

Ensamble Secuencial de Modelos para predicción de precios de vivienda y enfermedad de Parkinson

Joaquín Borja León

Dpto. Ciencias de la Computación e Inteligencia Artificial
Universidad de Sevilla
Sevilla, España
joaborleo@alum.us.es

Héctor Guerra Prada

Dpto. Ciencias de la Computación e Inteligencia Artificial
Universidad de Sevilla
Sevilla, España
hecguepra@alum.us.es

Resumen—Este trabajo presenta la implementación desde cero de un meta-algoritmo de ensamble secuencial para tareas de regresión. Se utilizan dos estimadores débiles: `DecisionTreeRegressor` y `LinearRegression`. Se aplican ambos enfoques a dos conjuntos de datos: precios de viviendas (`house_prices.csv`) y datos clínicos de la enfermedad de Parkinson (`parkinsons.csv`). La búsqueda manual de hiperparámetros se realiza mediante validación cruzada, optimizando `n_estimators`, `sample_size`, `lr`, `max_depth` para árboles y `fit_intercept` para regresión lineal. También se incluye un mecanismo opcional de parada temprana (*early stopping*). Los resultados demuestran que `DecisionTreeRegressor` capta relaciones no lineales en Parkinson ($R^2=0.9323$) y `LinearRegression` funciona bien en precios de vivienda ($R^2=0.7875$).

Palabras clave—Inteligencia Artificial, Ensamble Secuencial, Regresión, Boosting, Precios de Vivienda, Enfermedad de Parkinson.

I. INTRODUCCIÓN

Los métodos basados en ensambles han demostrado una capacidad predictiva superior a la de modelos individuales en múltiples dominios. En particular, los ensambles secuenciales construyen modelos débiles de forma iterativa, de modo que cada nuevo estimador corrige los errores residuales del conjunto previo. Este esquema es la base de algoritmos como AdaBoost y Gradient Boosting [1], [2].

El objetivo de este trabajo es implementar y estudiar un meta-algoritmo de ensamble secuencial para tareas de regresión, aplicándolo a dos conjuntos de datos con características muy distintas:

- Conjunto de datos de precios de viviendas (`house_prices.csv`): 560 registros con variables físicas y categóricas, prediciendo el precio de venta.
- Conjunto de datos de la enfermedad de Parkinson (`parkinsons.csv`): 2000 registros con datos clínicos y de señal de voz, prediciendo la puntuación total del UPDRS.

La métrica principal de evaluación es el coeficiente de determinación (R^2), complementado con el error absoluto medio (MAE). Los objetivos específicos son:

- 1) Implementar las fases de entrenamiento y predicción del ensamble secuencial sin usar librerías de boosting de alto nivel.

- 2) Comparar dos estimadores débiles: `DecisionTreeRegressor` y `LinearRegression`.
- 3) Realizar búsqueda manual de hiperparámetros (`n_estimators`, `sample_size`, `lr`, `max_depth` para árboles y `fit_intercept` para lineal) mediante validación cruzada.
- 4) Implementar un mecanismo opcional de *early stopping* basado en un subconjunto de validación interno.
- 5) Documentar todas las decisiones de diseño, exponer resultados reproducibles y extraer conclusiones.

La estructura del documento es la siguiente:

- Sección II: Preliminares.
- Sección III: Metodología.
- Sección IV: Resultados.
- Sección V: Conclusiones.
- Referencias bibliográficas.

II. PRELIMINARES

A. Métodos empleados

En este trabajo se utiliza el enfoque de ensamble secuencial (boosting). Cada estimador débil se entrena sobre los residuos del conjunto acumulado anterior. La función de pérdida es el error cuadrático medio, de modo que el residuo en la iteración m es:

$$r_{im} = y_i - F_{m-1}(x_i). \quad (1)$$

Se consideran dos familias de estimadores débiles:

- 1) `DecisionTreeRegressor` (árbol de decisión) con profundidad máxima (`max_depth`) variable.
- 2) `LinearRegression` (regresión lineal).

