

TRABAJO PRÁCTICO OBLIGATORIO

SUBE Transacciones diarias

Valentín Lanús, Evangelina Ovelar, Santino Lozano
Materia: Ciencias de Datos – UADE

Planteamiento del Problema y Objetivo

El transporte público, como servicio esencial, requiere una planificación rigurosa. Su eficiencia operativa (frecuencia, rutas y asignación de recursos) depende directamente de una comprensión precisa del comportamiento de los usuarios. La tarjeta SUBE genera un volumen masivo de datos transaccionales diarios, reflejando patrones de uso a nivel temporal y geográfico.

El desafío: La incertidumbre en la demanda diaria impacta directamente en la calidad del servicio y en la optimización de los costos operativos.

Objetivo Principal del Proyecto

El objetivo central es desarrollar un modelo predictivo robusto para **anticipar la cantidad de validaciones diarias de tarjetas SUBE** en la Provincia de Buenos Aires.

Este análisis busca:

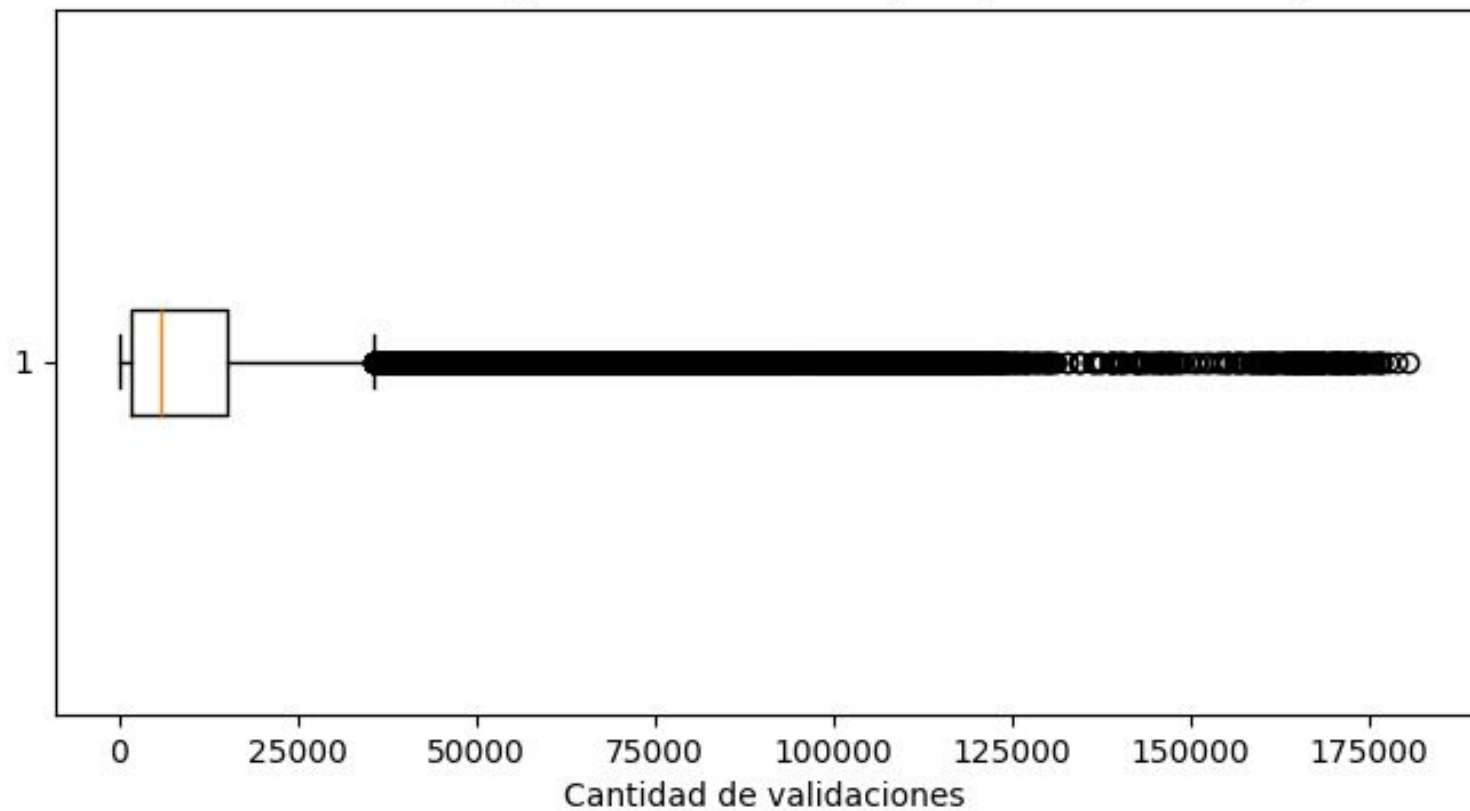
1. **Optimizar la planificación:** Permitir una asignación más eficiente de unidades y personal, minimizando demoras y sobrecargas.
2. **Generar información estratégica:** Ofrecer información valiosa a *stakeholders* para la toma de decisiones sobre tarifas, subsidios y expansión de rutas.
3. **Establecer un baseline:** Crear un modelo de referencia (*baseline*) basado en Machine Learning que capture los patrones de uso semanales y estacionales.

Limpieza y preparación de datos

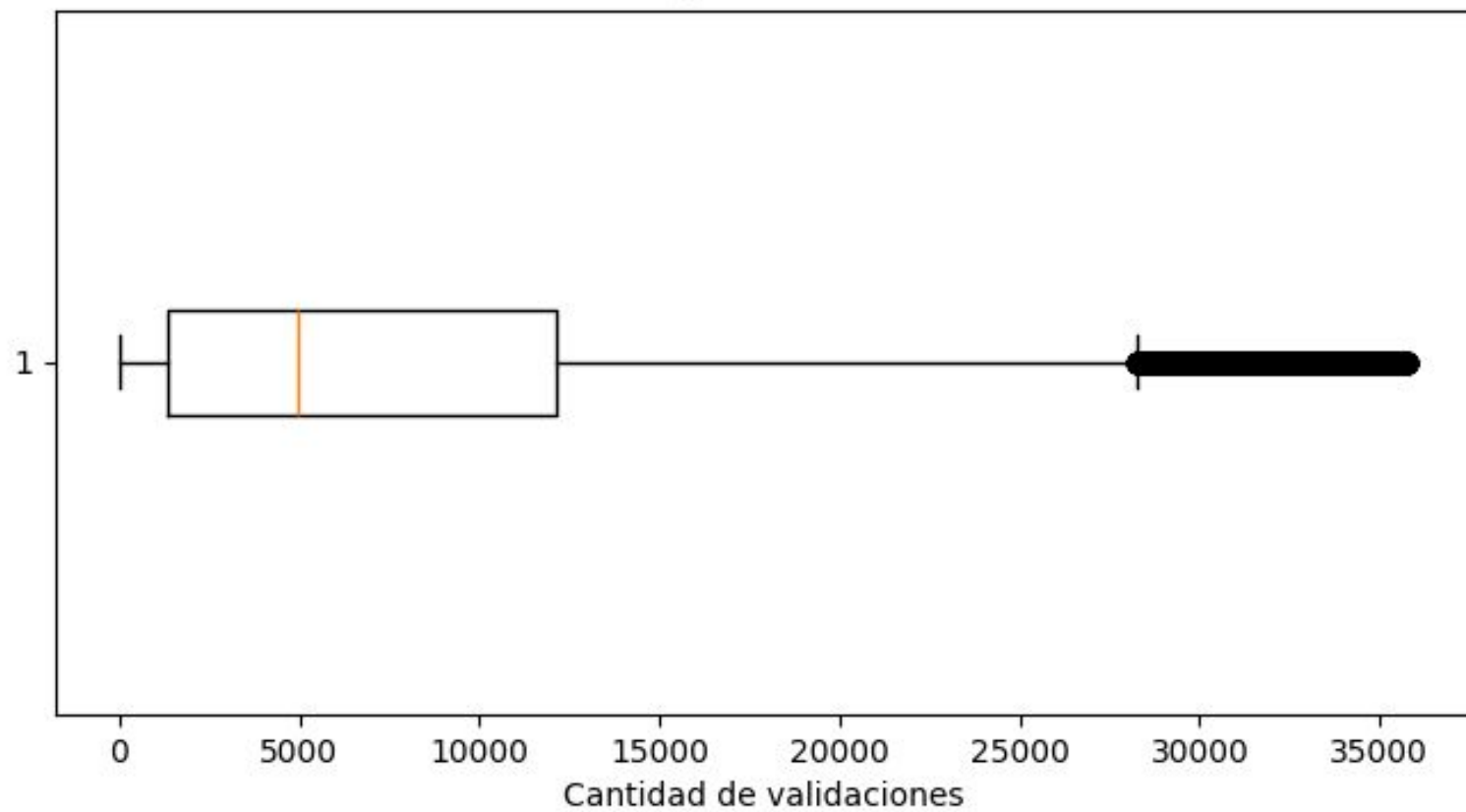
Flujo de Trabajo Aplicado:

1. **Filtro Geográfico (Scope):** El análisis se centró exclusivamente en los registros correspondientes a la **Provincia de Buenos Aires** (**PROVINCIA** = 'BUENOS AIRES').
2. **Conversión y Validación:** Se verificó que la variable clave **DIA_TRANSPORTE** fuera del tipo **datetime** y se eliminaron registros con valores nulos o con **CANTIDAD** (validaciones) $\leq 0\$$.
3. **Agrupación:** Los datos se agregaron a nivel diario (**DIA_TRANSPORTE**) para obtener la suma total de validaciones (**uso_diario**), creando nuestra serie de tiempo fundamental.
4. **Tratamiento de Outliers:** Se realizó un análisis exhaustivo de la variable **CANTIDAD**. Para mantener la robustez del modelo, se identificaron y eliminaron valores extremos (Outliers) en la serie de tiempo diaria utilizando el **método del Rango Intercuartílico (IQR)**, lo que asegura que las anomalías no distorsionen los patrones de uso normales.

