

# Act 10: Programando Regresión Lineal Múltiple en Python

Patricio Ricardí

March 2025

## ¿Qué es la regresión lineal múltiple?

Es un método estadístico que permite predecir un valor numérico (como el número de veces que se comparte un artículo) usando varias variables relacionadas. La fórmula básica es:

$$\text{Predicción} = \text{Intercepto} + (m_1 \times \text{Variable}_1) + (m_2 \times \text{Variable}_2) + \dots$$

Donde:

- $m_1, m_2, \dots$  = Coeficientes (indican la importancia de cada variable)
- Intercepto = Valor base cuando todas las variables son cero

## Pasos realizados

1. **\*\*Cargamos datos\*\***: Usamos un archivo CSV con información de artículos sobre Machine Learning, incluyendo: - Cantidad de palabras - Número de enlaces, comentarios e imágenes - Veces que fue compartido
2. **\*\*Limpiamos datos\*\***: - Quitamos artículos muy largos (> 3000 palabras)  
- Eliminamos casos extremos (> 80,000 compartidos)
3. **\*\*Creamos variable "interacciones"\*\***:

$$\text{interacciones} = \text{enlaces} + \text{comentarios} + \text{imágenes}$$

4. **\*\*Preparamos modelo\*\***: - Variables predictoras: `Word count` (palabras) y `interacciones` - Variable a predecir: `# Shares` (compartidos)
5. **\*\*Entrenamos modelo\*\***: Usamos la librería `scikit-learn` de Python.

## Código simplificado

```
[language=Python] Filtrar datos datos_filtrados = dataset[(dataset["Wordcount"] < 3000)(dataset["Shares"] < 80000)]
```

```

    Crear variable "interacciones" interacciones = datos_filtrados["ofLinks"] +
datos_filtrados["ofcomments"].fillna(0) + datos_filtrados["Imagesvideo"]
    Entrenar modelo modelo = LinearRegression() modelo.fit(X=variables_predictoras, y =
objetivo)
    Evaluar print("Error cuadrático medio:", mean_squared_error(y_true, y_pred))

```

## Resultados clave

- **Coefficientes:**
  - Por cada palabra: +6.63 compartidos
  - Por cada interacción: -483.41 compartidos
- **Precisión:**
  - Error medio: 352,122,816.48
  - $R^2 = 0.11$  (explica 11% de la variabilidad)
- **Ejemplo:**
  - Artículo con 2000 palabras y 20 interacciones:

$$(2000 \times 6.63) + (20 \times -483.41) = 3,591.8 \text{compartidos}$$

## Conclusiones

- El modelo es básico pero útil como introducción.
- Para mejorarlo podríamos:
  - Usar más variables (ej: tema del artículo)
  - Probar modelos más complejos (ej: redes neuronales)
  - Mejorar la limpieza de datos
- Muestra cómo pequeños cambios en las variables pueden afectar las predicciones.