Aplicada, Departamento de Matemáticas y Estadística,
Universidad del Norte, Barranquilla, Colombia

Predicción de los niveles de NO_2 en el aire de Seúl mediante el modelo Temporal Fusion Transformer (TFT)

Carlos López Pérez^a, Miguel Herrera Rocha^a

^a*Msc. en Estadística Aplicada, Departamento de Matemáticas y Estadística, Universidad del Norte, Barranquilla, Colombia.*

Resumen

El monitoreo de la calidad del aire es esencial para comprender los efectos de la contaminación atmosférica en la salud pública y el medio ambiente. En este estudio, se aborda la predicción de los niveles de dióxido de nitrógeno (NO_2) en el aire de Seúl, Corea del Sur, utilizando datos de 1988 a 2021. A pesar de la aplicación de modelos de predicción en series temporales, como los transformers en la predicción de contaminantes como el ozono (O_3), PM y CO_2 , se han encontrado pocos estudios que utilicen específicamente el modelo Temporal Fusion Transformer (TFT) para este propósito. En este trabajo, se emplea el modelo TFT para predecir los niveles de NO_2 , evaluando su capacidad para capturar dependencias temporales y realizar predicciones precisas. Los resultados muestran que el modelo TFT puede identificar patrones complejos y no lineales en los datos históricos, superando a otros modelos tradicionales en términos de precisión. Este estudio destaca la efectividad de los transformers, y en particular el TFT, para predecir contaminantes atmosféricos, sugiriendo su potencial para ser implementado en la predicción y gestión de la calidad del aire a nivel global.

Palabras Clave: NO_2 , Temporal Fusion Transformer, calidad del aire, serie temporal, transformers, modelos de predicción.

Email addresses: carlosdl@uninorte.edu.co (Carlos López Pérez),
mherrerar@uninorte.edu.co (Miguel Herrera Rocha)

1. Datos

Los datos utilizados en este estudio provienen del conjunto de datos histórico sobre la calidad del aire en Seúl, Corea del Sur desde 1988 hasta 2021. Este conjunto de datos contiene más de 5 millones de observaciones en diferentes puntos de medición por cada hora relacionadas con varios contaminantes atmosféricos tales como Ozono (O_3), PM2.5, PM10 y SO_2 y los niveles de dióxido de nitrógeno (NO_2). La variable de interés en este análisis es la concentración de NO_2 , que se mide en partículas por millón [ppm]. Este conjunto de datos está estructurado en formato de serie temporal lo que facilita su uso para modelos de predicción basados en este tipo de datos. Los autores proporcionan acceso completo al conjunto de datos, permitiendo su replicación para futuras investigaciones y puede ser consultado en la bibliografía (Hyun, 2021).

El dióxido de nitrógeno NO_2 es un contaminante atmosférico común que puede tener efectos negativos en la salud humana y el medio ambiente ya que la exposición a altos niveles ambientales puede aumentar el riesgo de infecciones del tracto respiratorio en la población debido a la interacción del contaminante con el sistema inmune (Stieb et al., 2021). El análisis de las concentraciones de NO_2 en el aire es crucial para entender los patrones de contaminación y tomar medidas adecuadas para su control.

1.1. Análisis exploratorio de los datos - EDA

El objetivo del presente análisis exploratorio de datos es comprender la estructura de los datos históricos de concentraciones de NO_2 en el aire de Seúl, identificar patrones temporales, detectar datos atípicos y determinar la existencia de tendencia o estacionalidad de la variable objetivo. Este proceso es fundamental para preparar los datos antes de aplicar modelos predictivos, ya que permite entender la distribución, tendencias y estacionalidad de la serie, así como manejar valores faltantes o atípicos.

El conjunto de datos inicial contenía un total de 5,984,782 registros, con mediciones de concentraciones. Una de las variables clave en la base de datos es

'loc', que hace referencia al punto de medición donde se tomaron las muestras de NO_2 . Esta variable indica el punto de cada medición, por tanto, se decidió trabajar con una única ubicación. Para ello, se identificaron los valores únicos de loc y se seleccionó la ubicación con el mayor número de registros, que resultó ser la ubicación 124. Esta locación contaba con 288,736 mediciones, lo que la convierte en la más representativa y adecuada para realizar un análisis detallado de las concentraciones de NO_2 a lo largo del tiempo.

