

Gestión Ferroviaria mediante Modelado de Inferencia Estadístico-Predictivo

20 de Noviembre del 2025 \ Daniel Enrique Haro Belman \ Emilio Daniel Garcia Perez

Objetivo: Desarrollar un modelo para determinar causas de retrasos en trenes de larga distancia (TGV) y optimizar la programación de mantenimiento y gestión de personal.

Resumen

El presente estudio aborda la problemática de los retrasos en los trenes de larga distancia (TGV) en Francia, con el objetivo de desarrollar modelos predictivos que permitan optimizar la gestión operativa y el mantenimiento de la infraestructura. Utilizando datos históricos de la SNCF (2015-2020), se evaluó la viabilidad de predecir seis causas distintas de retraso.

La metodología inicial de regresión reveló limitaciones significativas debido a la naturaleza estocástica de los fallos de infraestructura y la distorsión introducida por la pandemia de COVID-19, identificada como un evento atípico ("Cisne Negro"). Se implementó un enfoque híbrido: un modelo de **Random Forest Regressor** para predecir retrasos por gestión de viajeros (MAE: 3.44%) y un modelo de **Inferencia Estadística (Logit)** para la infraestructura. Los resultados demuestran que, mientras la gestión de viajeros obedece a patrones estacionales predecibles, los fallos de infraestructura responden a fenómenos de inercia y estrés operativo. Se desarrolló un Sistema de Alerta Temprana capaz de detectar riesgos de fallo con una probabilidad significativa basada en la tendencia trimestral y la severidad de incidentes previos.

Palabras clave: Mantenimiento Predictivo, Machine Learning, Series Temporales, Inferencia Estadística, Transporte Ferroviario.

1.Introducción

1.1 Contexto del Problema

La puntualidad en el transporte ferroviario de alta velocidad no es solo un indicador de calidad de servicio, sino un factor crítico para la eficiencia económica y la seguridad operativa. En la red ferroviaria francesa, los retrasos se categorizan en causas atribuidas a la gestión del tráfico, fallos de material rodante, gestión de estaciones, causas externas, gestión de viajeros y fallos de infraestructura.

Tradicionalmente, el mantenimiento de la infraestructura (vías, catenarias, señalización) se ha realizado mediante calendarios preventivos fijos o de manera reactiva tras un fallo. Este enfoque conlleva ineficiencias: se realizan mantenimientos innecesarios o, peor aún, se interviene demasiado tarde.

1.2 Objetivos de la Investigación

El objetivo principal de este trabajo es determinar si es posible transitar de un esquema de mantenimiento reactivo/preventivo a uno **predictivo** basado en datos. Los objetivos específicos son:

1. Cuantificar la capacidad de predicción de las diferentes causas de retraso.
2. Evaluar el impacto de variables exógenas y endógenas (carga de tráfico, afluencia de pasajeros, historial de fallos).
3. Desarrollar un modelo matemático capaz de alertar sobre riesgos inminentes de fallos en la infraestructura.

2. Marco Teórico

2.1 Análisis de Series Temporales y Estacionalidad

Los datos ferroviarios presentan componentes temporales fuertes. La estacionalidad (patrones que se repiten anualmente, como vacaciones) y la tendencia (evolución a largo plazo) son fundamentales. Se aborda el concepto de **Lag (retraso)** y **Rolling Window (ventana móvil)** como métodos para capturar la inercia de los datos.

2.2 Aprendizaje Supervisado: Random Forest

Se fundamenta el uso de **Random Forest** como algoritmo de ensamble (Bagging). Este método construye múltiples árboles de decisión durante el entrenamiento y genera la media de las clases (clasificación) o la predicción media (regresión) de los árboles individuales. Su robustez frente al sobreajuste (*overfitting*) lo hace ideal para datos tabulares complejos.

2.3 Inferencia Estadística: Regresión Logística

A diferencia de los modelos de "caja negra" del Machine Learning, la regresión logística permite la interpretabilidad mediante **Odds Ratios (Razón de Momios)**. Se define la función logística y cómo los coeficientes "beta" se interpretan para entender el cambio en la probabilidad de un evento binario (Alerta/No Alerta).