Distribución original de CANTIDAD (con posibles outliers)



Distribución después de eliminar outliers



Anomalías detectadas y ajustes de datos

Análisis de Patrones Anómalos (Heatmap Semanal)

Durante la exploración detallada (usando *heatmaps* semanales y análisis de promedios), se identificaron caídas atípicas en la cantidad de validaciones en dos fechas específicas:

- **Jueves 11 de abril**
- **Jueves 9 de mayo**

Estas caídas fueron consideradas **anomalías o shocks discretos**, ya que no se correlacionaban con feriados o patrones estacionales conocidos.

Anomalías detectadas y ajustes de datos

Estrategia de Normalización

Para evitar que estos valores atípicos sesgaran el entrenamiento del modelo de Machine Learning y distorsionaran la tendencia general:

- Se aplicó un **ajuste proporcional** a los valores de esos dos días.
- El ajuste se basó en el promedio de las validaciones de otros jueves cercanos.

Resultado: Se logró mantener la coherencia de la serie de tiempo, asegurando que las caídas detectadas no fueran interpretadas por el modelo como patrones recurrentes, sino como eventos aislados.

Análisis exploratorio de datos

Indicadores Clave del Dataset SUBE — Buenos Aires

1. Volumen general

- Total de validaciones registradas: **976,879,363**
- Días cubiertos por el dataset: **366**
- Outliers eliminados: **10452**

2. Cobertura operativa

- Empresas activas: **121**
- Líneas de transporte analizadas: **370**
- Empresa con mayor cantidad de validaciones: **EMPRESA DE TRANSPORTE PERALTA RAMOS SACI**

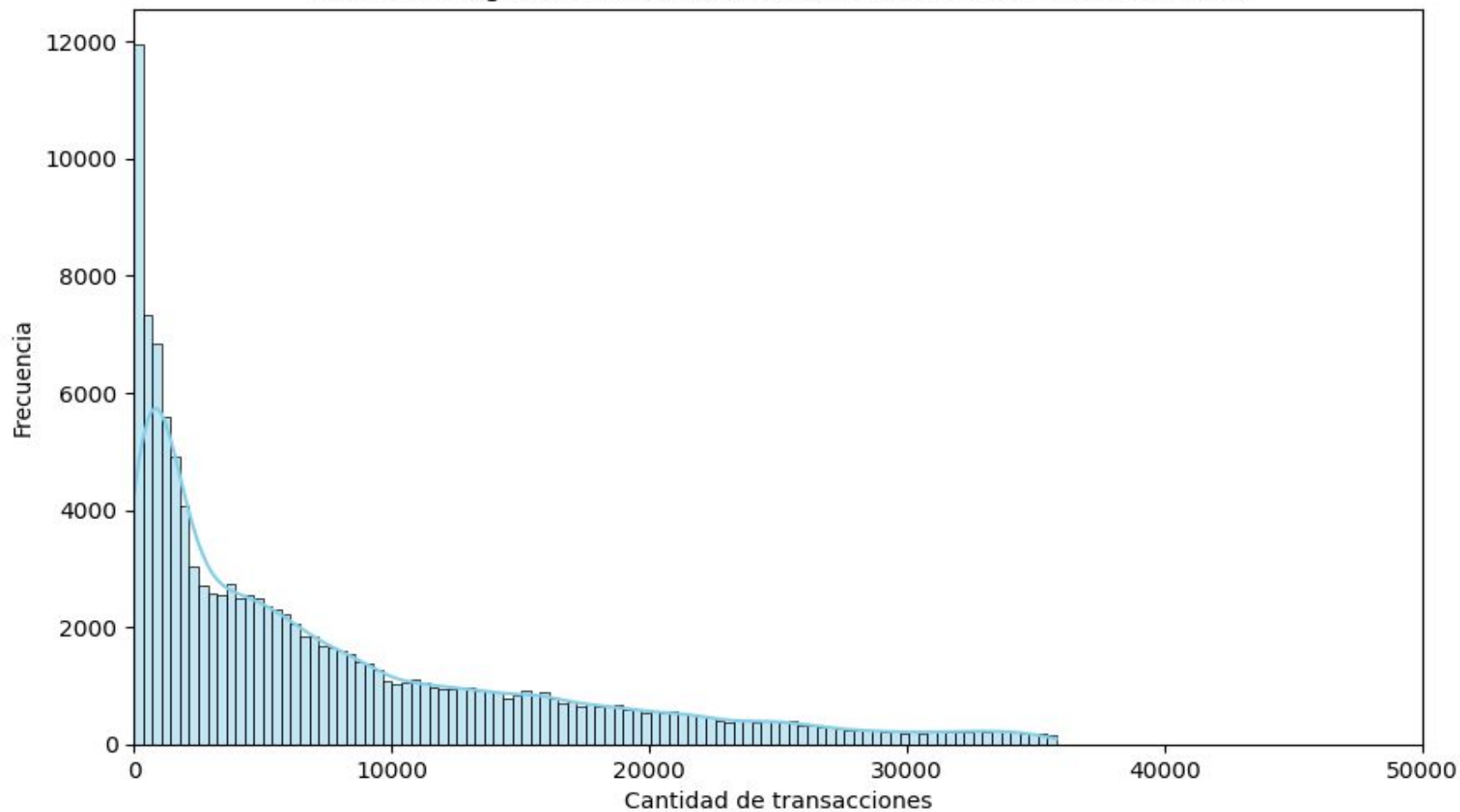
3. Promedios

- Promedio diario de transacciones: **2,669,069**
- Promedio por empresa: **9,031**