B. Trabajo relacionado

Los ensambles de boosting (AdaBoost, Gradient Boosting, XGBoost) han sido ampliamente estudiados en literatura reciente [1]–[3]. A diferencia de nuestro enfoque, se suele utilizar implementaciones optimizadas de alto nivel. En este trabajo implementamos desde cero la lógica de boosting.

III. METODOLOGÍA

En esta sección se describe el algoritmo de ensamble secuencial implementado. La implementación sigue el siguiente esquema:

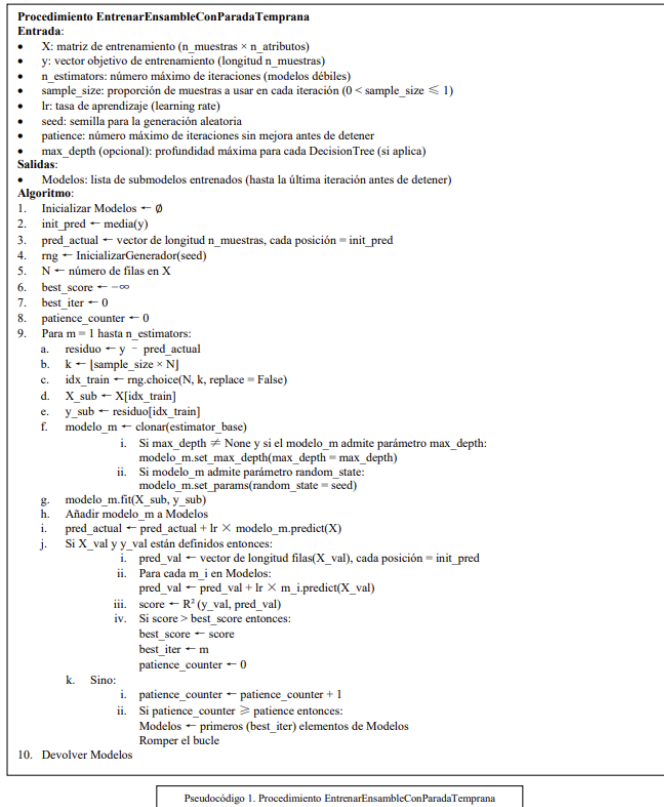


Fig. 1.

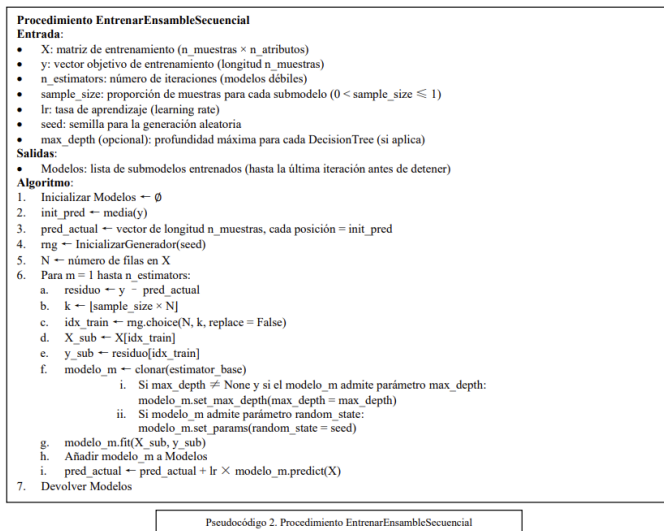
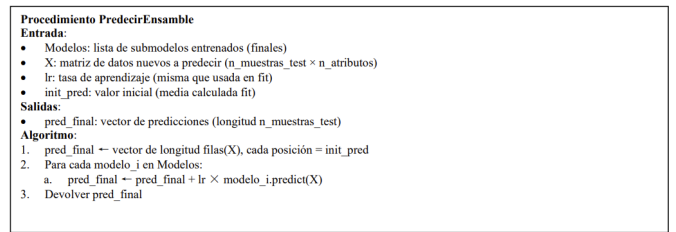


Fig. 2.