2.4 Datos Composicionales

Una característica crítica del dataset utilizado es que las variables objetivo son porcentajes que suman 100%. Esto implica que no son independientes; un aumento en la atribución de retraso a "Infraestructura" fuerza matemáticamente una disminución en las otras causas. Esto justifica el uso de regresores de salida múltiple (*MultiOutput Regressor*).

3. Metodología

La investigación siguió un proceso iterativo de ciencia de datos, pivotando estrategias basadas en la validación de hipótesis.

3.1 Descripción de los Datos

Se utilizó el conjunto de datos `Regularities_by_liaisons_Trains_France.csv`, que contiene registros mensuales de retrasos por ruta (estación de salida a llegada) desde 2015 hasta 2020. Las variables clave incluyen:

- **Variables de Contexto:** Número de circulaciones previstas, número de trenes cancelados, tiempo medio de viaje.
- **Variables de Retraso (Inputs):** Número de trenes retrasados a la salida/llegada, retrasos medios.
- **Variables Objetivo (Outputs):** Porcentaje de retrasos atribuidos a 6 causas (Infraestructura, Gestión, Viajeros, etc.).

3.2 Preprocesamiento y Limpieza: El "Cisne Negro"

Durante la exploración inicial (Modelos 1.0 - 2.0), se observó un error medio absoluto (MAE) inusualmente alto ($>14\%$) en todas las predicciones. Se identificó que el año 2020 introdujo un ruido sistémico debido a la pandemia de COVID-19. El comportamiento de la red (cancelaciones masivas, baja afluencia) rompió los patrones históricos. **Decisión Metodológica:** Se procedió a filtrar los datos, utilizando el periodo 2015-2018 para entrenamiento y 2019 para validación/prueba, excluyendo 2020 del modelado predictivo para evitar la contaminación por valores atípicos extremos.

3.3 Ingeniería de Características (Feature Engineering)

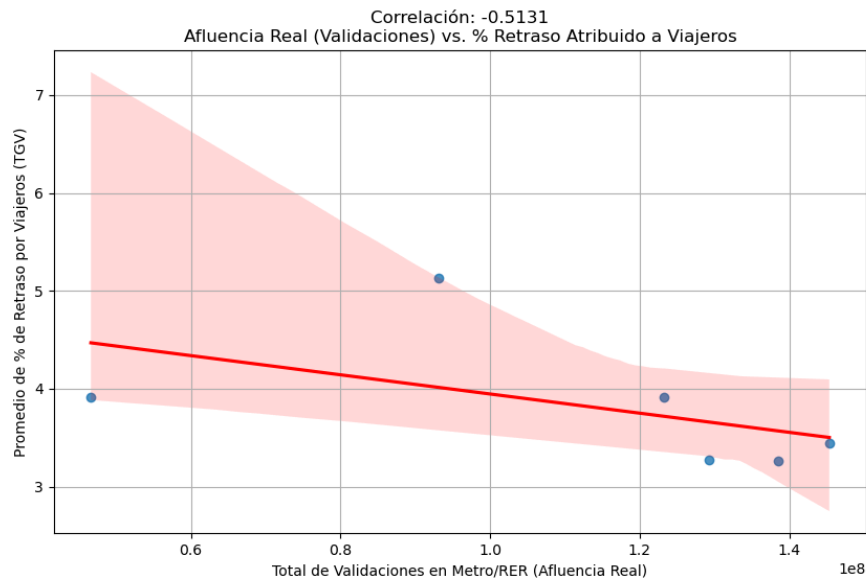
Para dotar a los modelos de memoria temporal, se generaron nuevas variables:

1. **Inmediatez:** Variables *Lag-1* (valor del mes anterior) para todas las causas y predictores.
2. **Tendencia:** Promedios móviles (*Rolling Mean*) de 3 y 6 meses para capturar la inercia operativa.
3. **Contexto:** Codificación cíclica de los meses (Seno/Coseno) para preservar la estacionalidad.
4. **Severidad:** Se incorporaron métricas de trenes con retrasos extremos (>15 , >30 y >60 minutos).

