

# Reto HidroPredictiva

*Equipo Overflow (4)*

# Objetivo

*Programa Nacional de Algoritmos Verdes (PNAV)* presenta el desafío de desarrollar modelos eficientes de predicción hidráulica.

Pronósticos para 5 puntos de la cuenta del Duero a 24 y a 48 horas.

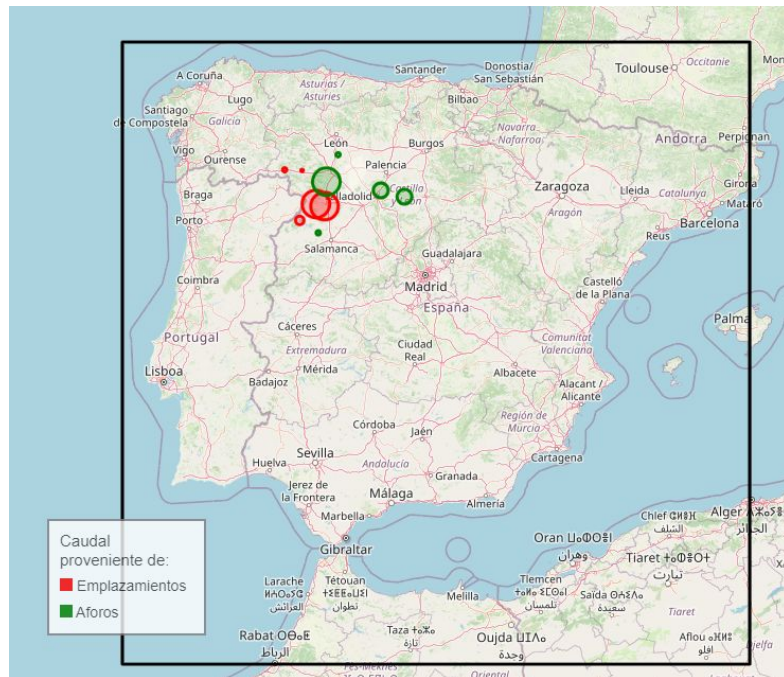


# Problema

- Volumen de datos elevados
- Eficiencia recursos computacionales.
- Precisión en las predicciones para planificación estratégica y operación de la cuenca hidrográfica.