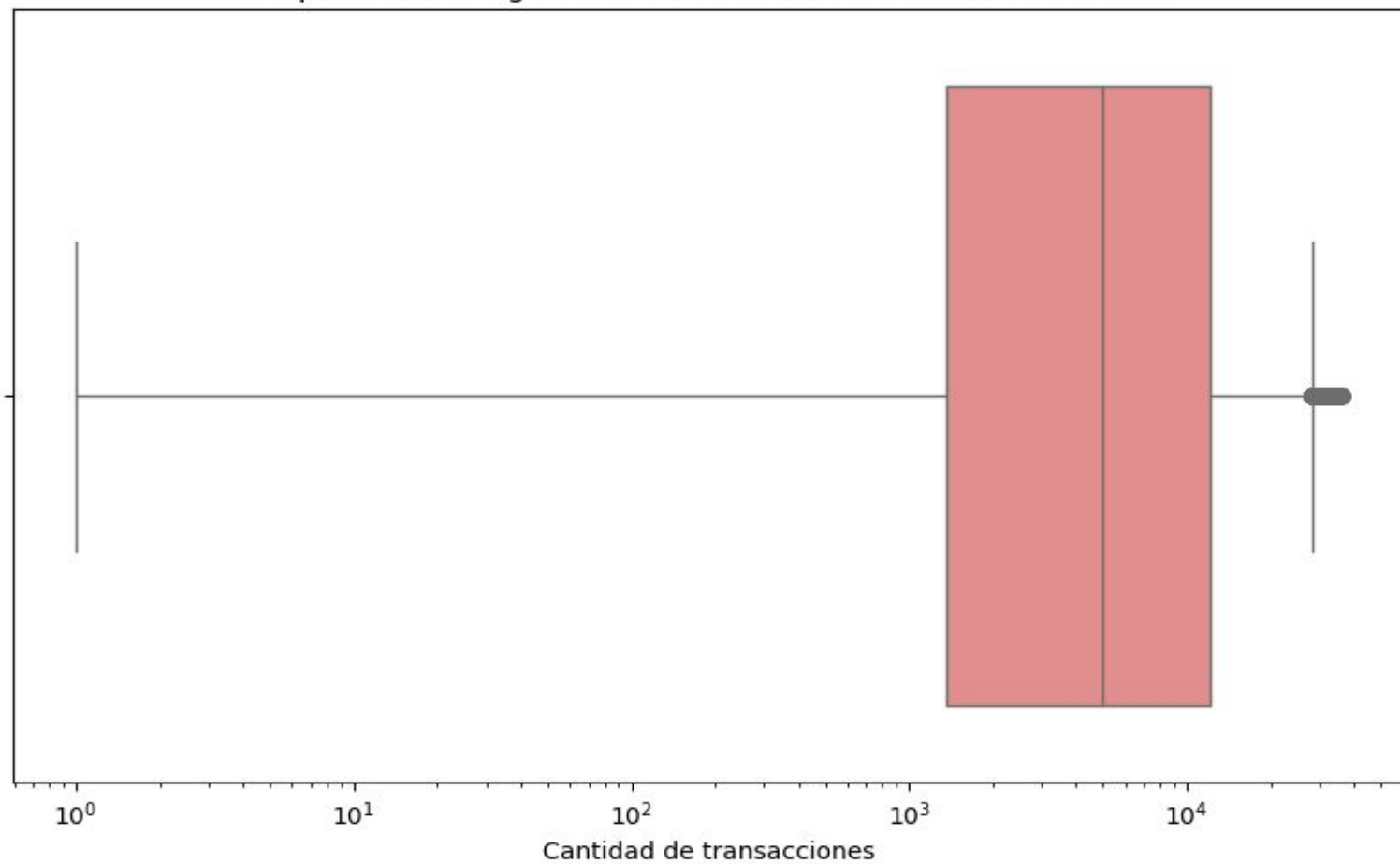
4. Distribución temporal

- Participación de fines de semana: **21.77%** del total

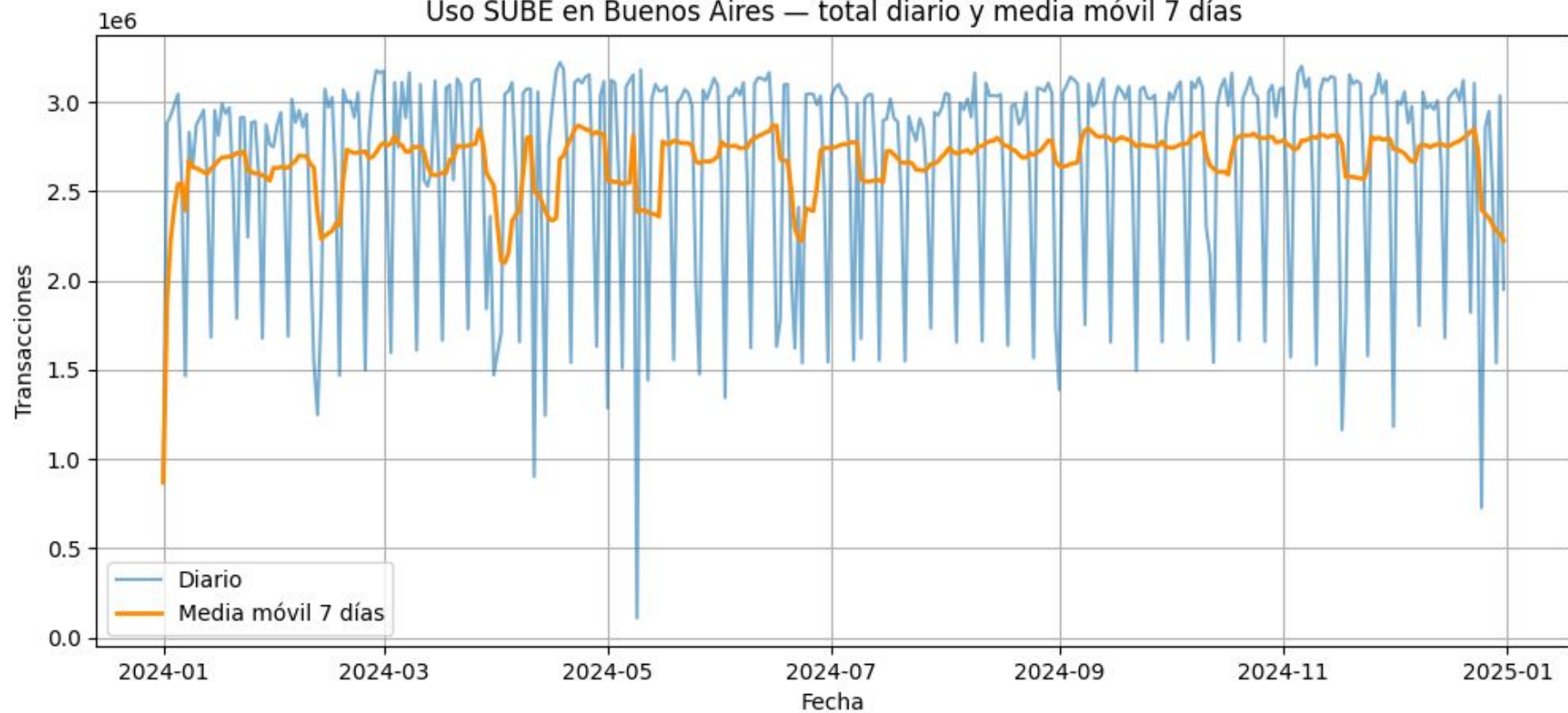
Distribución general de la cantidad de transacciones (Buenos Aires)



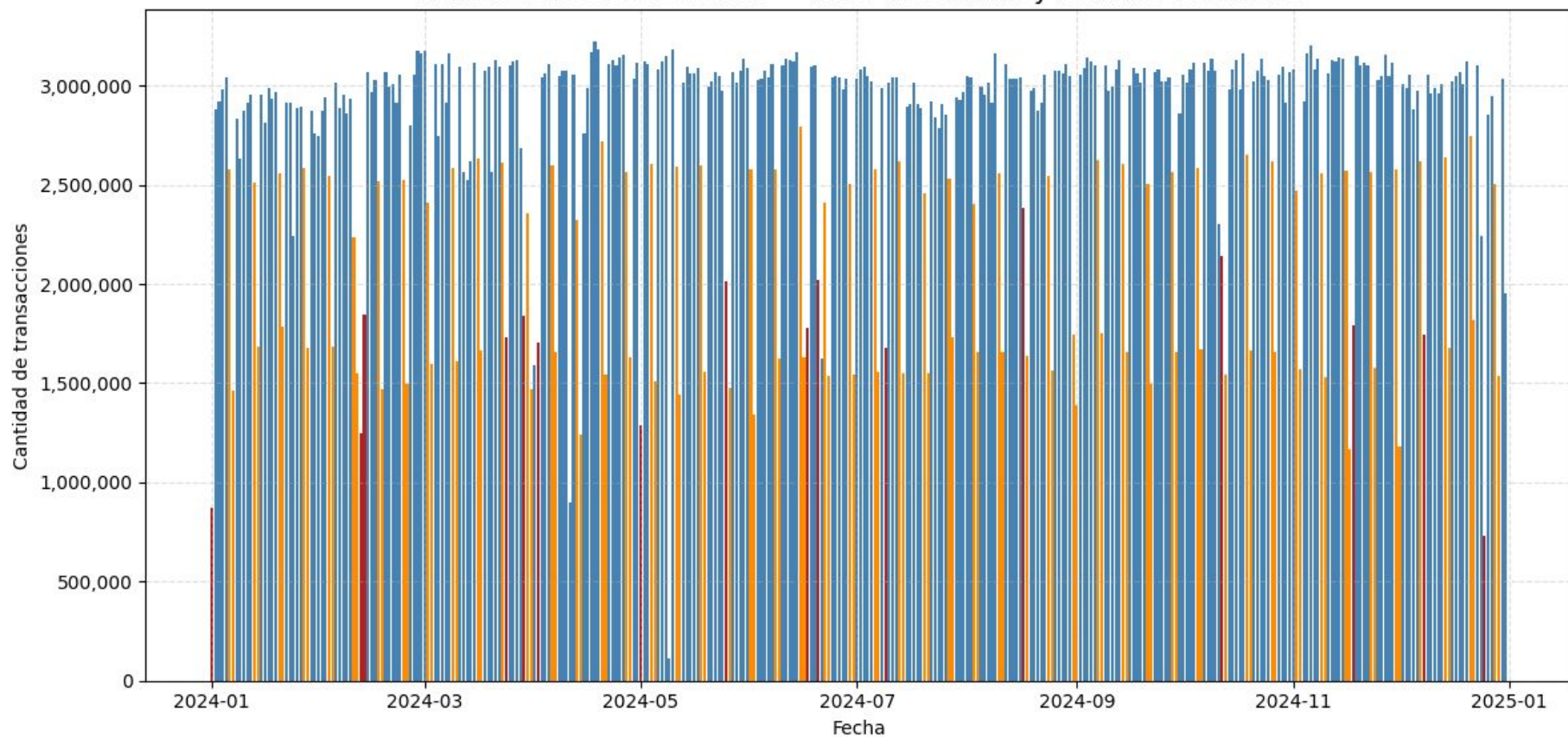
Boxplot (escala logarítmica) de transacciones en Buenos Aires

