

Proyecto Final

Luis Daniel Casais Mezquida

Análisis Inteligente de Datos 24/25

Universidad Carlos III de Madrid

homo
homini
SACRA
RES

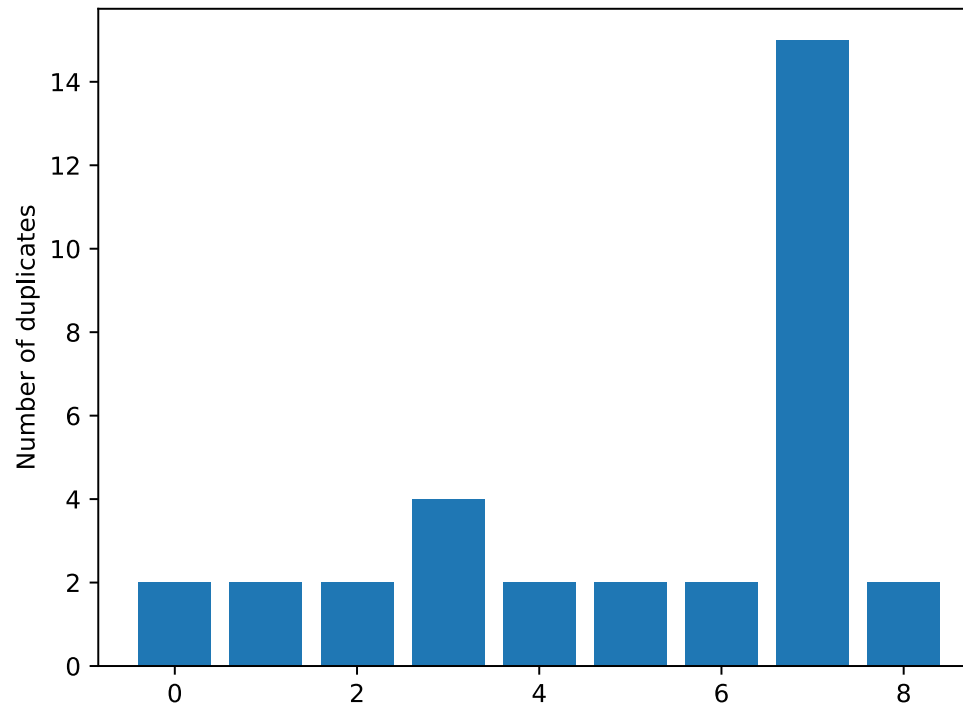
Dataset

["Obesity Prediction Dataset"](#), por Stephen Adeniran (2025)

- Estimar nivel de obesidad
- Hábitos alimenticios, físicos, y médicos
- México, Perú, y Colombia

Dataset

- 2112 observaciones
- 17 características
 - categóricas (5), numéricas (8), booleanas (4)
- 0 celdas vacías, 9 (33) duplicados (1.56%)



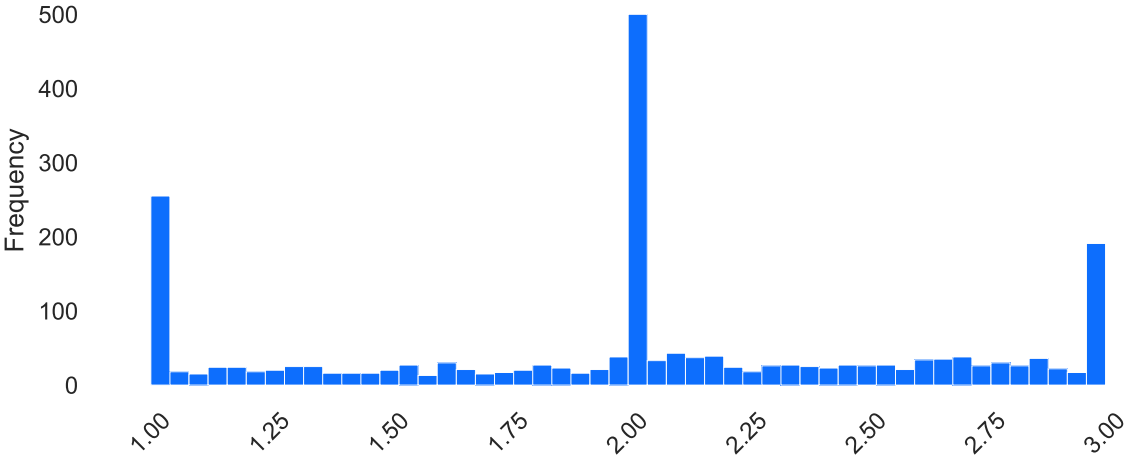
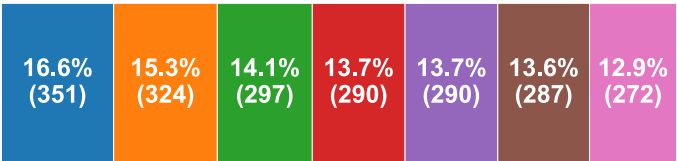
Dataset