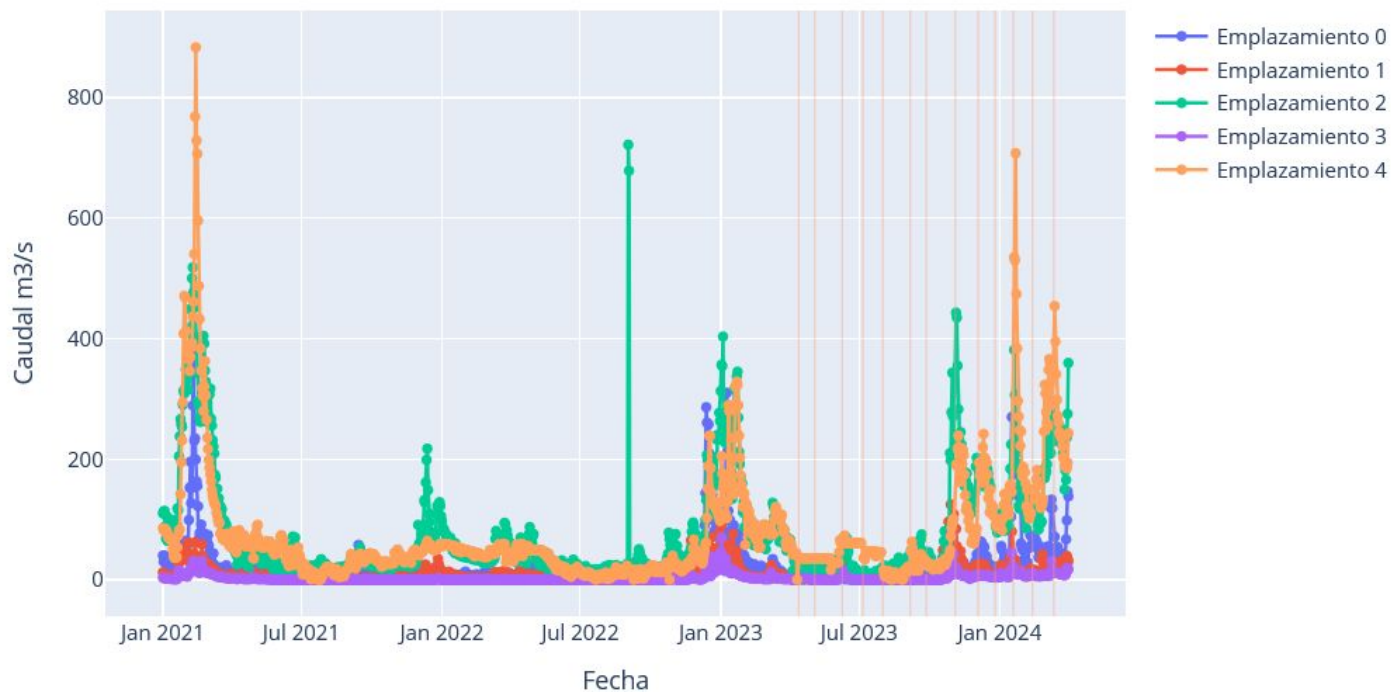
# Contexto



Datos aforo + GFS → Caudal emplazamiento

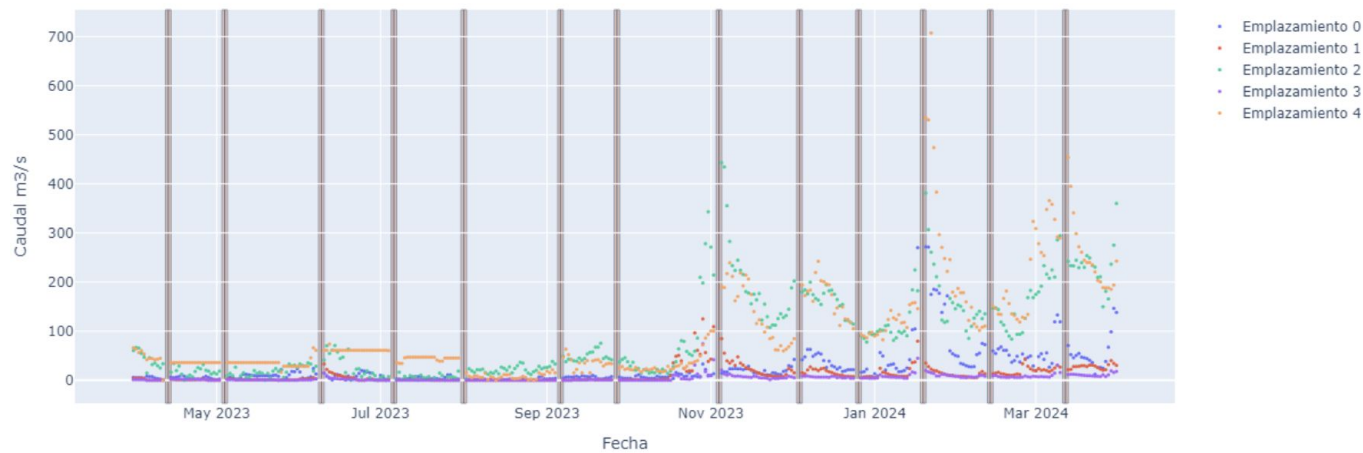
# Contexto

Target Set



# Contexto

Train Set del último año - las bandas señalan las fechas del Test Set



# Nuestra solución

2 (24 y 48h) modelo XGBoost por desplazamiento.

- Selección de features más importantes por desplazamiento.
- Lags para variables importantes.

¿Por qué XGBoost?

- Rápido
- Híbrido entre random forest y gradient boosting
- Redes neuronales (LSTM, transformers):
  - Datos insuficientes
  - Mayor complejidad computacional y espacial



# Preprocesamiento y selección de variables

Objetivo: Modelo más sencillo posible pero que mantenga precisión.