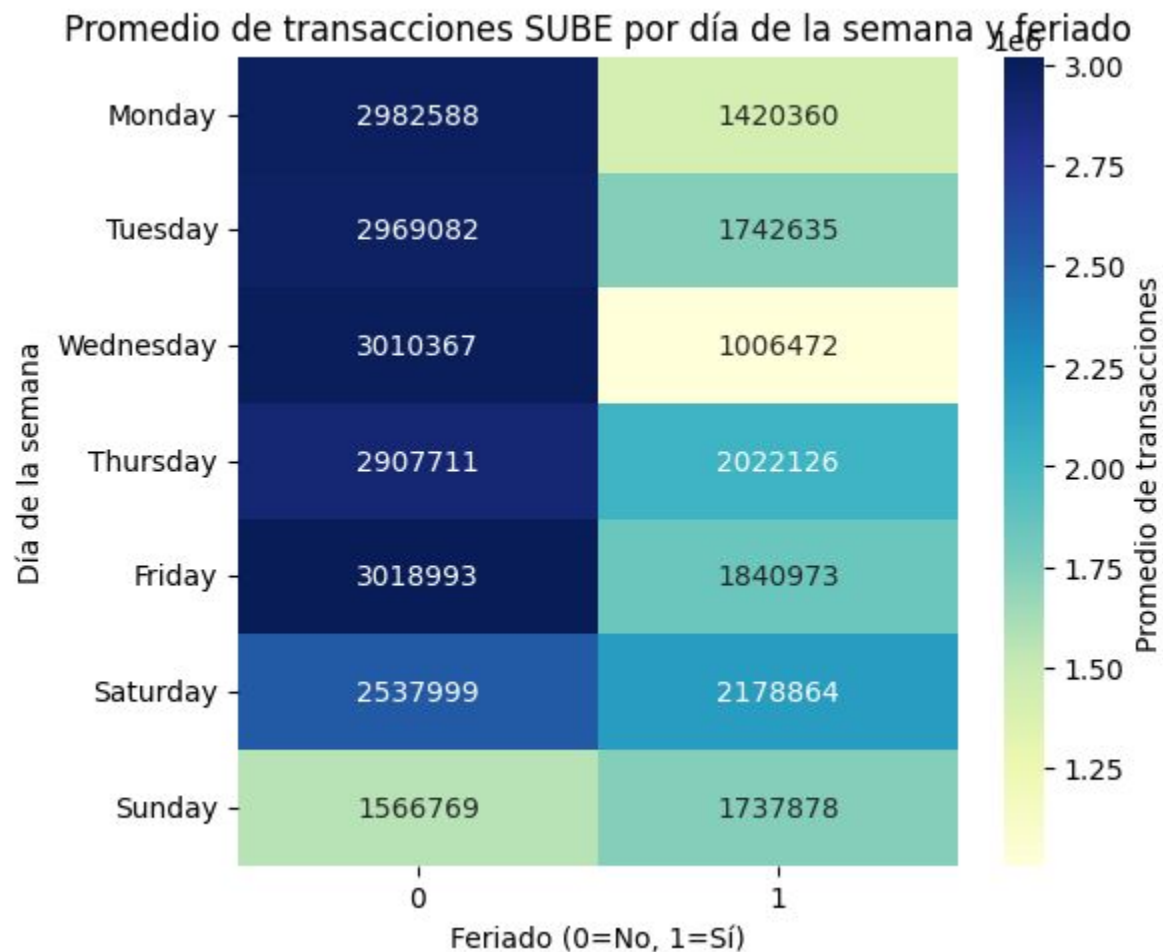


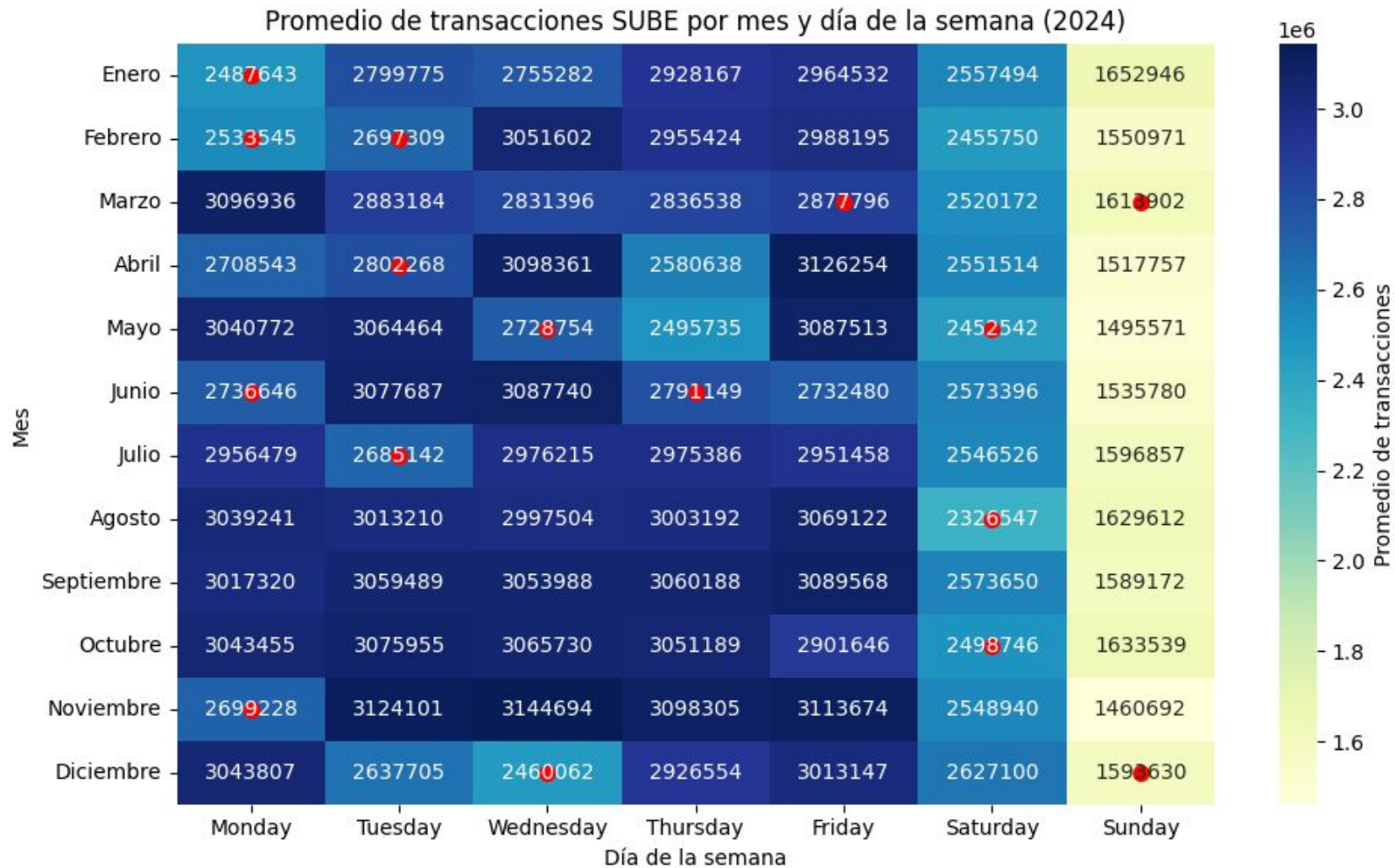
Uso SUBE en Buenos Aires — total diario y media móvil 7 días