Gender	Male
Age	21
Height	1.62
Weight	70.0
family_history_with_overweight	no
FAVC	yes
FCVC	2.0
NCP	1.0
CAEC	no
SMOKE	no
CH2O	3.0
SCC	no
FAF	1.0
TUE	0.0
CALC	Sometimes
MTRANS	Public_Transportation
NObeyesdad	Overweight_Level_I
# duplicates	15

Análisis univariable

- Características muy desbalanceadas (>80%): *SMOKE*, *family_history_with_overweight*
- Característica objetivo balanceada
- Características de rangos numéricos con "ruido"

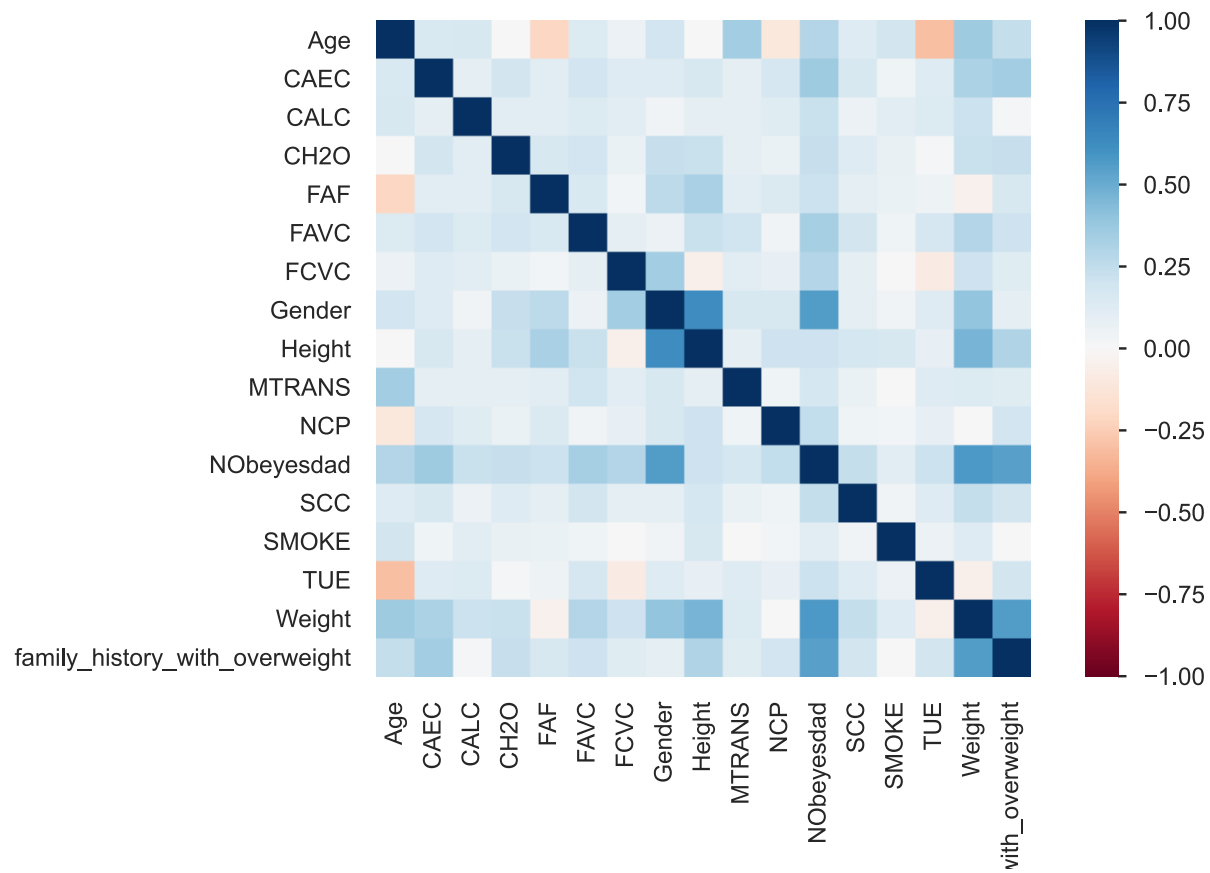
Dataset

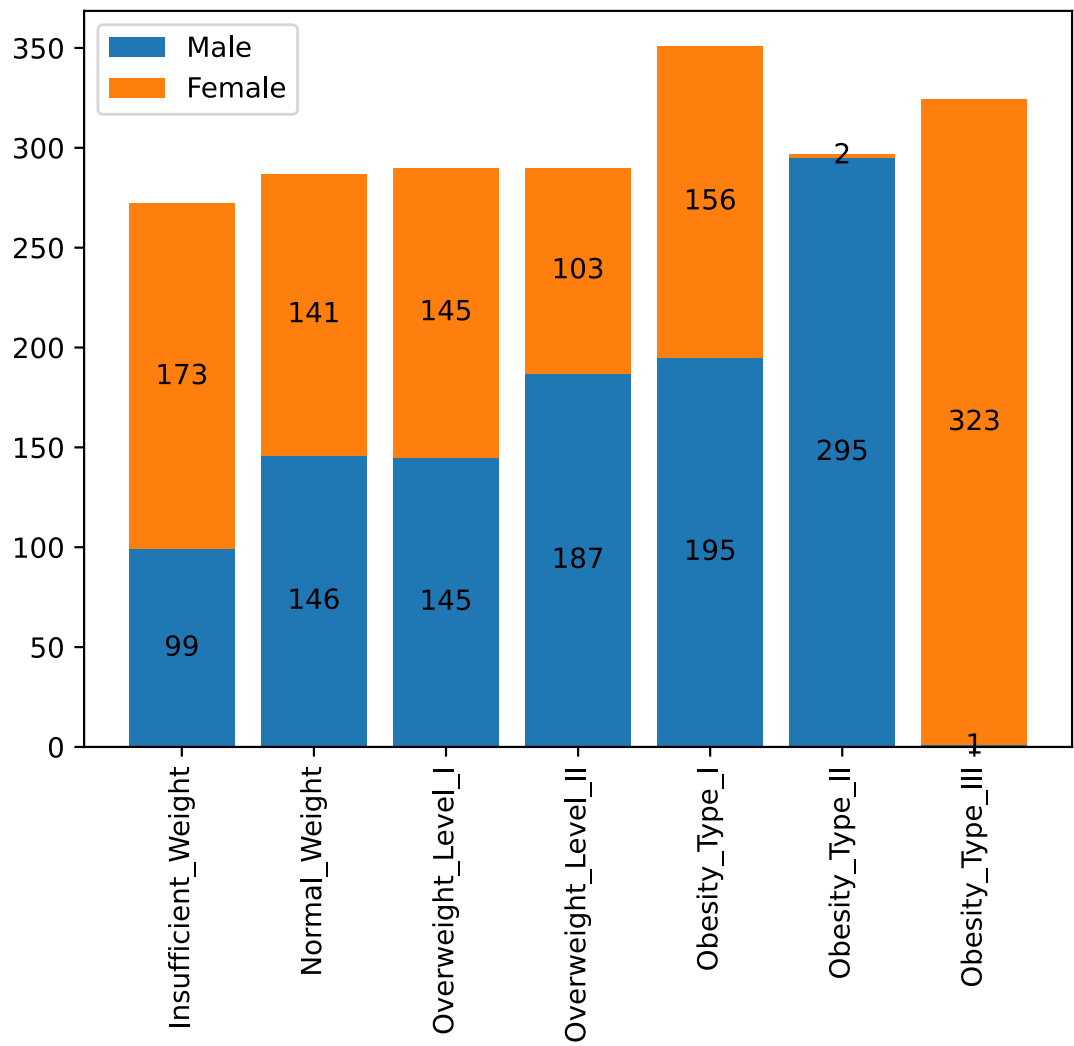


Análisis multivariable

- Alta correlación entre género, peso, e historial familiar con obesidad
- Obesidad por género:
 - Tipo II: apenas mujeres
 - Tipo III: apenas hombres

