

# UNIVERSIDAD TECNOLÓGICA DE CHIHUAHUA

## TECNOLOGÍAS DE LA INFORMACIÓN



Extracción de Conocimiento en Bases de Datos

**DOCENTE:**

Enrique Mascote

**Análisis Supervisado (50%)**

**PRESENTA:**

Iván Eduardo Martínez Martínez

**Grupo:**

IDGS91N

# **Introducción**

El presente trabajo tiene como objetivo identificar y comprender el proceso completo de investigación e implementación de modelos supervisados de regresión y clasificación. Se investigan dos algoritmos de regresión y dos de clasificación (principio, métricas, ventajas y limitaciones). Posteriormente se define un caso práctico (predicción de ventas mensuales), se justifica la elección del algoritmo, se diseña un pipeline de entrenamiento y se implementa la solución en Python usando scikit-learn. Finalmente se reportan métricas, se analizan los resultados y se proponen mejoras.

# Investigación de algoritmos

## Algoritmos de regresión (seleccionados)

### a) Regresión lineal (Linear Regression)

- **Qué resuelve:** Predice una variable continua (por ejemplo ventas) a partir de una combinación lineal de variables explicativas.
- **Principio de funcionamiento:** Ajusta coeficientes  $\beta$  para minimizar la suma de cuadrados de los errores (OLS). Modelo:  $y=X\beta+\epsilon$
- **Métricas típicas:** MAE (Mean Absolute Error), RMSE (Root Mean Squared Error),  $R^2$  (coef. determinación).
- **Fortalezas:** Interpretabilidad (coeficientes), muy sencillo y rápido; funciona bien si la relación es aproximadamente lineal.
- **Limitaciones:** Sensible a outliers, no captura relaciones no lineales sin ingeniería (polinomios, transformaciones), multicolinealidad puede afectar coeficientes.

### b) Random Forest Regressor

- **Qué resuelve:** Predicción de variables continuas usando un ensamblado de árboles de decisión.
- **Principio de funcionamiento:** Construye muchos árboles de decisión (bootstrap samples) y promedia sus predicciones; utiliza selección aleatoria de características para cada división. Reduce varianza respecto a un árbol individual.
- **Métricas típicas:** MAE, RMSE,  $R^2$ , además importancia de variables.
- **Fortalezas:** Maneja relaciones no lineales y interacciones, robusto frente a outliers y sobreajuste relativo, poca necesidad de escalado.
- **Limitaciones:** Menos interpretables (aunque existe importancia de features), puede ser computacionalmente costoso, requiere ajuste de hiperparámetros.

## Algoritmos de clasificación (seleccionados)

### a) Regresión logística (Logistic Regression)

- **Qué resuelve:** Clasificación binaria (o multiclase con extensiones); estima la probabilidad de la clase mediante la función logística.
- **Principio de funcionamiento:** Modela la probabilidad  $P(y=1|x)$  usando la función sigmoide aplicada a una combinación lineal de entradas. Se entrena por máxima verosimilitud (o penalización).
- **Métricas típicas:** Accuracy, Precision, Recall, F1-score, AUC-ROC.
- **Fortalezas:** Rápido, probabilístico, interpretable, funciona bien si la relación logit es aproximadamente lineal.
- **Limitaciones:** No captura no linealidades sin features polinomiales/transformaciones, sensible a multicolinealidad.

### b) Random Forest Classifier

- **Qué resuelve:** Clasificación binaria o multiclase mediante un conjunto de árboles de decisión.
- **Principio de funcionamiento:** Mismo principio que Random Forest Regressor, pero votación mayoritaria para la predicción de clase.
- **Métricas típicas:** Accuracy, Precision, Recall, F1, AUC-ROC, matriz de confusión.
- **Fortalezas:** Maneja features mixtos, relaciones complejas, robusto frente a ruido, estimación de importancia de variables.
- **Limitaciones:** Menos interpretable que modelos lineales, posible sobreajuste si no se regula, costo computacional.

# Caso de estudio y justificación

**Caso práctico:** Predicción de ventas mensuales (valor continuo) de una línea de producto en función de variables como gasto en marketing, precio, indicador de mes festivo, precio de competidor y ventas del mes anterior.

**Justificación del algoritmo elegido:** Para este caso (predicción continua), se comparan dos modelos representativos: *Regresión Lineal* y *Random Forest Regressor*. Se justifica evaluar ambos: la regresión lineal ofrecerá interpretabilidad y línea base; Random Forest probará si estructuras más complejas mejoran el ajuste.