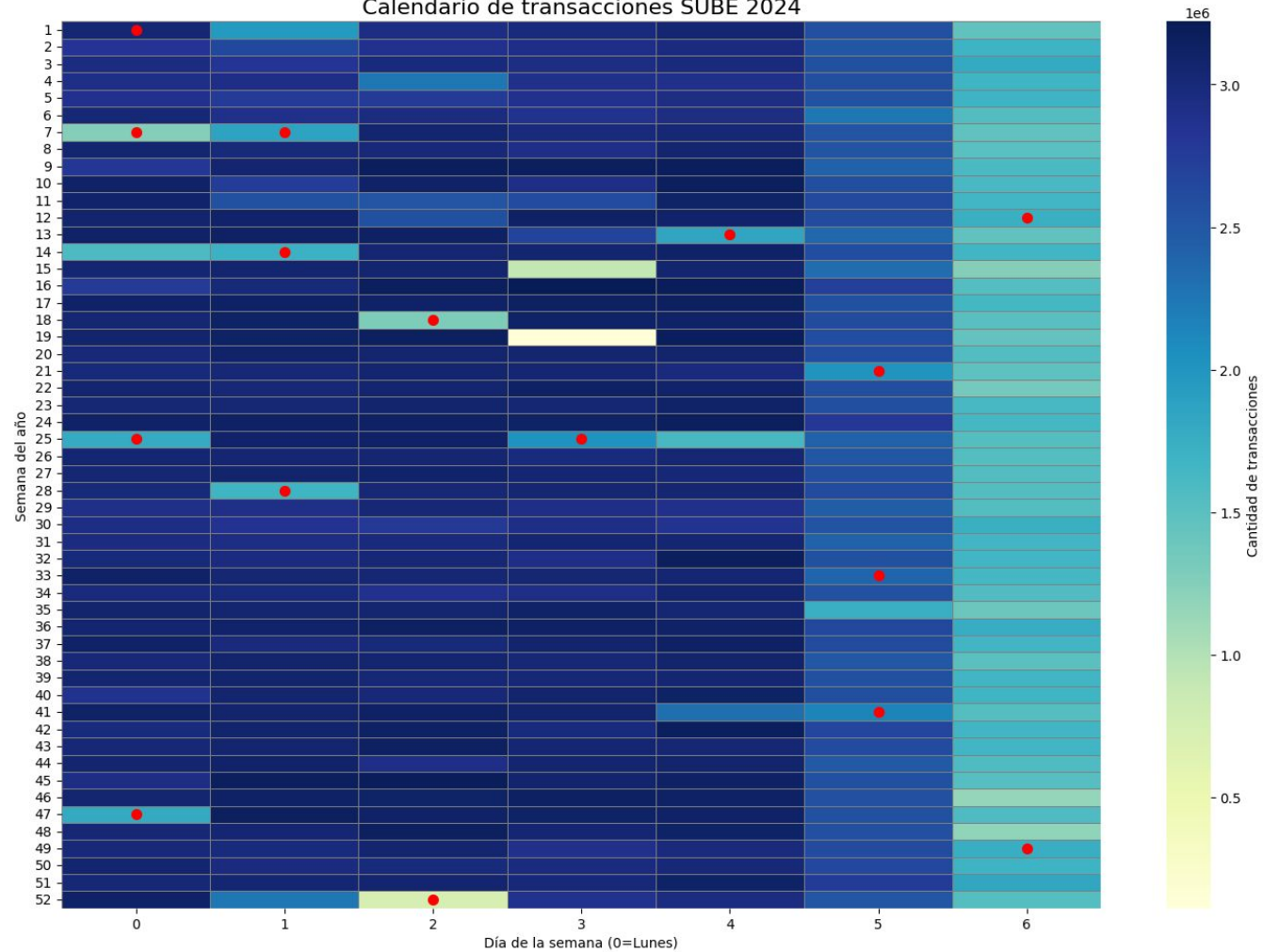
Transacciones SUBE diarias — fines de semana y feriados resaltados