Dataset





Metodología

Métrica de evaluación: F1-score, macro

Preprocesado:

1. Redondear rangos numéricos a entero
2. Escalado para variables numéricas
3. Encoding para categóricas

Experimentación

- 20% test
- *Cross Validation* de 5
- Ajuste hiperparámetros

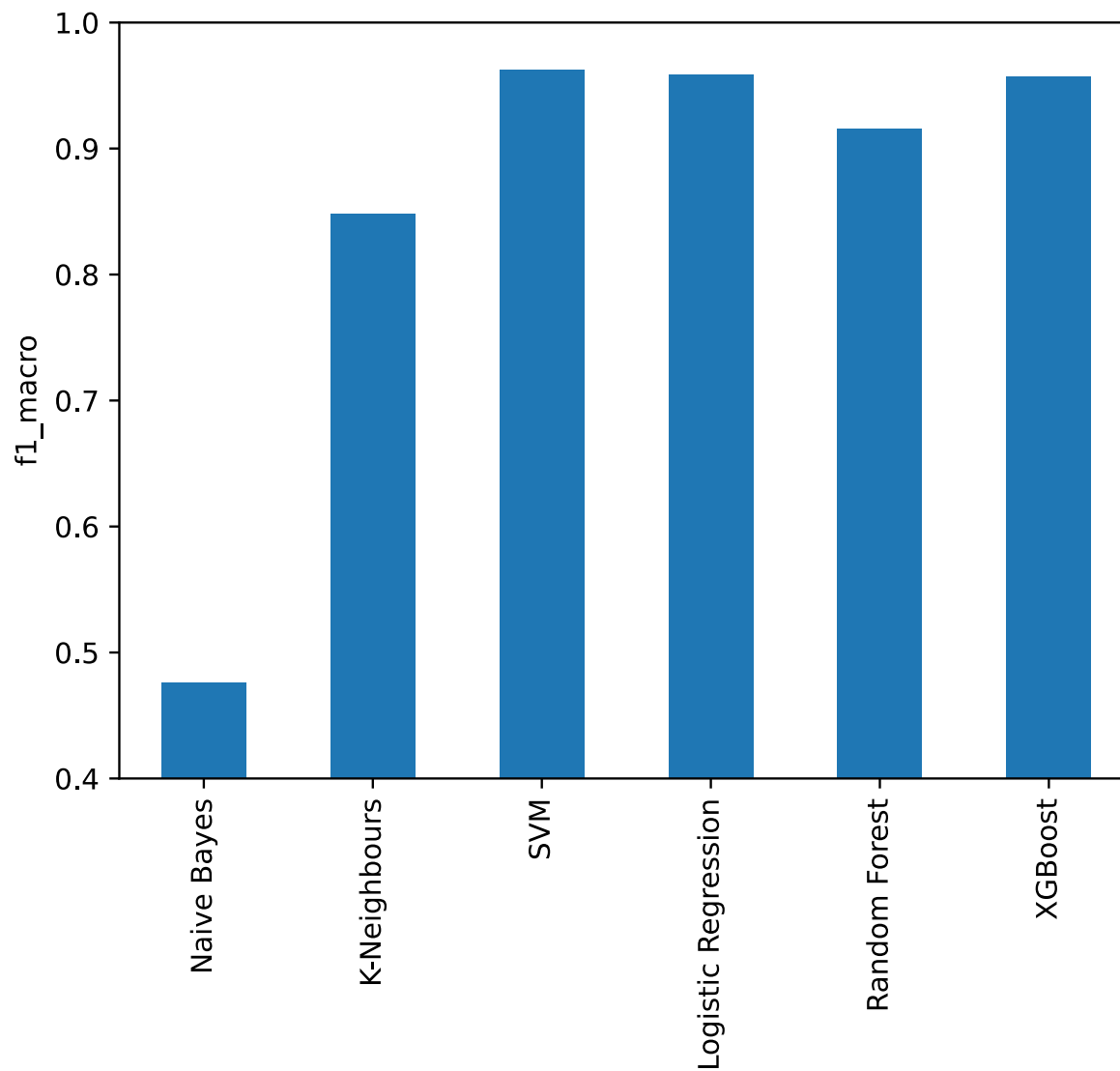
Modelos:

1. Base

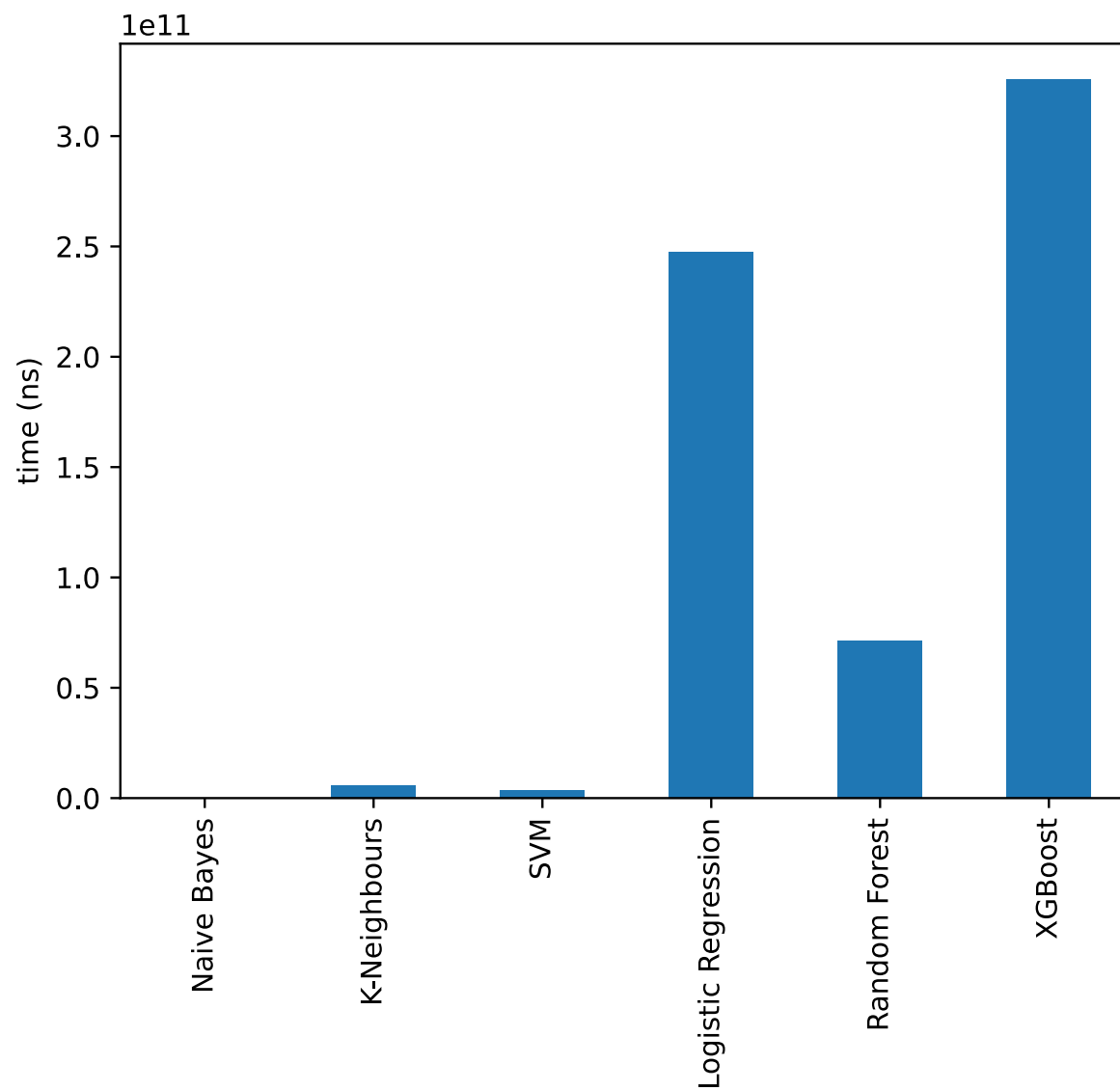
- *Naive Bayes*
- *K-Neighbours*
- *SVM*
- *Logistic Regression*
- *Random Forest*
- *XGBoost*