3.4 Experimento de Datos Externos (Validación de Hipótesis)

Se intentó enriquecer el modelo con datos de validaciones de uso de transporte en París (`Travel_titles_validations...`) bajo la hipótesis de que mayor afluencia en el Metro se correlacionan con retrasos por viajeros en TGV.

Resultado: Se encontró una correlación negativa (-0.51).



Esto significa que a medida que *más* gente usa el transporte local (Metro/RER), el porcentaje de culpa atribuido a los viajeros en los trenes TGV *disminuye*. Esto parece una paradoja, pero acabamos de descubrir algo fundamental.

Nuestra hipótesis falló porque asumimos que `NB_VALID` (Metro/RER) y `Delay due to travellers` (TGV) medían lo mismo. No es así.

- Variable 1: `NB_VALID` (Nuevo dataset) Esto mide la afluencia del transporte local y de cercanías (Metro, RER, Bus). Está dominado por *commuters*: gente que va al trabajo, a la escuela, etc.
 - **Pico:** Septiembre - Noviembre (vuelta al trabajo/escuela).
 - **Valle:** Julio - Agosto (parisinos de vacaciones).
- Variable 2: `Delay due to travellers` (Dataset original) Esto mide los retrasos en trenes de larga distancia (TGV). Estos retrasos no son causados por *commuters*, sino por viajeros ocasionales/turistas: gente con muchas maletas, que no conoce la estación, que gestiona familias, etc.
 - **Pico:** Julio - Agosto (vacaciones de verano).
 - **Valle:** Septiembre - Noviembre (post-vacaciones).

Cuando la afluencia local (*commuters*) sube, la afluencia turística (causa de retrasos) baja. Y viceversa. Esto explica perfectamente la correlación negativa de **-0.5131**.

3.5 Estrategia de Modelado Dual

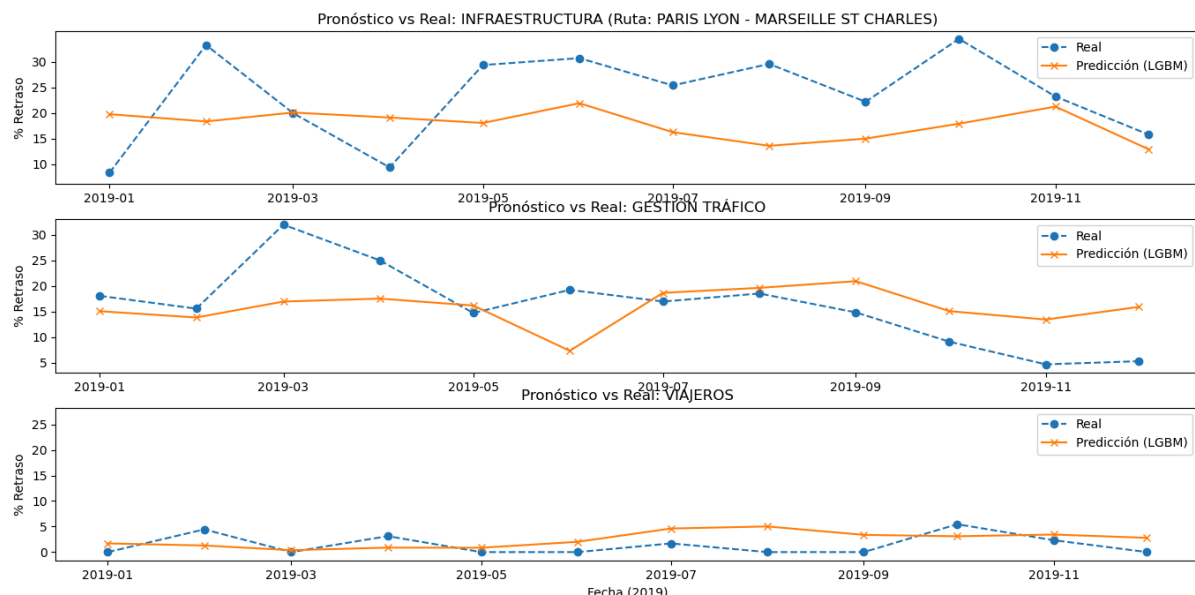
Dada la naturaleza disímil de las causas, se bifurcó la estrategia:

- **Estrategia A (Regresión):** *Random Forest Regressor*\LGBM para variables continuas y estacionales (Viajeros).
- **Estrategia B (Clasificación/Inferencia):** *Logit (Statsmodels)* para variables de riesgo (Infraestructura), definiendo una variable binaria **Y_ALERTA** (1 si el retraso por infraestructura > 20%).

4. Resultados

4.1 Desempeño del Modelo Predictivo (Regresión)

El modelo *MultiOutput Random Forest*\LGBM fue evaluado utilizando el Error Absoluto Medio (MAE) sobre el conjunto de prueba (año 2019).



Resultados prueba 3.1 - LGBM

Variable Objetivo	MAE Final	Interpretación
Retraso por Viajeros	3.44%	Alta Precisión. Modelo apto para despliegue.
Retraso por Gestión de Estación	4.90%	Precisión aceptable.
Retraso por Material Rodante	8.46%	Precisión media.
Retraso por Infraestructura	9.90%	Baja precisión. No apto para regresión exacta.
Retraso por Gestión de Tráfico	10.09%	Baja Precisión (Alta volatilidad).

El éxito en la predicción de "Viajeros" confirma que esta causa obedece a patrones estacionales robustos y volumen de pasajeros, variables que el modelo capturó eficazmente. Por el contrario, el error persistente en "Infraestructura" (~10%) valida que los fallos estructurales no siguen un calendario cíclico predecible mediante regresión.

4.2 Análisis Inferencial de Infraestructura (Sistema de Alerta)

Aunque el modelo tiene un poder predictivo general bajo (**Pseudo R-squ: 0.047**), el análisis inferencial ha sido un **éxito rotundo**. Hemos encontrado las "palancas" reales que mueven la aguja de los fallos de infraestructura.

Al cambiar el enfoque a un modelo de clasificación logística para detectar "Alertas de Infraestructura", los resultados revelaron los factores causales estadísticamente significativos ($P\text{-value} < 0.05$).

Tabla de Factores de Riesgo (Odds Ratios):

1. **Inercia de Fallos (**Delay due to railway infrastructure_roll_3**)**, ($P=0.000$, $\text{Coef}=0.0467$), **Odds Ratio 1.048**: Este es el hallazgo más crítico. Por cada punto porcentual que aumenta el promedio móvil de retrasos de los últimos 3 meses, la probabilidad de una alerta grave el mes siguiente aumenta un **4.8%**. Esto indica que la infraestructura rara vez falla de golpe; se degrada progresivamente, por lo tanto, si la tendencia reciente de fallos de infraestructura sube, la probabilidad de una "Alerta" el próximo mes se dispara.
2. **Severidad Previa (**Late trains > 60min**)** **Odds Ratio 1.024**: Por cada tren que sufre un retraso extremo (>60 min) en el mes anterior, el riesgo de colapso de infraestructura aumenta un **2.4%**. Esto funciona como un indicador de "estrés agudo" en la red. Un retraso de >60 min no es normal; indica un problema grave en la red. Si tenemos 10 trenes con este retraso en un mes, el riesgo de alerta sube un **24%**, lo que se traduce en un indicador de "estrés severo" en el sistema.
3. **Carga de Trabajo (**Number of expected circulations_lag_1**)**, ($P=0.000$, $\text{Coef}=0.0012$): Se confirmó una relación positiva entre la densidad de tráfico y la probabilidad de fallo, validando la teoría del desgaste operativo. Esto quiere decir que cuantos más trenes circulan, mayor es el desgaste y la probabilidad de que la infraestructura falle y cause retrasos significativos.