Modelos base: DecisionTreeRegressor, LinearRegression y Early Stopping:

a) **DecisionTreeRegressor:** El DecisionTreeRegressor es un modelo de aprendizaje supervisado que predice valores numéricos dividiendo el espacio de características en regiones, siguiendo una estructura jerárquica en forma de árbol. En cada nodo, se



Pseudocódigo 3. Procedimiento PredecirEnsamble

Fig. 3.

elige la variable y el umbral que mejor separen los datos para minimizar el error de predicción. Es un modelo *no lineal*, capaz de capturar relaciones complejas entre variables, aunque puede sobreajustar si no se controla su profundidad.

b) **LinearRegression:** El modelo LinearRegression ajusta una función lineal a los datos, minimizando el error cuadrático medio entre las predicciones y los valores reales. Es un modelo *lineal y simple*, adecuado cuando hay una relación aproximadamente lineal entre las variables predictoras y la variable objetivo. No capta relaciones no lineales ni interacciones, pero tiene bajo riesgo de sobreajuste y es computacionalmente eficiente.

c) **Early Stopping:** Para implementar el mecanismo de *early stopping*, en cada iteración se reserva un 20% del conjunto de entrenamiento (subconjunto de validación). Se calcula el error en validación después de cada nuevo modelo y se detiene el entrenamiento si no hay mejora tras *patience* iteraciones consecutivas.

IV. RESULTADOS

A. Ensemble con DecisionTreeRegressor

Configuración de partida: DecisionTreeRegressor($\text{max_depth}=3$), $n_{\text{estimators}}=50$, $\text{sample_size}=0.8$, $\text{lr}=0.1$, $\text{random_state}=42$. Validación cruzada 5-fold:

- **House Prices:** R^2 por fold = [0.72356832 0.75864956 0.76608286 0.78143207 0.77203406], media = 0.7603533739391801.
- **Parkinsons:** R^2 por fold = [0.63276246 0.66040138 0.63794864 0.65623765 0.6812172], media = 0.6537134652991631.

Primera búsqueda manual de hiperparámetros (Grid Search simple):

$$\begin{aligned} n_{\text{estimators}} &\in \{10, 50, 100\}, \\ lr &\in \{0.01, 0.1, 0.2\}, \\ \text{sample_size} &\in \{0.6, 0.8, 1.0\}, \\ \text{max_depth} &\in \{3, 5, 7\}. \end{aligned}$$

Evaluación final en test:

- **House Prices (óptimo):** $n_{\text{estimators}} = 150$, $lr = 0.3$, $\text{sample_size} = 1.0$, $\text{max_depth} = 3$.
Test $R^2 = 0.7191232968866706$, MAE = 29074.149712968017.

TABLA I
TOP-10 HOUSE PRICES – PRIMERA BÚSQUEDA

$n_estimators$	lr	$sample_size$	max_depth	r^2_mean
100	0.2	0.8	3	0.790717
100	0.1	0.6	3	0.790011
100	0.2	1.0	3	0.788360

TABLA II
TOP-10 PARKINSONS – PRIMERA BÚSQUEDA

$n_estimators$	lr	$sample_size$	max_depth	r^2_mean
100	0.1	1.0	7	0.919854
100	0.2	1.0	7	0.916761
50	0.2	1.0	7	0.915149

- **Parkinsons (óptimo):** $n_estimators = 100$, $lr = 0.1$, $sample_size = 1.0$, $max_depth = 7$.
Test $R^2 = 0.936545781009529$, MAE = 1.615012580924182.

B. Ensemble con LinearRegression

Configuración inicial: LinearRegression
($fit_intercept=True$), $n_estimators=50$,
 $sample_size=0.8$, $lr=0.1$, $random_state=42$.
Validación cruzada 5-fold:

- **House Prices:** R^2 por fold = [0.60978187 0.77420332 0.65521267 0.82114321 0.44958941], media = 0.6619860969721221.
- **Parkinsons:** R^2 por fold = [0.20806173 0.2282059 0.12500756 0.12577325 0.1387287], media = 0.16515542749688544.