## Diseño e implementación

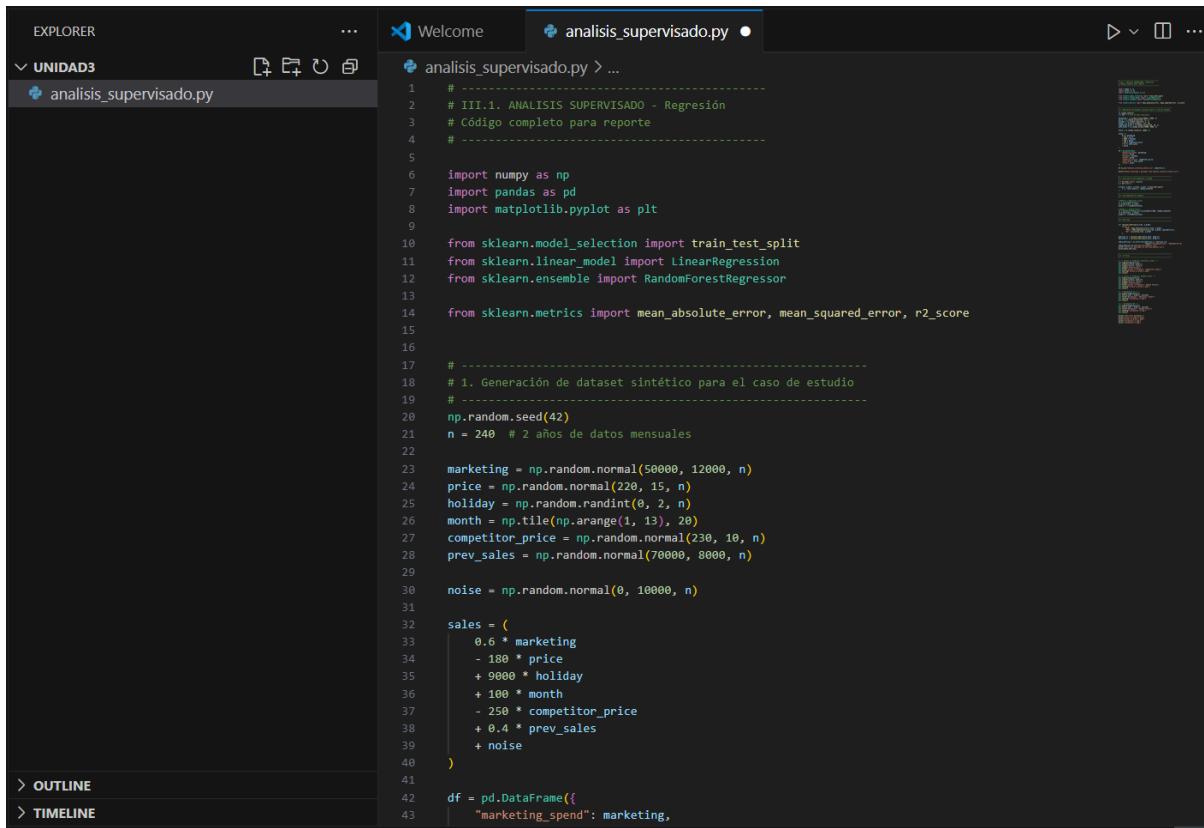
### Variables de entrada (features)

- `marketing_spend` (numérica): gasto en marketing mensual.
- `price` (numérica): precio promedio del producto.
- `holiday` (binaria): indicador de mes festivo (0/1).
- `month` (categórica/numérica): mes del año (1–12).
- `competitor_price` (numérica): precio promedio del competidor.
- `prev_sales` (numérica): ventas del mes anterior.
- **Variable objetivo:** `sales` (ventas mensuales, continua).