2. *Stacking*

Modelos base



Experimentación

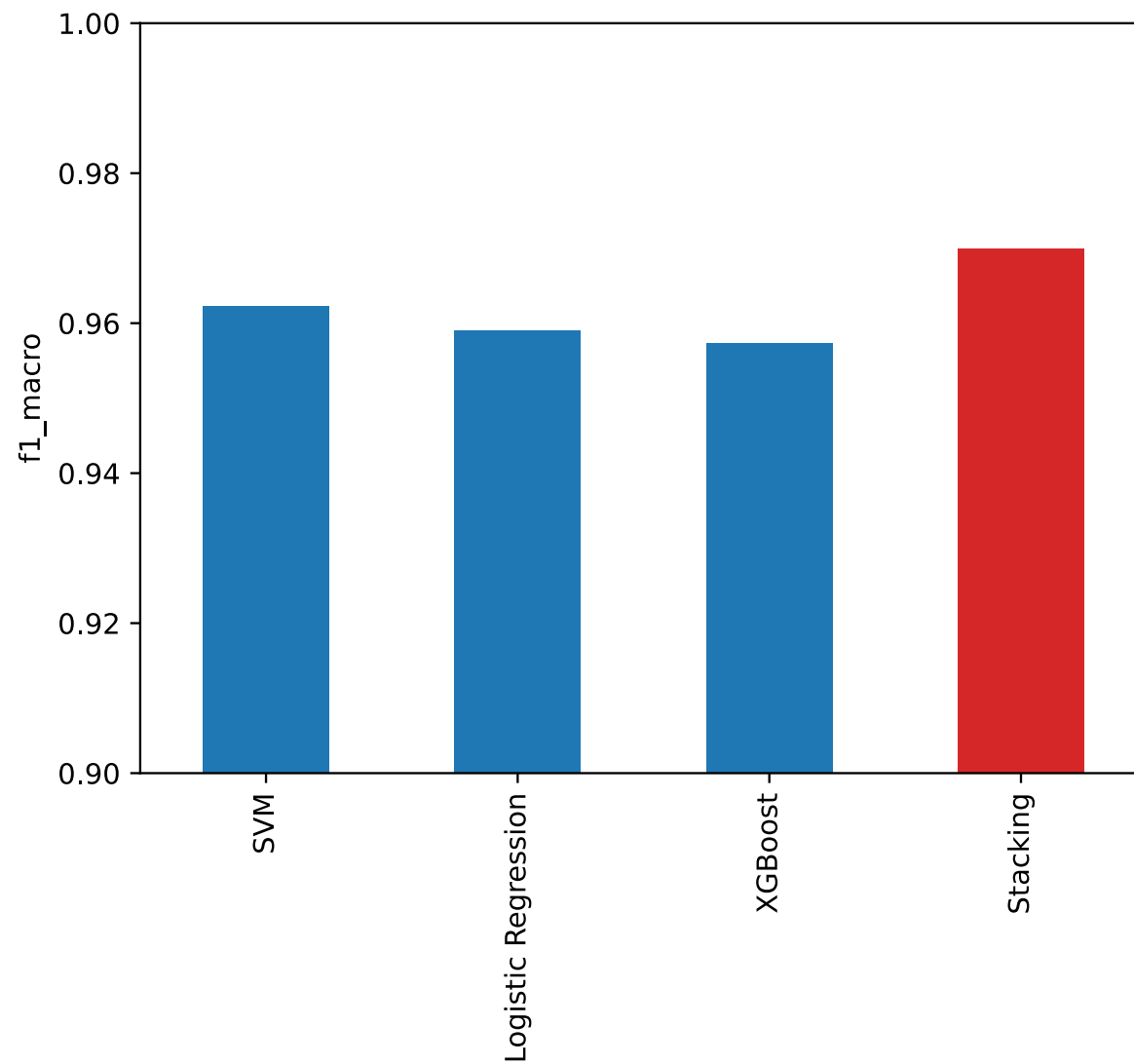


Stacking

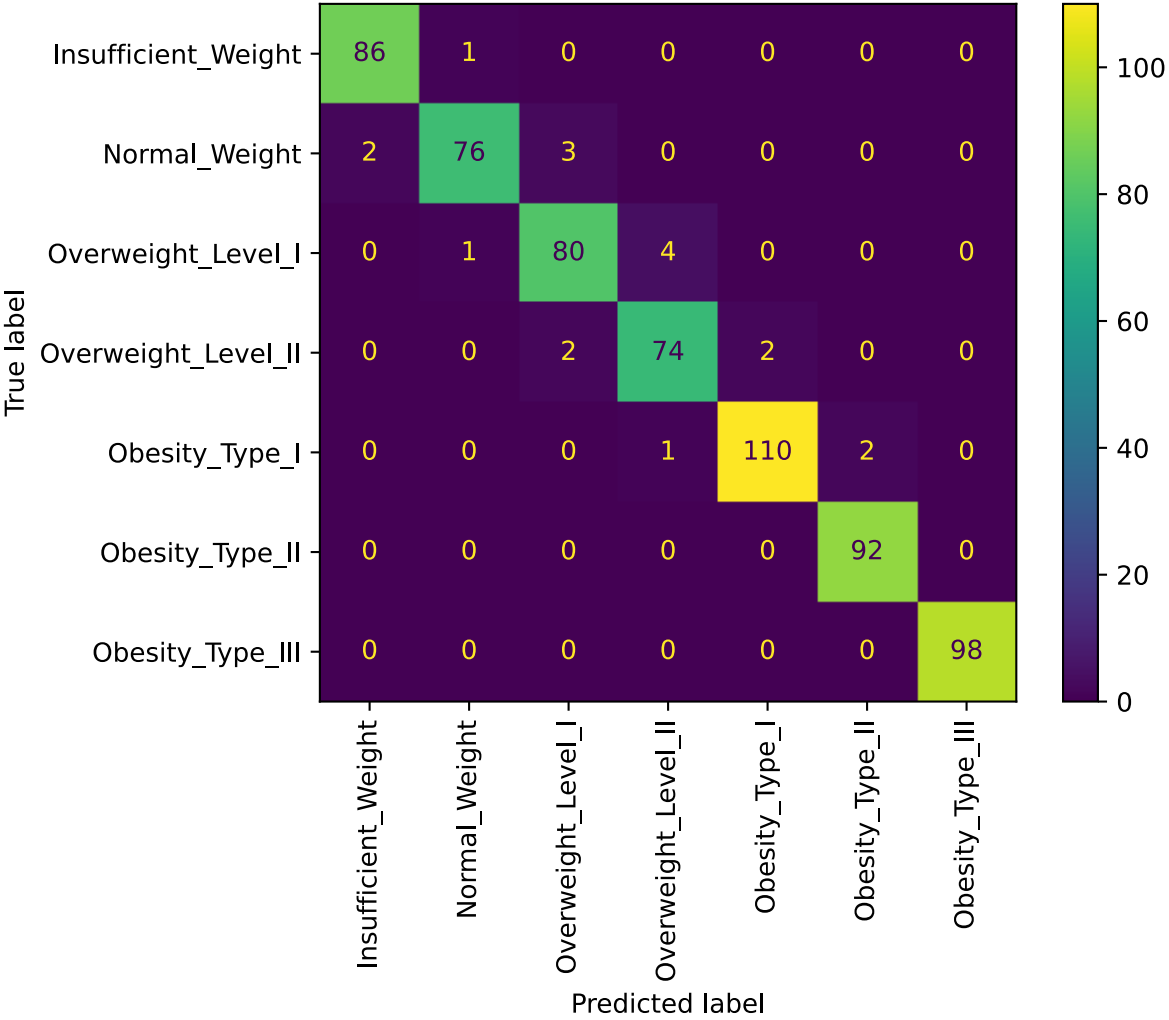
- *SVM + Logistic Regression + XGBoost*
- Usando mejores hiperparámetros de búsqueda anterior
- Estimador final: *Logistic Regression*

Modelo elegido.

