

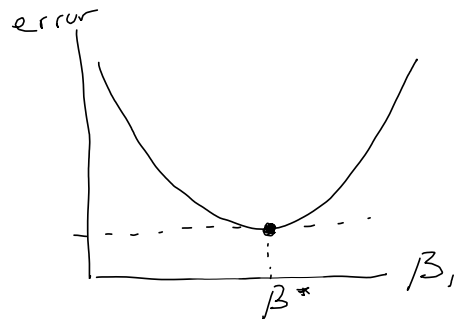
Métricas de evaluación regresión

Nombre	Fórmulas $e = y - \hat{y}$ <small>i: filas</small>
Mean squared error (<u>MSE</u>) [0, +∞) ✓ ✗	$\frac{1}{n} \sum_{i=1}^n e_i^2$
Root MSE [0, +∞) ✓ ✗	$\sqrt{\frac{1}{n} \sum_{i=1}^n e_i^2}$
Mean absolute error (<u>MAE</u>) [0, +∞) ✓ ✗	$\frac{1}{n} \sum_{i=1}^n e_i $
Mean absolute % error [0, +∞) ✓ ✗	$\frac{100\%}{n} \sum_{i=1}^n \left \frac{e_i}{y_i} \right $
R^2 [0, 1] ✗ ✓	$1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2}$

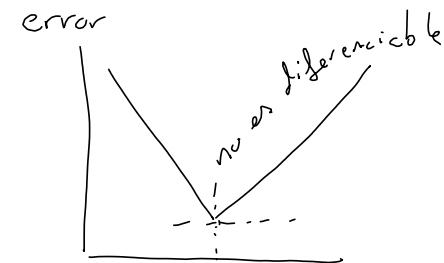
- ## Ventajas
- misma que Reg. Lineal
 - magnitud es más comparable con el target
 - visión exacta del error medio
 - error relativo
 - comparable con otros datasets

- ## Desventajas
- e^2 perdemos magnitud error
 - no es el error medio exacto
 - Reg lineal no se ajusta con el MAE.
 - no funciona si target cercano 0
 - no da magnitud del error

note: $\min_{\beta_1} \sum_{i=1}^n e_i^2$



$\min_{\beta_1} \sum_{i=1}^n |e_i|$



Workflow

Data Understanding → Data Preparation → Modelling → Evaluation

Data Understanding

- relación atributo y target → Modelling / Evaluation
- entender atributo y posible limpieza → Data Preparation

Data Preparation

- "traducir" todas las variables a numéricas
 - Fechas: extraer inputs mes, año, día, trimestre, ...
 - nulos: 1. Tengo info en otra variable, 2. Imputación inteligente, 3. Imputación "a saw" (media, moda, ...)
 - categóricas
 - ↳ Boolean → orden intrínseco (versiones) → orden de fuentes externas → OHE → agrupar (<100) → Frec. encoding
- preparar según el algoritmo (e.g. outliers)
 - Reglas: all in, sin cambios
 - Geométricas: transformación logarítmica de target y atributo numérico (relación lineal)
 - Vecindad: scaling

Modelling

1. Instanciar: escoger algoritmo y parámetros
2. Entrenar: `.fit`
↳ X_{train} , y_{train}

Evaluation

1. Overfitting: comparar (y_{train} , $pred-y_{train}$) vs (y_{test} , $pred-y_{test}$)

Métricas dependen de la tarea

Clasificación: Accuracy, F1 Score, Confusion matrix, AUC

Regresión: MSE, RMSE, MAE, MAPE, R^2

2. Desbalanceo: objetivo NO es tener proporción igualitaria

→ Oversampling

→ Undersampling

→ prediga mejor clase minoritaria



development, NO validation

¿Por qué mi modelo funciona mal?

1. Memorización: el modelo contiene reglas que solo aplican a los datos observados
2. Realidad tiene una distribución distinta

Original

dev val → realidad: decisión consciente

train / test aleatoria: quiero misma distribución train / test

↳ random hold out: una sola partición

↳ si tengo pocos datos: múltiples versiones

↳ K-Fold Cross Val.

↳ Bootstrap