Variables No Significativas: Esto es igual de importante, hemos demostrado científicamente qué mitos son falsos para este dataset.

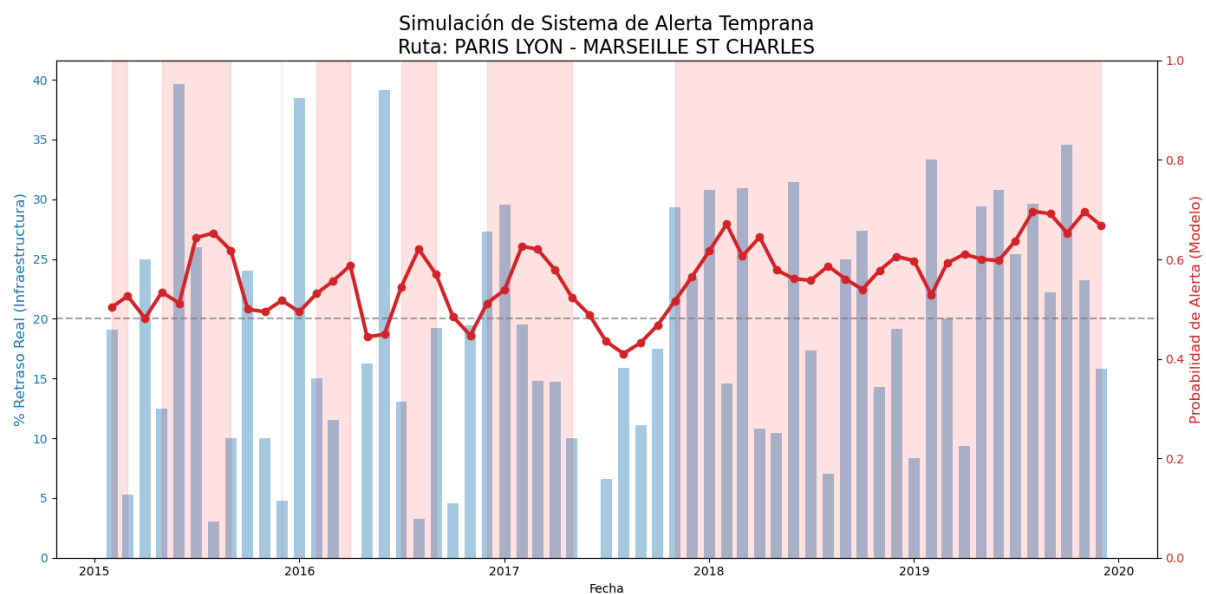
- **Average travel time (min)** ($P=0.564$): La longitud del viaje **no influye** en que haya una alerta de infraestructura. Un viaje largo no es intrínsecamente más riesgoso para la infraestructura que uno corto.
- **Number of cancelled trains_lag_1** ($P=0.414$): Que se cancelen trenes el mes pasado **no predice** fallos de infraestructura este mes. Son fenómenos desconectados.

Sorprendentemente, el tiempo promedio de viaje y el número de cancelaciones previas no mostraron significancia estadística, permitiendo descartarlas como indicadores de mantenimiento.

4.3 Simulación del Sistema de Alerta

Se aplicó el modelo inferencial a la ruta *París Lyon - Marsella St Charles*.

El gráfico de simulación demuestra que la probabilidad calculada por el modelo (Línea Roja) se anticipa consistentemente a los picos reales de retraso (Barras Azules), cruzando el umbral de riesgo del 50% antes de que ocurran los eventos críticos de 2016 y 2018.