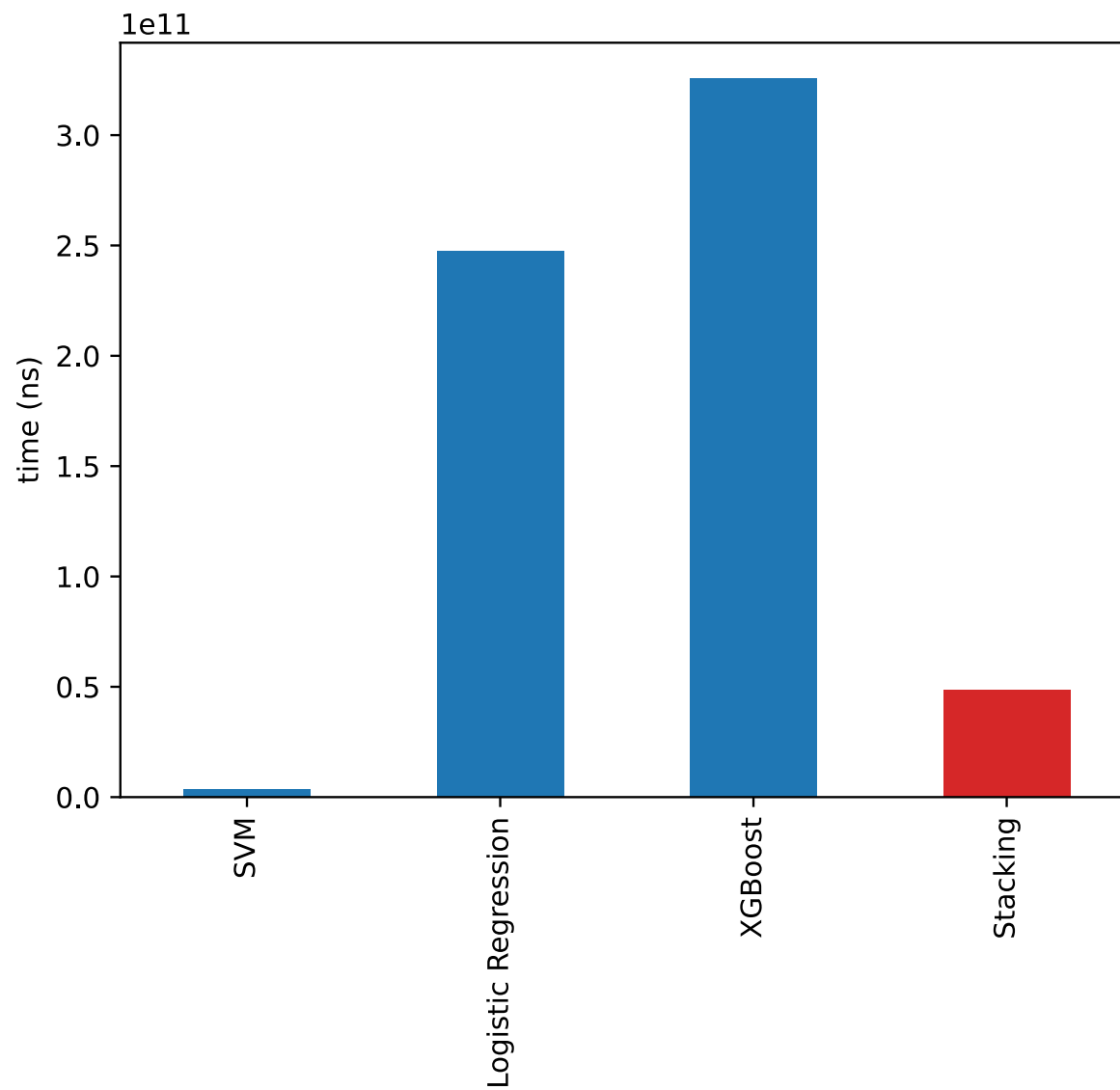
Experimentación



Experimentación



Experimentación



Gracias por su atención