Durante el análisis, se identificó una observación atípica en la concentración de NO_2 , cuyo valor era mucho más alto que el resto de los datos, registrando 0.687 ppm. Este valor (contrastaba notablemente con el percentil 99 que era de 0.08 ppm.) Dada esta discrepancia, se decidió tratar este valor como un outlier y reemplazarlo por un valor NaN para evitar distorsiones en el análisis. Posteriormente, este valor faltante fue imputado junto con otros valores faltantes en la serie. En total, se encontraron 6,898 valores faltantes en la variable de estudio, lo que representó aproximadamente el 2.39% del total de observaciones. Estos valores fueron imputados utilizando el método de interpolación cuadrática, lo que permitió mantener la integridad de la serie temporal sin perder información relevante. A continuación, se presenta un resumen estadístico de las concentraciones de NO_2 después de la imputación:

Métrica	Valor
Cantidad de datos	288736.00
Media	0.0318
Desviación estándar	0.0171
Valor Mínimo	0.0000
Cuartil 1	0.0190
Mediana	0.0290
Cuartil 3	0.0420
Valor Máximo	0.2040
Rango Intercuartílico (IRQ)	0.0230

Table 1: Resumen estadístico de las concentraciones de NO_2 .

La distribución de las concentraciones de NO_2 , visualizada en la gráfica de densidad, muestra un sesgo a la izquierda, con la mayoría de los valores concentrados en el rango inferior (entre 0.00 y 0.05 ppm), pero con una cola que se extiende hacia valores más altos. Esto indica que, aunque la mayoría de las mediciones son bajas, existen muy ocasionales picos de concentración que elevan la media.

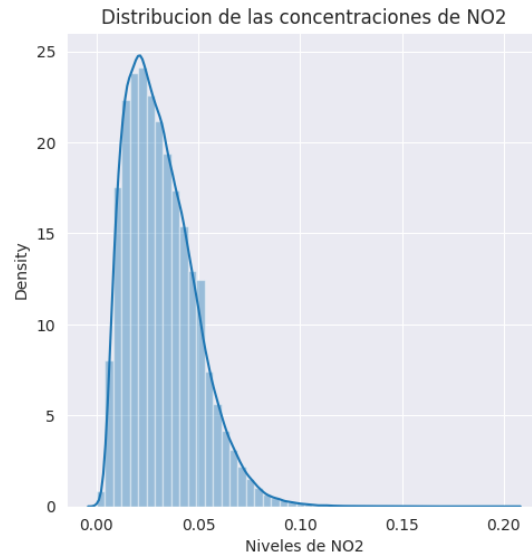


Figure 1: Distribución de la concentración de NO_2 en ppm

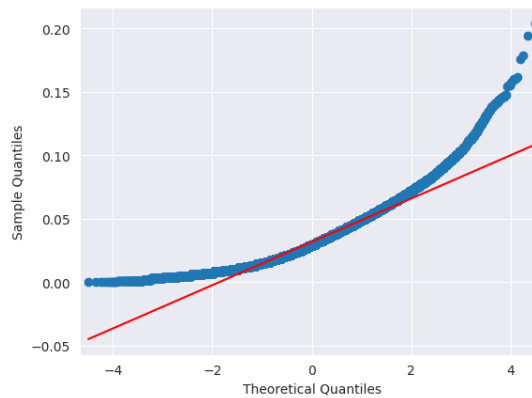


Figure 2: Grafico de cuantiles - NO_2

Por otro lado, el gráfico de cuantiles revela una desviación significativa de la normalidad, especialmente en los extremos, donde los valores observados no siguen la línea teórica de una distribución normal. Esto se confirma con la prueba de Shapiro-Wilk, que arrojó un p-valor extremadamente bajo ($p < 0.0001$) y un estadístico de 0.949, lo que lleva a rechazar la hipótesis nula de normalidad. En conclusión, las concentraciones de NO_2 no siguen una distribución normal.

A continuación se presenta en las gráficas 3 y 4 una visualización de la serie temporal con el objetivo de observar la tendencia general y posibles anomalías. En la siguiente gráfica se observan tres series superpuestas que representan diferentes niveles de agregación temporal de las concentraciones de NO_2 .