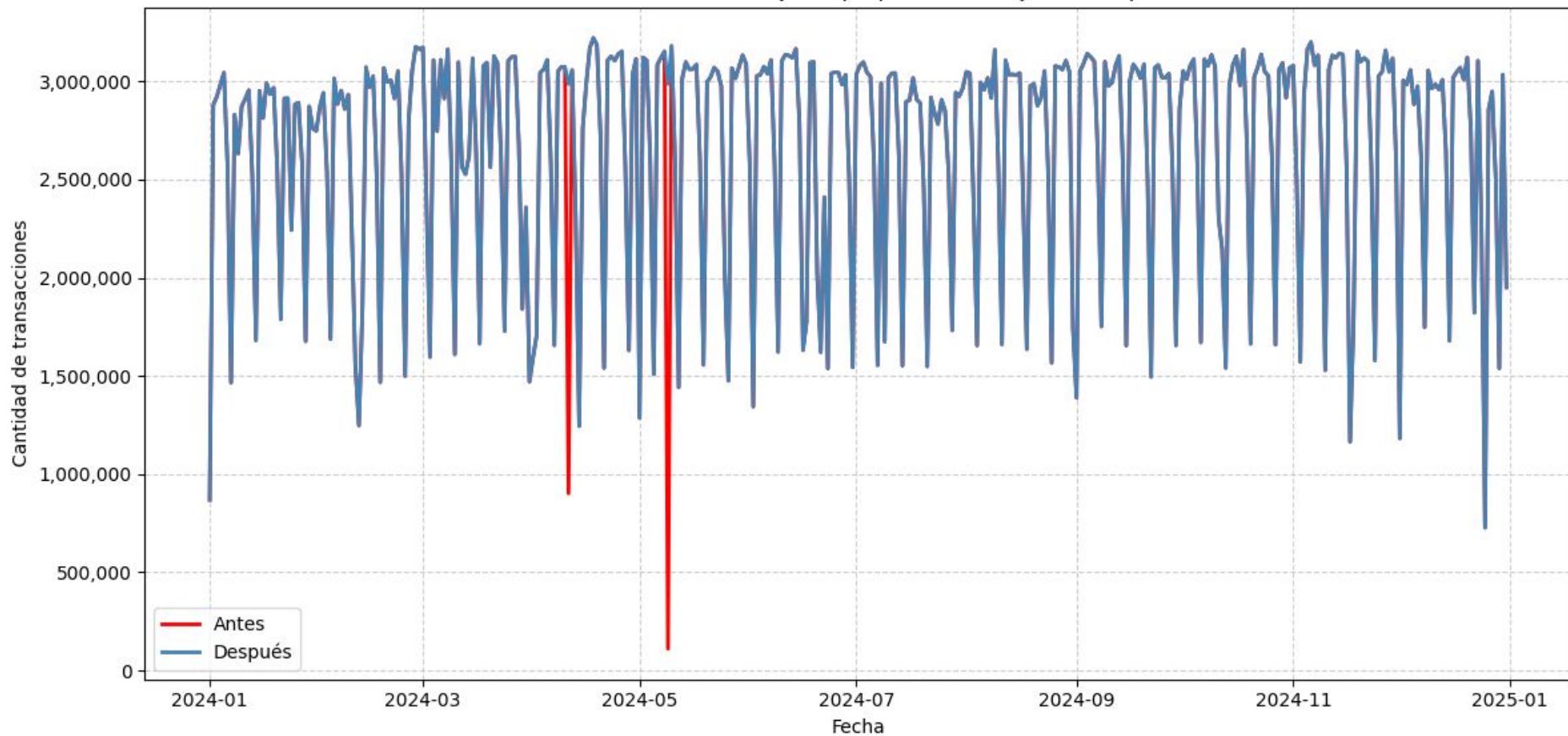




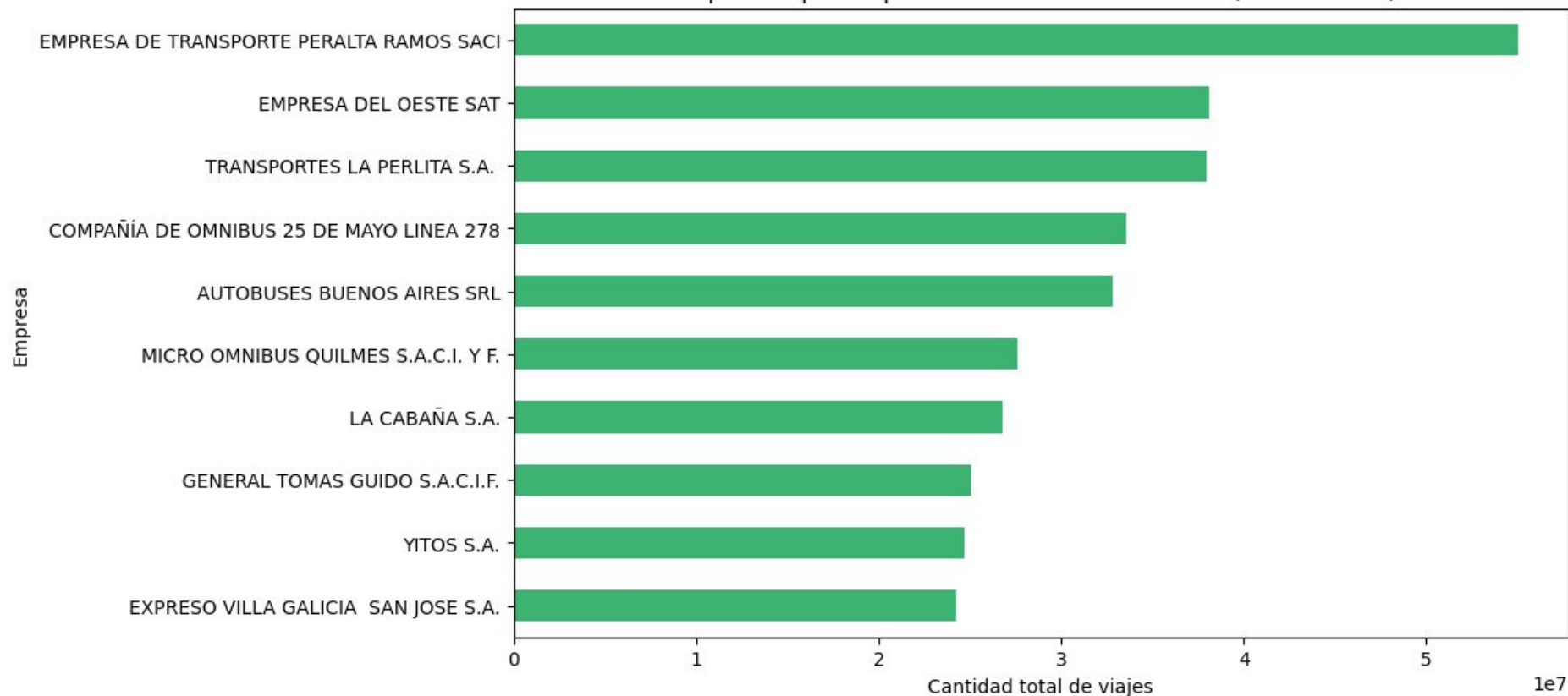
Calendario de transacciones SUBE 2024



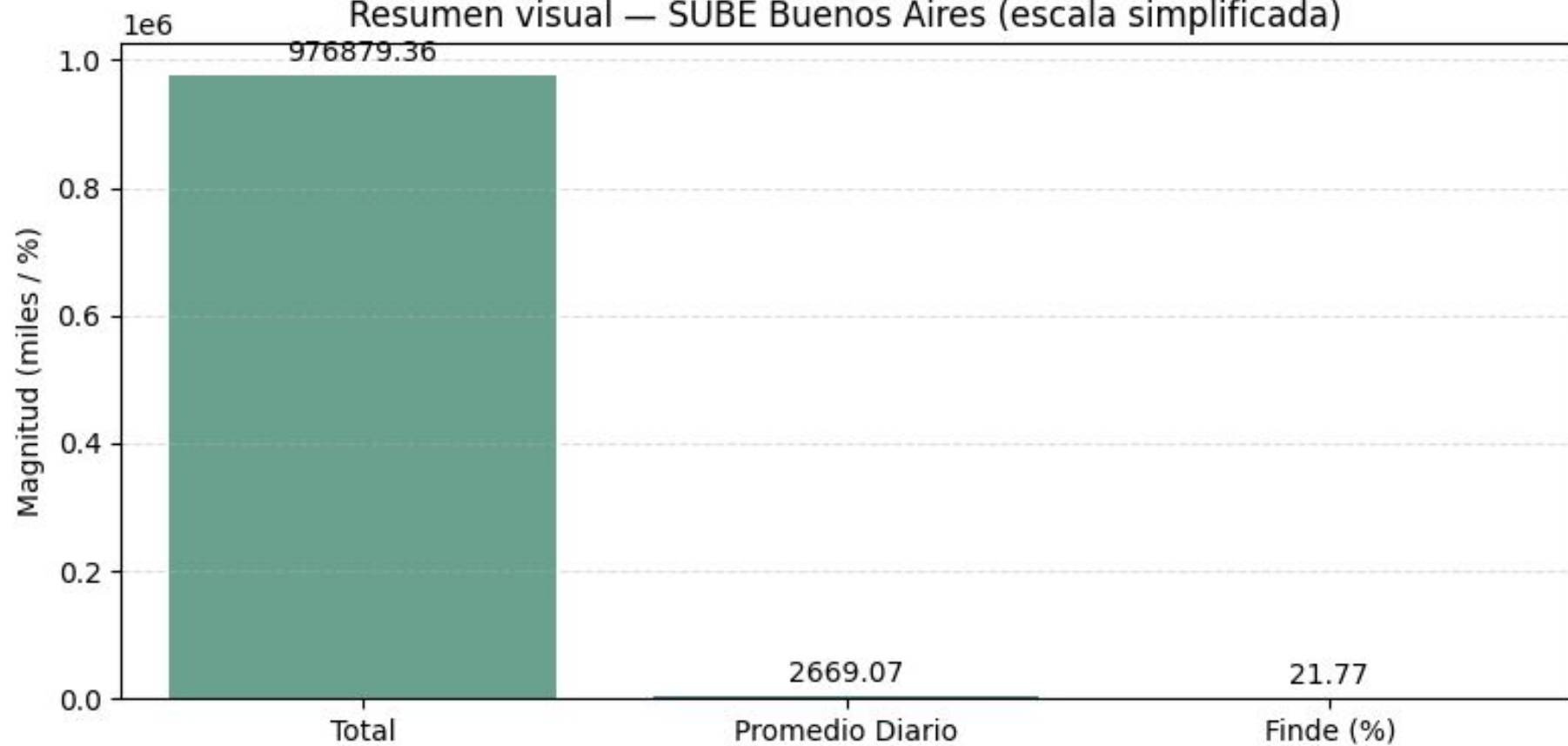
Transacciones SUBE: ajuste proporcional de jueves atípicos



Top 10 empresas por cantidad de transacciones (Buenos Aires)



Resumen visual — SUBE Buenos Aires (escala simplificada)



Modelado del Machine Learning

El Enfoque de Machine Learning

Para predecir las validaciones desde el último día disponible en el dataset, adoptamos un enfoque de Machine Learning utilizando el algoritmo **RandomForestRegressor (RF)**.

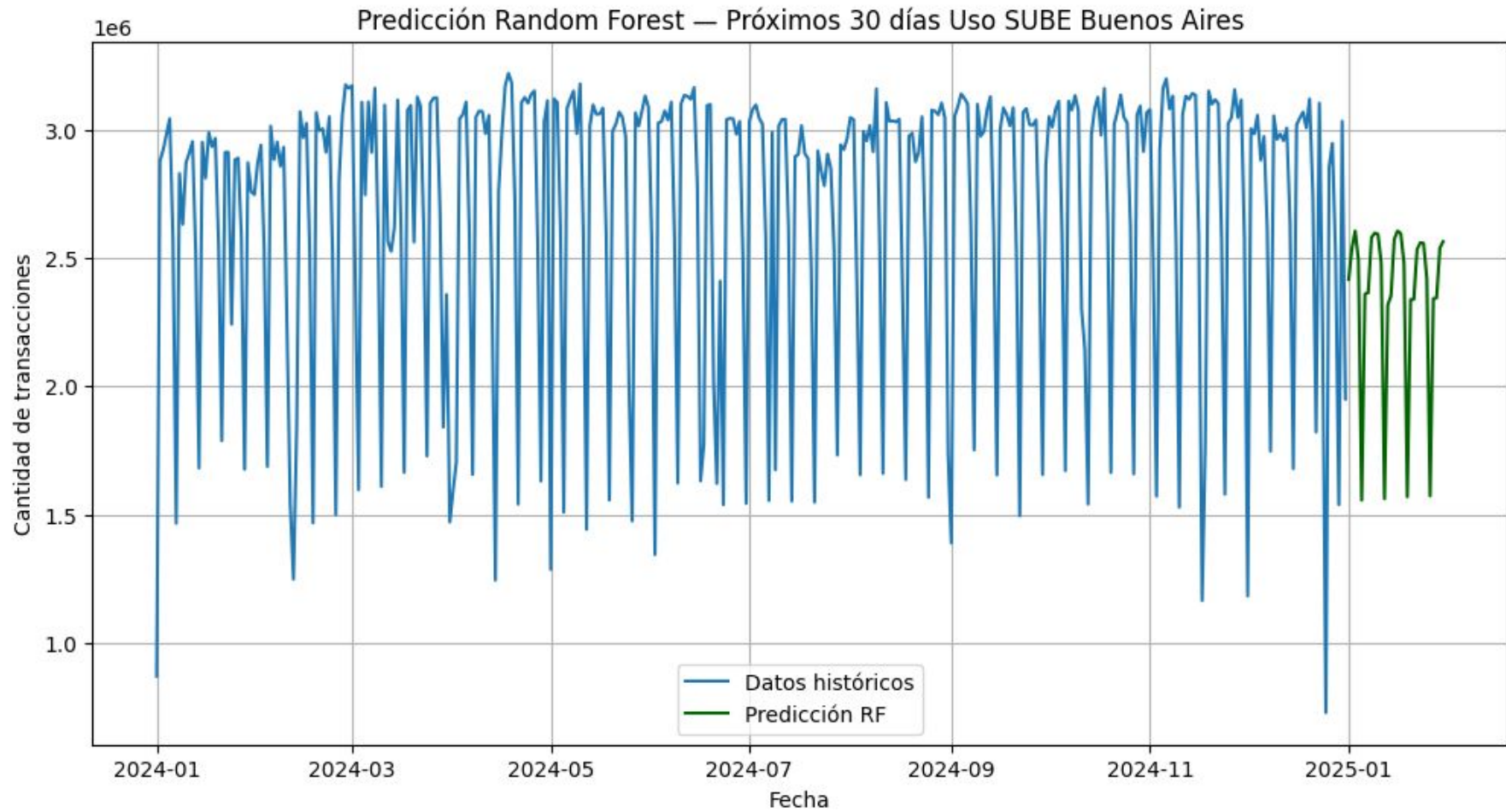
¿Por qué RandomForestRegressor?