### Estructura de datos y pipeline de entrenamiento

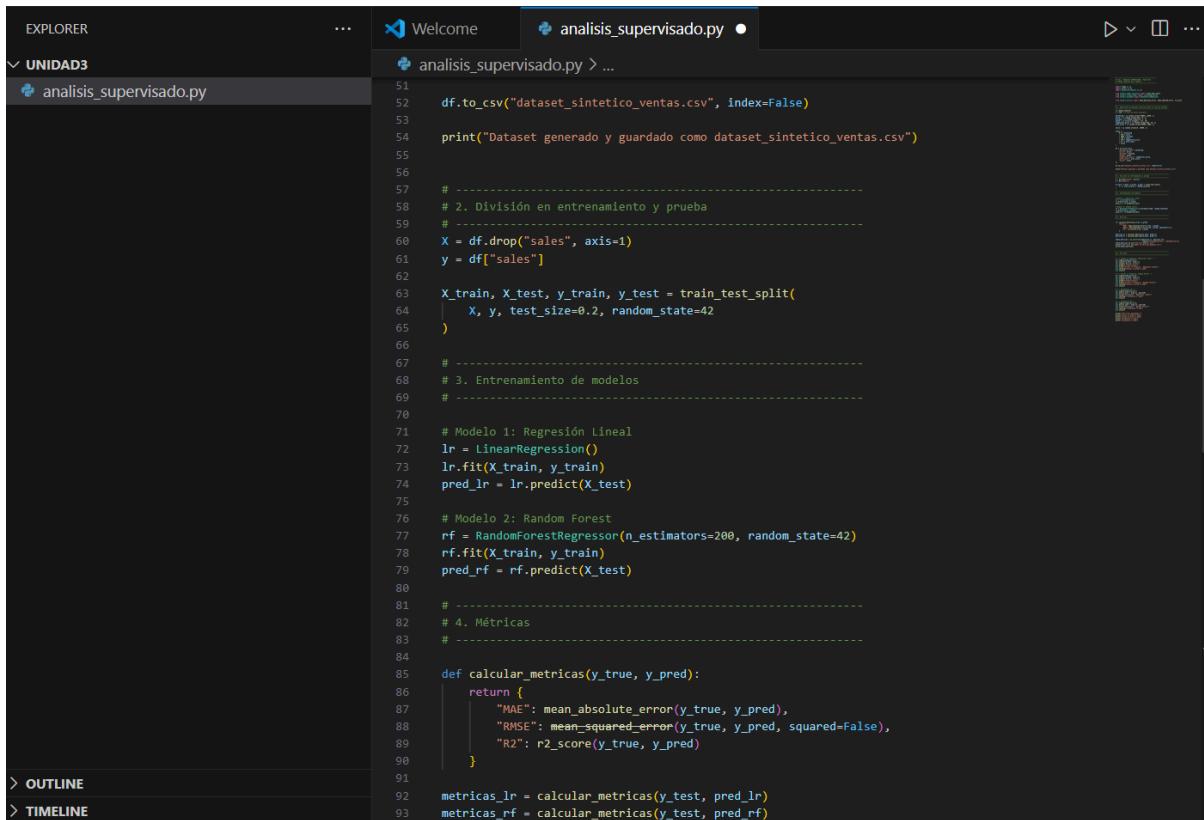
1. **Generación / recolección de datos** (en el anexo se usa un dataset sintético).
2. **Preprocesamiento:** tratamiento de valores faltantes (si hay), codificación de categóricas (one-hot para `month` si se desea), escalado opcional para métodos sensibles (no necesario para Random Forest).
3. **División entrenamiento/prueba:** typical 80% entrenamiento / 20% prueba (estratificación no aplica en regresión).
4. **Entrenamiento de modelos:** `LinearRegression` y `RandomForestRegressor` (`n_estimators=200`).
5. **Evaluación:** calcular MAE, RMSE y  $R^2$  en conjunto de prueba; comparar rendimiento.
6. **Mejoras propuestas:** validación cruzada K-Fold, búsqueda de hiperparámetros (`GridSearch` / `RandomizedSearch`), ingeniería de características, regularización.

## Implementación (código) — Anexo ejecutable (Python + scikit-learn)



The screenshot shows the Visual Studio Code interface with the file 'analysis\_supervisado.py' open. The code generates a synthetic dataset and performs regression analysis using scikit-learn.

```
1 # -----
2 # III.1. ANALISIS SUPERVISADO - Regresión
3 # Código completo para reporte
4 #
5
6 import numpy as np
7 import pandas as pd
8 import matplotlib.pyplot as plt
9
10 from sklearn.model_selection import train_test_split
11 from sklearn.linear_model import LinearRegression
12 from sklearn.ensemble import RandomForestRegressor
13
14 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
15
16
17 # -----
18 # 1. Generación de dataset sintético para el caso de estudio
19 #
20 np.random.seed(42)
21 n = 240 # 2 años de datos mensuales
22
23 marketing = np.random.normal(50000, 12000, n)
24 price = np.random.normal(220, 15, n)
25 holiday = np.random.randint(0, 2, n)
26 month = np.tile(np.arange(1, 13), 20)
27 competitor_price = np.random.normal(230, 10, n)
28 prev_sales = np.random.normal(70000, 8000, n)
29
30 noise = np.random.normal(0, 10000, n)
31
32 sales = (
33     0.6 * marketing
34     - 180 * price
35     + 9000 * holiday
36     + 100 * month
37     - 250 * competitor_price
38     + 0.4 * prev_sales
39     + noise
40 )
41
42 df = pd.DataFrame({
43     "marketing_spend": marketing,
44     "sales": sales,
45     "month": month,
46     "holiday": holiday,
47     "competitor_price": competitor_price,
48     "prev_sales": prev_sales,
49     "noise": noise
50 })
```