Paso 1:

Análisis de correlación entre features e  $y$  → 4 variables más importantes

Paso 2:

Análisis de correlación tamaño lag óptimo → hasta 4 días

Paso 3:

Añadir variables temporales

¿Por qué hemos elegido este método para reducción de dimensionalidad?

- Explicabilidad y control.

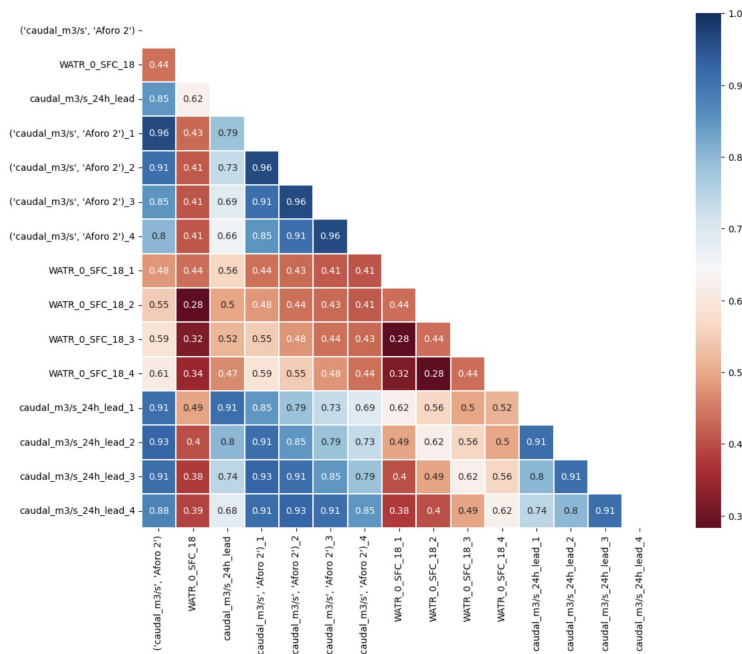




# Análisis de correlación

## Emplazamiento 0

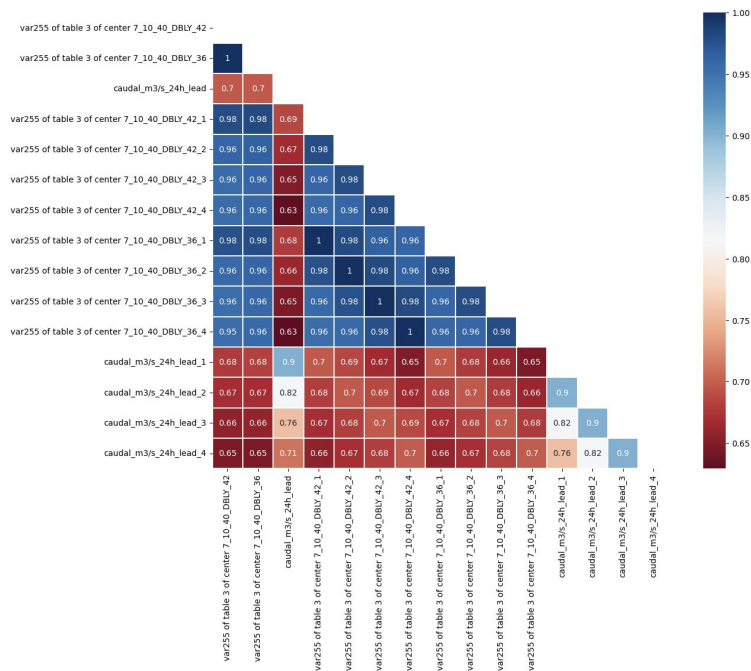
\*Hemos eliminado del modelo las referencias a los valores del caudal



# Análisis de correlación

## Emplazamiento 1

\*Hemos eliminado del modelo las referencias a los valores del caudal



# Análisis de correlación

Emplazamiento i



# Testing y validación

Validación cruzada:

- Time Series Split
- 5 folds
- Aprovechamos todo el dataset



# Código eficiente

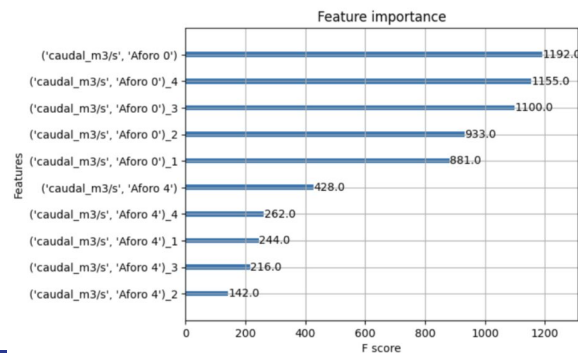
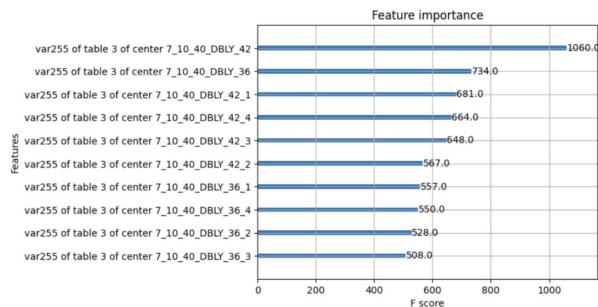
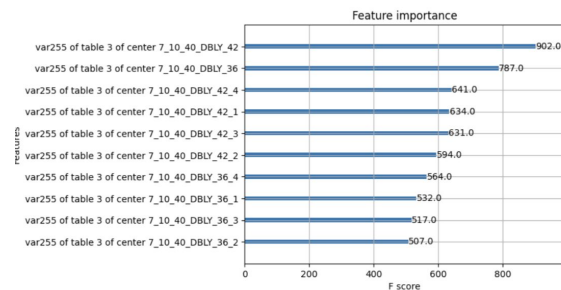
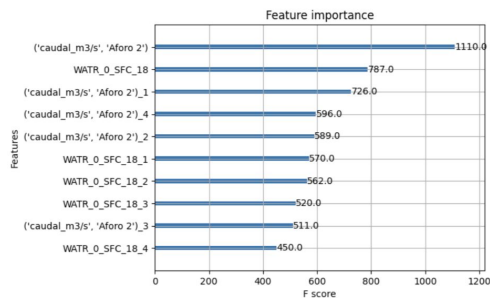
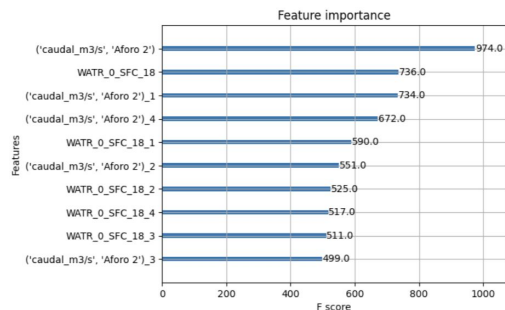
Hemos empleado librerías optimizadas a bajo nivel como:

- XGBboost
- pandas
- numpy
- ...



# Beneficios

## Sacamos máximo provecho a las variables: F score



# Beneficios

- **0.003544** kWh de consumo de electricidad (incluyendo código que análisis).
- **0.002089** g.CO<sub>2</sub>eq/s estimación media anual 65.8649 kg.CO<sub>2</sub>eq/año.
- Sin selección de variables y lags (370%):
  - 0.013351 kWh
  - 0.014929 g.CO<sub>2</sub>eq/s anual



# Resultados

- MAE: 29 (Sin el uso de la variable del caudal de los emplazamientos)

Metricas del conjunto de validación:

Emplazamientos	MAE caudal_m3/s_24h	MAE caudal_m3/s_48h
Emplazamiento 0	9.7629	9.8054
Emplazamiento 1	3.8009	3.9101
Emplazamiento 2	20.50	23.33
Emplazamiento 3	1.5718	1.0964
Emplazamiento 4	22.700	23.60



# Mejoras

- Entrenar con más datos
  - Considerar franjas del testset
- Crear modelos específicos para franjas irregulares.
- Probar juntar modelos.
- Seguir experimentando con hiperparametros.
- Analizar otras métricas de ajuste