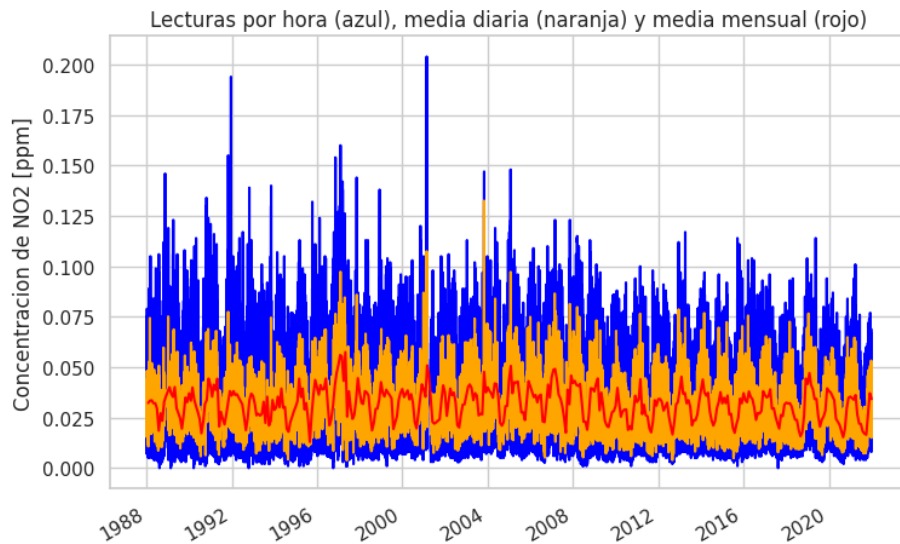


Figure 3: Visualización concentraciones de NO_2 en el tiempo

La serie graficada en azul muestra las mediciones originales de NO_2 tomadas por hora. Se puede apreciar una alta variabilidad, con fluctuaciones significativas a lo largo del tiempo. Para suavizar la variabilidad horaria y resaltar patrones a más largo plazo, se calculó la media diaria (naranja) y mensual (rojo) de las concentraciones de NO_2 , en esta gráfica podemos sospechar que no se encuentra una tendencia marcada (ya sea creciente o decreciente) pero que

con las medias mensuales (rojo) se sospecha de un patrón estacional, esto será verificado mas adelante.

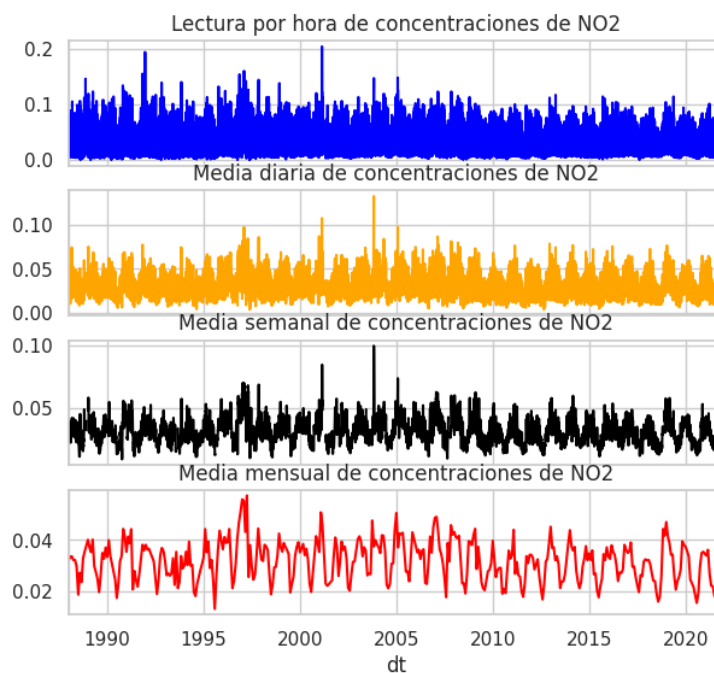


Figure 4: Concentraciones de NO_2 variando su agrupación temporal

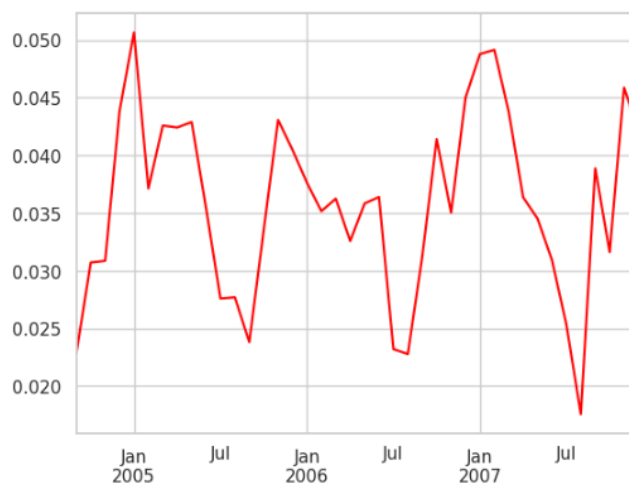


Figure 5: Ampliación de la serie de medias mensuales entre 2005-2008

Si se amplia en una sección de la serie con las medias mensuales, es posible ver un patrón que se repite de forma anual en donde las concentraciones de NO_2 parecen llegar a su punto mas bajo en los meses de Julio y agosto (ver gráfica 5)