- **Robusto a Ruido:** El RF es un ensamble de múltiples árboles de decisión que promedia sus salidas, lo que lo hace **robusto a no linealidades y a pequeñas contaminaciones** en los datos.
- **Captura de Interacciones:** Es excelente para captar interacciones complejas entre variables (por ejemplo, el efecto combinado de ser "lunes" y, a la vez, "vacaciones de invierno").

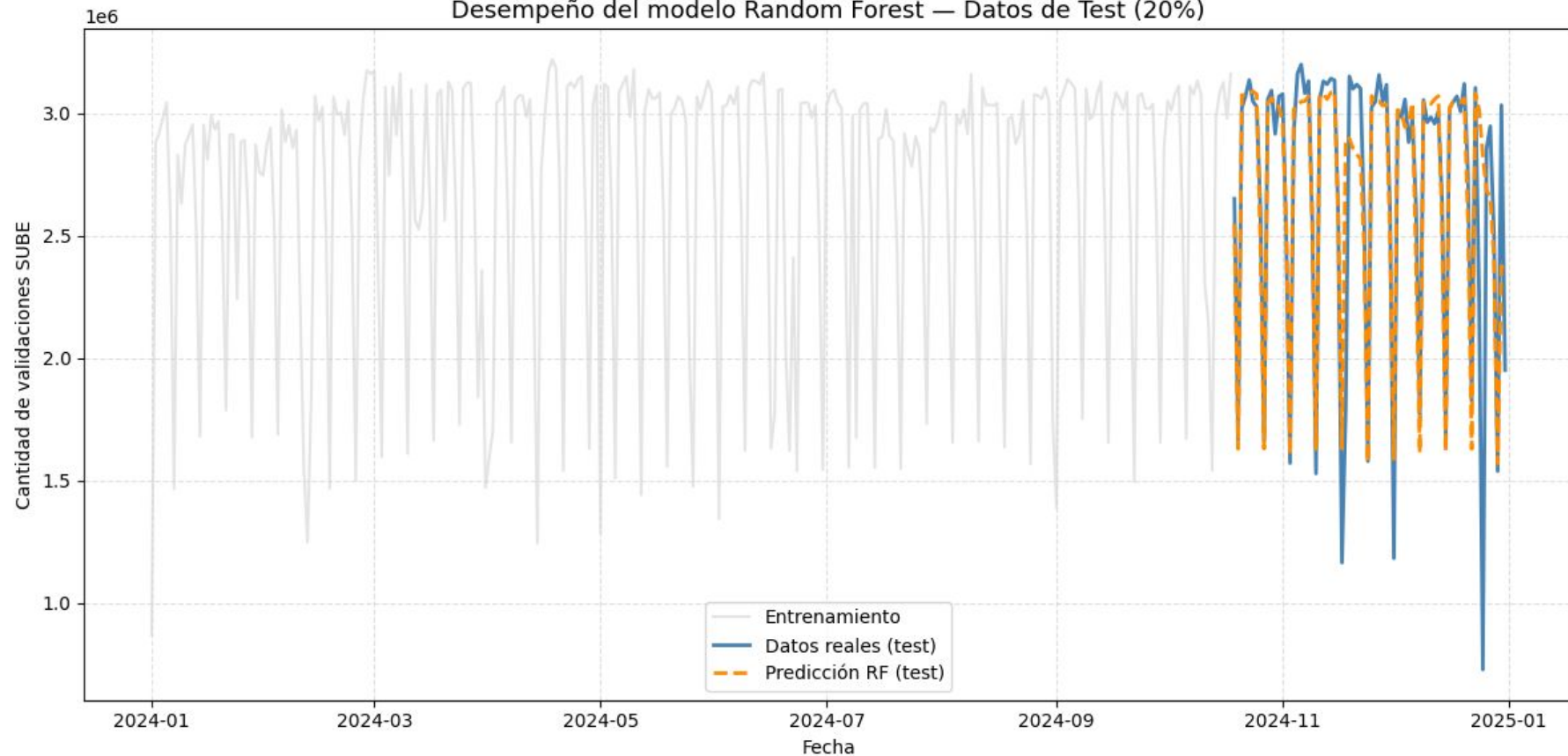
Modelado del Machine Learning

Variables (Features) Predictivas: El modelo se construye a partir de un conjunto de variables diseñadas para capturar la tendencia y la estacionalidad del uso de SUBE:

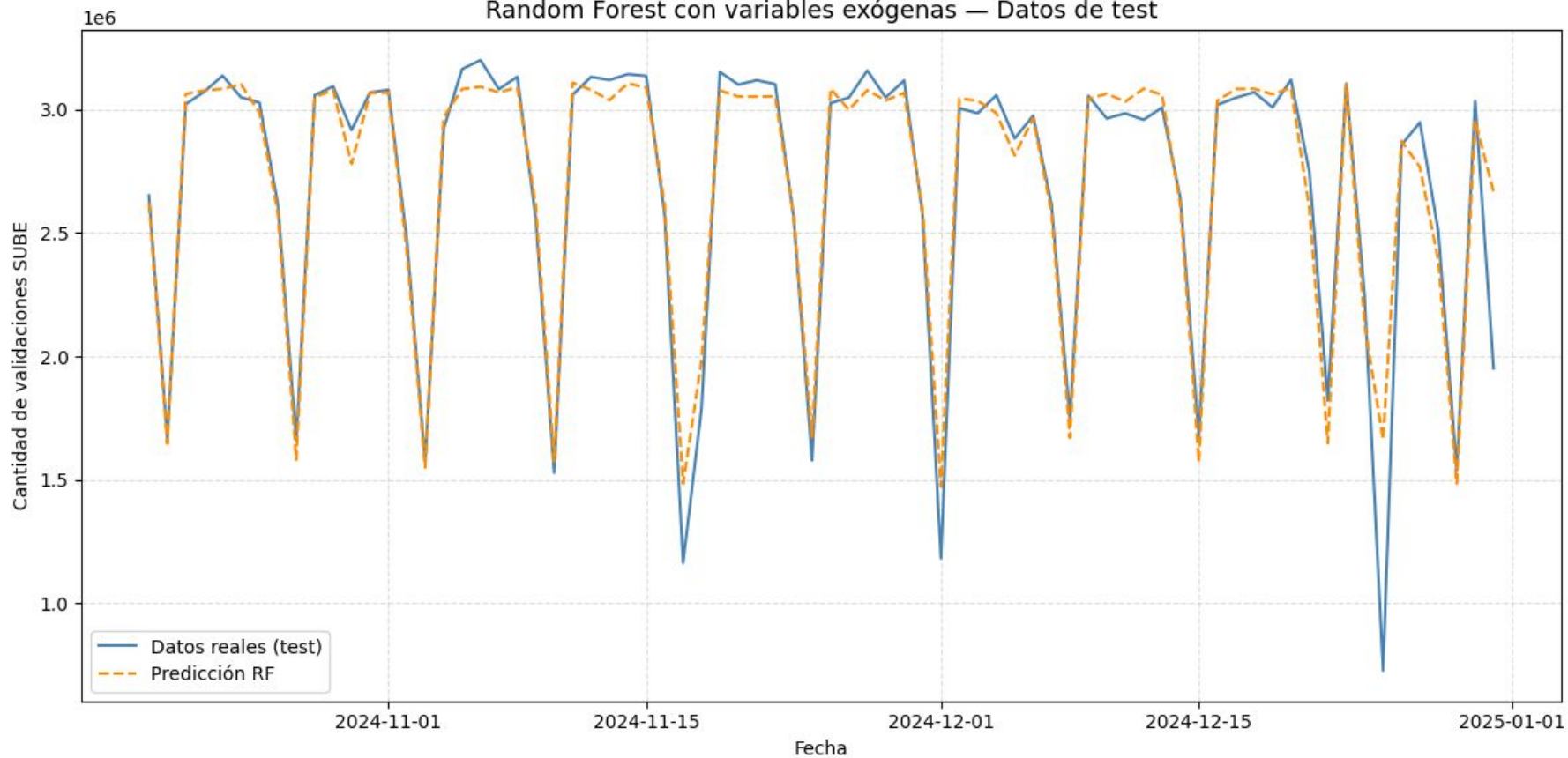
1. **Señales de Calendario:** Día de la semana (**dow**), indicador de fin de semana (**es_finde**), mes, día del mes, y número de semana ISO.
2. **Tendencia:** Contador de tiempo (**dia_num**) y la Media Móvil de 7 días (**rolling_7**).
3. **Shocks Discretos:** Indicadores binarios para **feriados** y **vacaciones de invierno**.