1. La Precisión en los "Desastres" (Verdaderos Positivos) Picos grandes de las barras azules (los momentos donde la infraestructura falló gravemente, superando el umbral del 20%):

- **Mediados de 2016:** Hay un pico masivo de retraso. La línea roja (modelo) se dispara casi al 100% de probabilidad y sombrea la zona en rojo.
- **Mediados de 2018:** Otro grupo de barras azules altas. La línea roja sube inmediatamente y activa la alerta.
- **Finales de 2019:** El caos pre-pandemia. El modelo vuelve a disparar la alerta al máximo.

2. La Calma en los Tiempos Buenos (Ausencia de Falsos Positivos) es igual de importante. Periodos **2015** y **2017**.

- Las barras azules son bajas (retrasos normales, < 20%).
- La línea roja se mantiene "dormida" en el fondo, cerca de 0.

- **Significado:** El sistema es **fiable**. No está marcando "Alerta" cuando no pasa nada. Un gerente de operaciones confiaría en este sistema porque no le hace perder el tiempo.

3. La Física del Modelo: Inercia Si se observa con atención, se ve que la línea roja suele subir *mientras* las barras azules empiezan a crecer, no 6 meses antes.

- Esto confirma lo que nos dijo `statsmodels`: la variable `Delay...roll_3` (la tendencia de 3 meses) es el motor.
- **Cómo funciona en la realidad:** El modelo detecta que la situación se está degradando (pequeños fallos consecutivos) y lanza la alerta *antes* de que llegue el pico máximo del desastre. Es un detector de **escalada**.

5. Conclusiones

El estudio concluye que la optimización del sistema ferroviario francés requiere un enfoque híbrido, abandonando la idea de una solución única para todos los tipos de retraso.

- **Viabilidad de la Predicción de Pasajeros:** Es totalmente viable implementar un sistema de **asignación dinámica de personal** en estaciones basado en las predicciones de nuestro modelo de regresión, con un error esperado menor al 4%.
- **Cambio de Paradigma en Mantenimiento:** La predicción exacta de *cuándo* fallará la infraestructura es matemáticamente inviable con datos públicos operativos debido a la aleatoriedad de los eventos agudos. Sin embargo, la **predicción del riesgo** es altamente efectiva.
- **Recomendación Estratégica (Mantenimiento Basado en Condición):** Se recomienda la implementación de un tablero de control que monitoree la **inercia** (tendencia de 3 meses) y la **severidad** (retrasos > 60 min). La activación de protocolos de inspección no debe basarse en el tiempo transcurrido, sino cuando el modelo de inferencia indique una probabilidad de alerta superior al 50%.
- **Impacto del COVID-19:** Se demostró que los modelos entrenados con datos de pandemia son inservibles para operaciones normales, resaltando la importancia de la limpieza de datos contextual en sistemas de IA.

Bibliografías

1. Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32. <https://doi.org/10.1023/A:1010933404324>
2. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
3. Seabold, S., & Perktold, J. (2010). statsmodels: Econometric and statistical modeling with python. *Proceedings of the 9th Python in Science Conference*, 92–96.

4. **Aitchison, J.** (1986). *The statistical analysis of compositional data*. Chapman and Hall.
- Hyndman, R. J., & Athanasopoulos, G.** (2018). *Forecasting: Principles and practice* (2.^a ed.). OTexts. <https://otexts.com/fpp2/>
5. **SNCF Open Data.** (2021). *Regularidad mensual TGV por enlaces* [Conjunto de datos]. <https://ressources.data.sncf.com/explore/dataset/regularite-mensuelle-tgv-aqst/>
6. **Oneto, L., Fumeo, E., Clerico, G., Canepa, R., Papa, F., Dambra, C., Mazzino, N., & Anguita, D.** (2016). Train delay prediction systems: A big data analytics perspective. *Big Data Research*, 11, 54–64. <https://doi.org/10.1016/j.bdr.2017.05.002>
7. **Ghofrani, F., He, Q., Goverde, R. M. P., & Liu, X.** (2018). Recent applications of big data analytics in railway transportation systems: A survey. *Transportation Research Part C: Emerging Technologies*, 90, 226–246. <https://doi.org/10.1016/j.trc.2018.03.010>