Se realizó una descomposición de la serie temporal mediante medias móviles separa la serie en sus componentes principales: tendencia, estacionalidad y residuales. La estacionalidad muestra patrones recurrentes, con picos y valles que se repiten anualmente, mientras que los residuales capturan fluctuaciones no explicadas por los otros componentes. Este análisis ayuda a entender la estructura de la serie y a identificar patrones clave para futuros modelos predictivos.

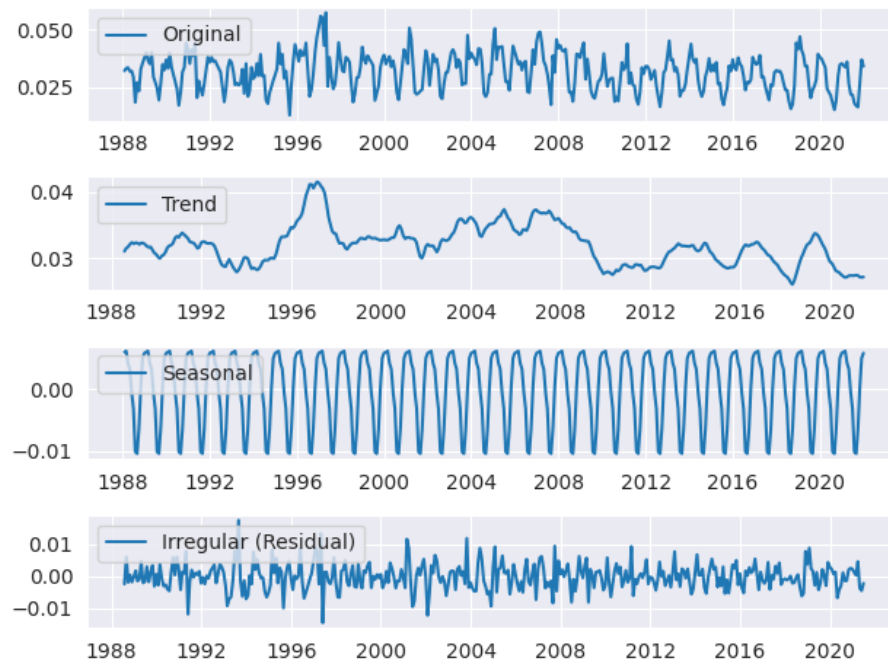


Figure 6: Descomposición de la serie temporal

Se realizó un gráfico autocorrelación ACF de las concentraciones medias mensuales de dióxido de nitrógeno el cual muestra correlaciones significativas hasta ciertos rezagos, con un comportamiento oscilante senoidal cuya significancia decrece lentamente a medida que aumentan los rezagos. Este patrón es

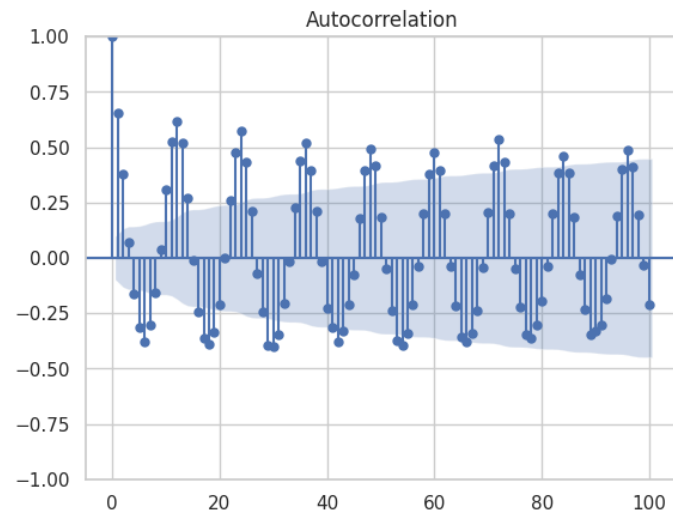


Figure 7: Grafico de Autocorrelacion (ACF)

indicativo de una componente estacional o tendencia cíclica en las concentraciones de NO_2 . Las correlaciones que se mantienen significativas sugieren que las concentraciones mensuales de NO_2 están fuertemente influenciadas por los valores previos.