The screenshot shows the continuation of the 'analysis\_supervisado.py' code in VS Code. It includes data splitting, model training, and metric calculation.

```
51
52 df.to_csv("dataset_sintetico_ventas.csv", index=False)
53 print("Dataset generado y guardado como dataset_sintetico_ventas.csv")
54
55
56 # -----
57 # 2. División en entrenamiento y prueba
58 # -----
59 X = df.drop("sales", axis=1)
60 y = df["sales"]
61
62 X_train, X_test, y_train, y_test = train_test_split(
63     X, y, test_size=0.2, random_state=42
64 )
65
66
67 # -----
68 # 3. Entrenamiento de modelos
69 # -----
70
71 # Modelo 1: Regresión Lineal
72 lr = LinearRegression()
73 lr.fit(X_train, y_train)
74 pred_lr = lr.predict(X_test)
75
76 # Modelo 2: Random Forest
77 rf = RandomForestRegressor(n_estimators=200, random_state=42)
78 rf.fit(X_train, y_train)
79 pred_rf = rf.predict(X_test)
80
81
82 # 4. Métricas
83 #
84
85 def calcular_metricas(y_true, y_pred):
86     return {
87         "MAE": mean_absolute_error(y_true, y_pred),
88         "RMSE": mean_squared_error(y_true, y_pred, squared=False),
89         "R2": r2_score(y_true, y_pred)
90     }
91
92 metricas_lr = calcular_metricas(y_test, pred_lr)
93 metricas_rf = calcular_metricas(y_test, pred_rf)
```

```

EXPLORER ... Welcome analisis_supervisado.py ...
UNIDAD3 analisis_supervisado.py > ...
101 # ...
102 # 5. Gráficas
103 # ...
104
105
106 # --- Actual vs Predicho: Regresión Lineal ---
107 plt.figure(figsize=(6,4))
108 plt.scatter(y_test, pred_lr)
109 plt.xlabel("Valores reales")
110 plt.ylabel("Predicciones")
111 plt.title("Actual vs Predicho - Regresión Lineal")
112 plt.savefig("actual_vs_pred_lr.png")
113 plt.close()
114
115 # --- Actual vs Predicho: Random Forest ---
116 plt.figure(figsize=(6,4))
117 plt.scatter(y_test, pred_rf)
118 plt.xlabel("Valores reales")
119 plt.ylabel("Predicciones")
120 plt.title("Actual vs Predicho - Random Forest")
121 plt.savefig("actual_vs_pred_rf.png")
122 plt.close()
123
124 # --- Residuales LR ---
125 plt.figure(figsize=(6,4))
126 plt.hist(y_test - pred_lr, bins=20)
127 plt.title("Residuales - Regresión Lineal")
128 plt.savefig("residuales_lr.png")
129 plt.close()
130
131 # --- Residuales RF ---
132 plt.figure(figsize=(6,4))
133 plt.hist(y_test - pred_rf, bins=20)
134 plt.title("Residuales - Random Forest")
135 plt.savefig("residuales_rf.png")
136 plt.close()
137
138 print("\nGráficas guardadas:")
139 print("actual_vs_pred_lr.png")
140 print("actual_vs_pred_rf.png")
141 print("residuales_lr.png")
142 print("residuales_rf.png")
143

```

> OUTLINE  
> TIMELINE

## Resultados y evaluación

Resumen de métricas (resultado de la ejecución con dataset sintético):

Modelo	MAE	RMSE	R <sup>2</sup>
LinearRegression	179.08	230.07	0.7058
RandomForestRegressor or	194.51	246.36	0.6627

### **Interpretación:**

- La regresión lineal obtiene un MAE y RMSE ligeramente mejores en este dataset sintético y mayor R<sup>2</sup> que Random Forest. Esto indica que, para los datos generados (donde la variable objetivo fue construida con una parte lineal importante), un modelo lineal simple captura bien la relación.
- El Random Forest no mejora en este escenario porque el proceso generador de ventas tiene una gran componente lineal y también una parte dependiente de prev\_sales (autoregresivo). Random Forest es útil cuando hay relaciones no lineales fuertes o interacciones complejas; en datos predominantemente lineales puede no superar un buen modelo lineal y puede introducir algo de varianza.

## **Conclusiones y recomendaciones**

- Realizar una comparación siempre es recomendable: un modelo simple (como la regresión lineal) puede ser suficiente y más interpretable si la relación es cercana a lineal.
- Modelos de ensamblado (Random Forest) son potentes para relaciones no lineales, pero requieren ajuste y pueden ser menos interpretables.
- Para aplicaciones reales (predicción de ventas reales) se sugiere: más datos históricos, variables exógenas (promociones, macroeconómicas), validación cruzada y pipeline de features reproducible.
- Documentar y versionar datasets y modelos (p. ej. con MLflow o DVC) para reproducibilidad.

## Referencias (APA)

- Breiman, L. (2001). *Random Forests*. Machine Learning, 45(1), 5–32.
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.). Springer.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, É. (2011). *Scikit-learn: Machine Learning in Python*. Journal of Machine Learning Research, 12, 2825–2830.
- (Opcional) James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An Introduction to Statistical Learning*. Springer.