Primera búsqueda manual de hiperparámetros:

$$\begin{aligned} n_estimators &\in \{10, 50, 100\}, \\ lr &\in \{0.01, 0.1, 0.2\}, \\ sample_size &\in \{0.6, 0.8, 1.0\}, \\ fit_intercept &\in \{True, False\}. \end{aligned}$$

TABLA III
TOP-5 HOUSE PRICES – PRIMERA BÚSQUEDA

$n_estimators$	lr	$sample_size$	$fit_intercept$	r^2_mean
10	0.2	1.0	True	0.720497
10	0.2	0.8	True	0.715898
50	0.1	1.0	True	0.703330

TABLA IV
TOP-5 PARKINSONS – PRIMERA BÚSQUEDA

$n_estimators$	lr	$sample_size$	$fit_intercept$	r^2_mean
10	0.2	1.0	True	0.166249
10	0.2	0.8	True	0.166206
50	0.2	0.8	True	0.166206

Evaluación final en test:

- **House Prices (lineal óptimo):** $n_estimators = 10$, $lr = 0.2$, $sample_size = 1.0$, $fit_intercept = True$.

Test $R^2 = 0.766222639788844$, MAE = 25257.28627191641.

- **Parkinsons (lineal óptimo):** $n_estimators = 10$, $lr = 0.2$, $sample_size = 1.0$, $fit_intercept = True$.
Test $R^2 = 0.14895427840774145$, MAE = 8.073492674590106.

C. Early Stopping (Parada Temprana)

- **House (LR + early stopping):** Test $R^2 = 0.766222639788844$, MAE = 25257.28627191643.
- **Parkinsons (LR + early stopping):** Test $R^2 = 0.14887395829827899$, MAE = 8.077265793598144.

V. CONCLUSIONES

- Este trabajo ha abordado la implementación desde cero de un meta-algoritmo de ensamble secuencial para tareas de regresión, aplicándolo a dos conjuntos de datos de naturaleza diferente: precios de viviendas y progresión de la enfermedad de Parkinson. Se han comparado dos tipos de modelos base (DecisionTreeRegressor y LinearRegression) y se ha evaluado su rendimiento mediante validación cruzada y métricas como el coeficiente de determinación (R^2) y el error absoluto medio (MAE). Asimismo, se ha incorporado un mecanismo de parada temprana para optimizar el proceso de entrenamiento.
- Los resultados muestran que el desempeño del meta-modelo depende en gran medida del tipo de estimador base utilizado. El modelo basado en árboles obtuvo un rendimiento excelente en el conjunto de Parkinson ($R^2 = 0.936545781009529$), gracias a su capacidad para capturar relaciones no lineales complejas. En cambio, su rendimiento fue inferior en el conjunto de precios de vivienda, donde alcanzó $R^2 = 0.678402654711951$, posiblemente debido al tamaño limitado del dataset, que restringe la capacidad de generalización del ensamble. Por otro lado, la regresión lineal se comportó mejor en este conjunto ($R^2 = 0.766222639788844$), mostrando mayor estabilidad con pocos datos y relaciones más simples. Sin embargo, en el conjunto de Parkinson, su desempeño fue muy bajo ($R^2 = 0.14895427840774145$), debido a que los valores no siguen una relación lineal clara con las variables predictoras. La búsqueda manual de hiperparámetros permitió identificar configuraciones óptimas en ambos casos.

REFERENCIAS

- [1] Y. Freund and R. Schapire, "A Decision-Theoretic Generalization of On-Line Learning and an application to boosting," *Journal of Computer and System Sciences*, 1997.
- [2] J. Friedman, T. Hastie, and R. Tibshirani, "The Elements of Statistical Learning," Springer, 2001.
- [3] T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016.
- [4] Scikit-learn Developers, "LinearRegression," https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html.
- [5] Scikit-learn Developers, "DecisionTreeRegressor," <https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeRegressor.html>.