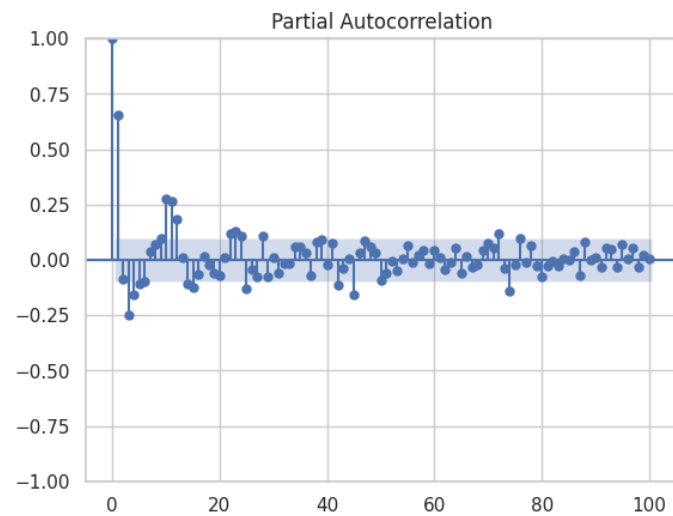


Figure 8: Grafico de Autocorrelación Parcial (PACF)

La PACF muestra una dependencia significativa en los primeros lags, pero esta dependencia se corta rápidamente después de unos pocos. Esto sugiere que la influencia de los valores pasados en el presente es más fuerte en el corto plazo, lo que es típico en series con componentes estacionales y de tendencia. Este comportamiento sugiere que un modelo AR podría capturar bien la estructura del proceso.

El gráfico de concentración de NO_2 por trimestre muestra una variación estacional, con los niveles más altos en el primer trimestre y los más bajos y menos dispersos en el tercer trimestre. Esto sugiere que las concentraciones de NO_2 están influenciadas por factores estacionales que varían a lo largo del año.

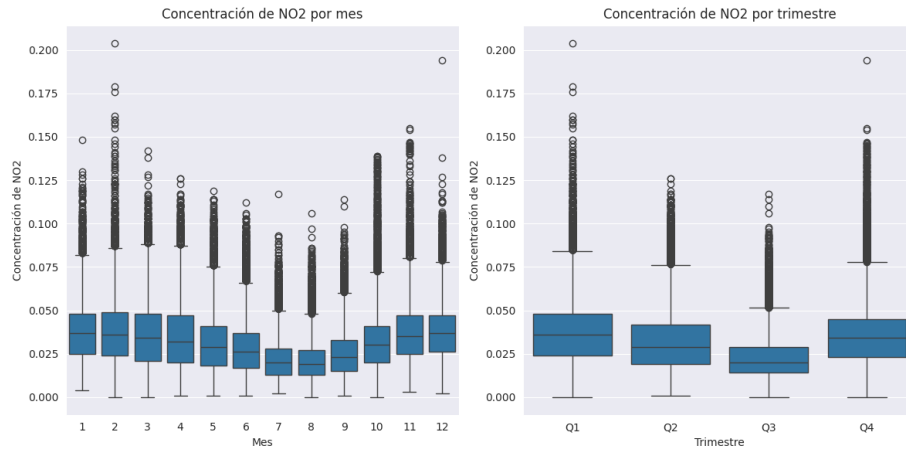


Figure 9: Concentración de NO_2 por trimestre

Al hacer esta comparación por mes se obtienen los siguientes resultados y se ve con mayor claridad cuales son los meses con mayores y menores concentraciones de dióxido de nitrógeno.

Por ultimo, se desea verificar la estacionalidad de las medias mensuales de las concentraciones de NO_2 por lo que se realiza una diferenciación a un mes y a 12 meses de rezago y una prueba Dickey fuller para confirmar la estacionalidad. En los resultados de la prueba de Dickey-Fuller para la concentración de NO_2 , obtuvimos un valor p de 0.021 sin diferenciar los datos, lo que indica que pode-

mos rechazar la hipótesis nula de que la serie tiene una raíz unitaria (es decir, que la serie no es estacionaria) al nivel de significancia del 5%. Esto sugiere que la serie podría ser estacionaria, aunque el valor p está cerca del umbral de 0.05, por lo que es un resultado algo marginal.