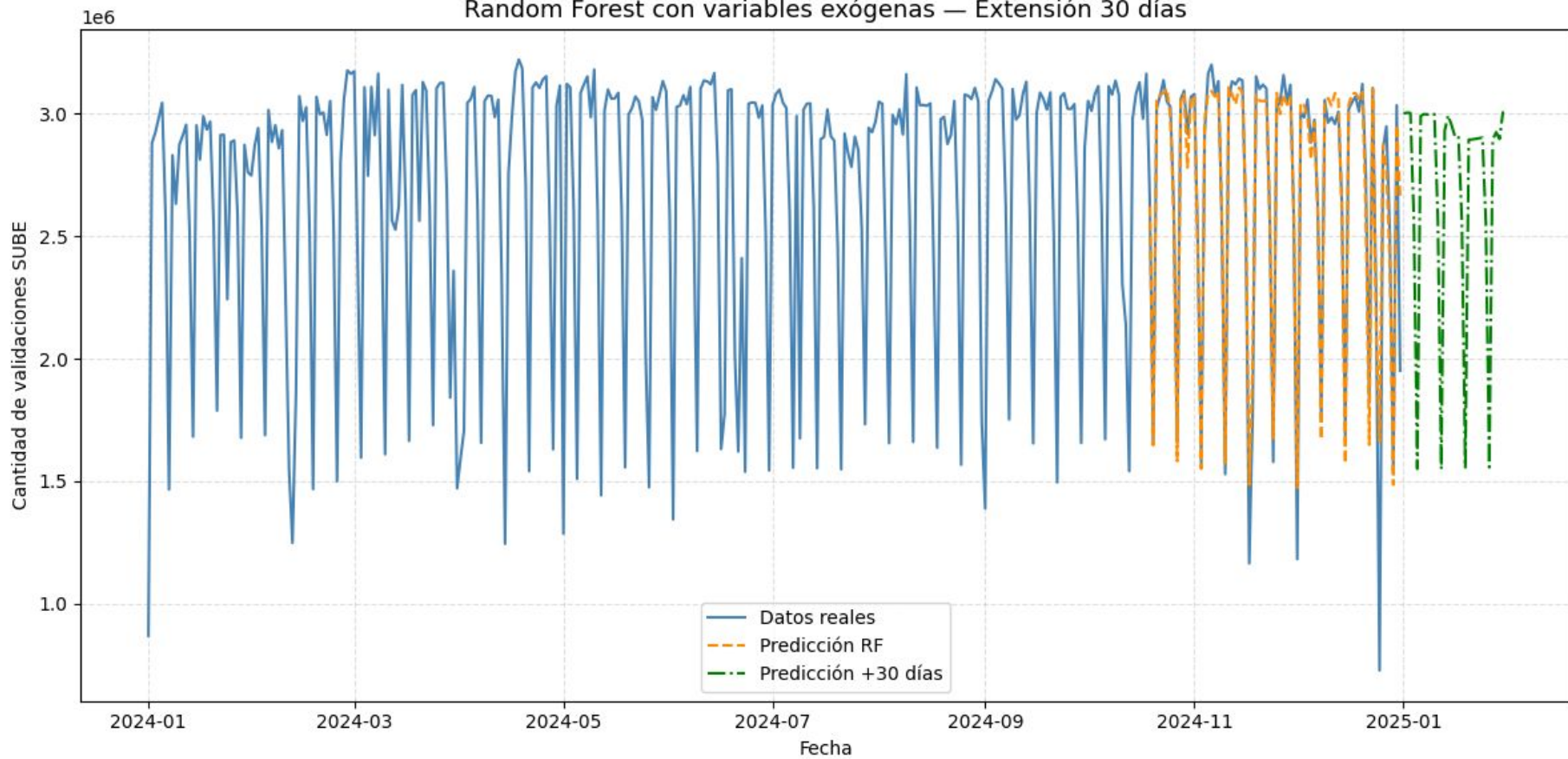
Desempeño del modelo Random Forest — Datos de Test (20%)



Random Forest con variables exógenas — Datos de test

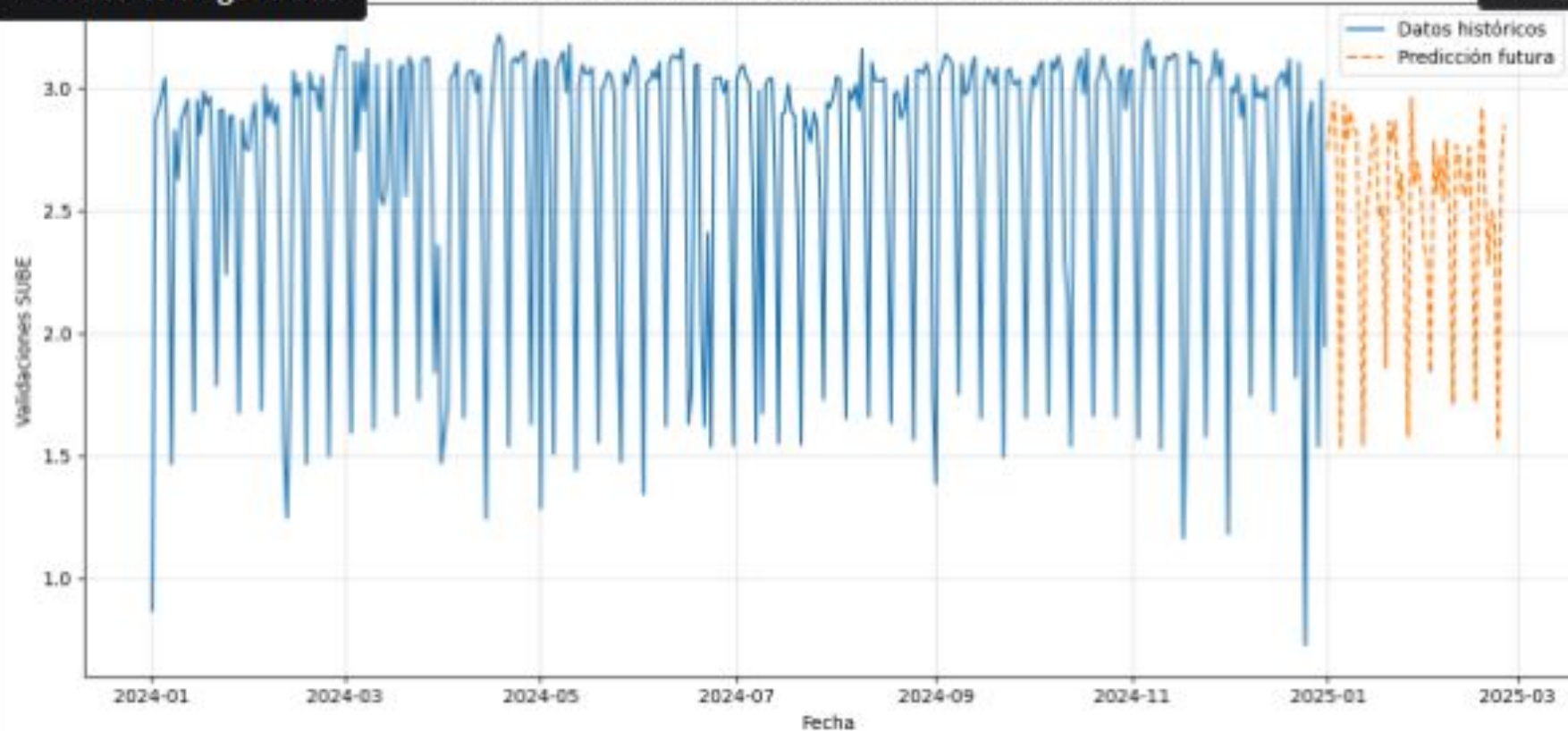


Random Forest con variables exógenas — Extensión 30 días



Predicción generada

Predicción a 56 días (MAE=233,514, RMSE=355,761, $R^2=0.607$)



Hipótesis

La demanda diaria de validaciones SUBE en la Provincia de Buenos Aires está determinada por un patrón temporal estable y predecible, dominado por la estacionalidad semanal (picos en días hábiles y caídas en fines de semana), shocks de calendario (feriados y vacaciones) y variaciones estacionales moderadas a lo largo del año. Este comportamiento se sostiene siempre que el dato sea consistente y depurado y que se consideren explícitamente los eventos de calendario en el análisis.