Por otro lado, al diferenciar los datos a 12 meses, el valor p disminuyó significativamente a $3.87e-07$, lo que indica una fuerte evidencia para rechazar la hipótesis nula de no estacionalidad bajo cualquier nivel de significancia. Esto sugiere que, después de aplicar una diferenciación a 12 meses, la serie temporal se vuelve más estacionaria, lo que implica que la tendencia a largo plazo ha sido eliminada y la serie muestra características más estacionarias.

En el siguiente gráfico se muestra el cambio en las graficas ACF y PACF con las diferenciaciones.

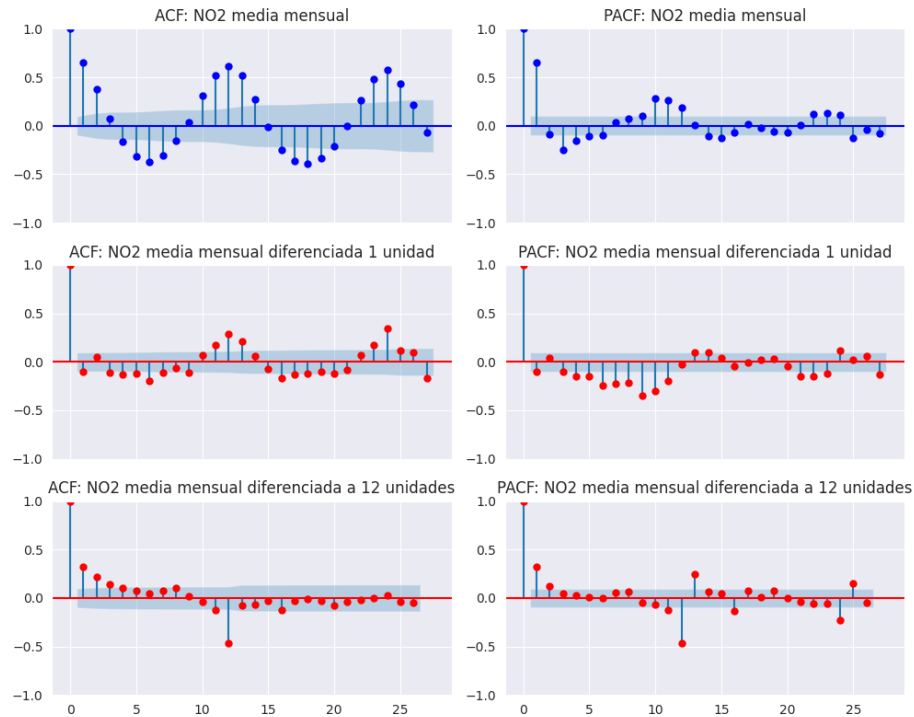


Figure 10: ACF y PACF de las concentraciones diferenciadas a 1 y 12 unidades

El análisis gráficos de las ACF y PACF sugiere que la serie original de la

media mensual de NO_2 no es estacionaria (*Aunque en Dickey Fuller se comprobó que si bajo ciertos niveles de significancia*), ya que presenta correlaciones significativas en múltiples rezagos. La diferenciación en una unidad reduce la autocorrelación, pero persisten algunos patrones, sin embargo, según los patrones hallados en el presente EDA, se decidió aplicar una diferenciación estacional de 12 unidades, al hacerlo, la serie parece estabilizarse, con autocorrelaciones más controladas.

2. Métodos

(Catalano, 2024) (Mattioli, 2024)

3. Resultados

4. Conclusiones

5. Apéndices

References

- Catalano, L. (2024). *A Transformer-based approach to air quality prediction in Milan through satellite imagery combined with meteorological and morphological data*. Ph.D. thesis Politecnico di Torino.
- Hyun, W. (2021). Seoul air quality historical data. URL: <https://www.kaggle.com/datasets/williamhyun/seoulairqualityhistoricdata/data>.
- Mattioli, F. (2024). *Enhancing Urban Air Quality Predictions with Temporal Fusion Transformers: A Methodological Evaluation*. Master's thesis Universidad Politecnica de Madrid Madrid, España.
- Stieb, D. M., Berjawi, R., Emode, M., Zheng, C., Salama, D., Hocking, R., Lyrette, N., Matz, C., Lavigne, E., & Shin, H. H. (2021). Systematic review and meta-analysis of cohort studies of long term outdoor nitrogen dioxide exposure and mortality. *PloS one*, 16, e0246451.