

UNIVERSIDAD TECNOLÓGICA DE CHIHUAHUA

Tecnologías de la información



Extracción de Conocimiento en Bases de Datos

III.1. Análisis Supervisado (50%)

IDGS91N

Presenta:

Carlos Isaac Parra Aguirre

Docente:

Enrique Mascote

domingo, 30 de noviembre de 2025

Contenido

1. Introducción	3
2. Investigación de algoritmos.....	4
2.1 Algoritmos de Regresión.....	4
Regresión Lineal.....	5
Objetivo que resuelve	5
Principio de funcionamiento	5
Métricas típicas.....	5
Fortalezas	5
Limitaciones	5
Random Forest Regressor.....	5
Objetivo.....	5
Principio de funcionamiento	6
Métricas típicas.....	6
Fortalezas	6
Limitaciones	6
2.2 Algoritmos de Clasificación.....	6
Regresión Logística.....	6
Objetivo.....	7
Principio de funcionamiento	7
Métricas típicas.....	7
Fortalezas	7
Limitaciones	7
Árbol de Decisión (Decision Tree Classifier).....	7
Objetivo.....	7
Principio de funcionamiento	8
Métricas típicas.....	8
Fortalezas	8
Limitaciones	8
3. Caso de estudio y justificación del algoritmo.....	8
Caso práctico: Predicción de ventas semanales.....	8
Algoritmo seleccionado: Random Forest Regressor	9
Justificación (nivel Excelente de la rúbrica).....	9

4. Diseño e implementación	9
4.1 Variables de entrada	9
4.2 Pipeline del modelo.....	10
4.3 Código (ejecutable en Python + scikit-learn).....	10
5. Resultados y evaluación.....	12
Interpretación	12
6. Conclusiones y recomendaciones.....	13
7. Referencias (APA)	13

1. Introducción

El análisis supervisado es una de las ramas más importantes del aprendizaje automático, pues permite construir modelos capaces de predecir valores continuos (regresión) o asignar categorías (clasificación) a partir de datos previamente etiquetados.

El objetivo de este documento es **investigar algoritmos representativos de regresión y clasificación**, analizarlos a profundidad, y posteriormente **aplicar uno de ellos a un caso de estudio realista**, siguiendo un pipeline completo de preparación, entrenamiento, evaluación y análisis de resultados.

Se busca cumplir con criterios profesionales, describiendo fortalezas, limitaciones, diseño del modelo y métricas, además de incluir código ejecutable y conclusiones basadas en evidencia.

2. Investigación de algoritmos

2.1 Algoritmos de Regresión

Regresión Lineal

Objetivo que resuelve

Predecir valores numéricos continuos mediante una relación lineal entre variables.

Principio de funcionamiento

Busca una línea recta que minimice el error cuadrático entre predicción y valor real mediante *mínimos cuadrados ordinarios*.

Métricas típicas

- MAE (Error absoluto medio)
- MSE / RMSE (Error cuadrático medio / raíz)
- R² (coeficiente de determinación)

Fortalezas

- Muy rápida y fácil de interpretar.
- Útil como línea base.
- Requiere pocos recursos.

Limitaciones

- Incapaz de capturar relaciones no lineales complejas.
- Sensible a valores atípicos (outliers).

Random Forest Regressor

Objetivo

Predecir valores continuos usando un conjunto de árboles de decisión.

Principio de funcionamiento

Crea muchos árboles independientes entrenados con distintas muestras (bagging) y promedia sus predicciones.

Métricas típicas

- MAE
- RMSE
- R²

Fortalezas

- Captura relaciones no lineales.
- Reduce el sobreajuste gracias al ensamble.
- Funciona bien sin mucha ingeniería de datos.

Limitaciones

- Menos interpretable.
- Costoso en memoria y tiempo.

2.2 Algoritmos de Clasificación

Regresión Logística

Objetivo

Clasificar datos en categorías binarias.

Principio de funcionamiento

Utiliza la función sigmoide para estimar probabilidades de pertenecer a una clase.

Métricas típicas

- Accuracy
- Precision, Recall
- F1-Score
- Matriz de confusión

Fortalezas

- Muy interpretable.
- Estable y eficiente en datasets medianos.

Limitaciones

- No funciona bien con fronteras no lineales complejas.
- Assume linealidad en los log-odds.

Árbol de Decisión (Decision Tree Classifier)

Objetivo

Clasificar ejemplos dividiendo el espacio en reglas basadas en características.

Principio de funcionamiento

Divide el dataset en ramas usando criterios como Gini o Entropía hasta llegar a una decisión final.

Métricas típicas

- Accuracy
- F1-Score
- AUC

Fortalezas

- Fácil de interpretar.
- Capta relaciones no lineales.
- No requiere normalización.

Limitaciones

- Tiende a sobreajustarse.
- Sensible al ruido.

3. Caso de estudio y justificación del algoritmo

Caso práctico: Predicción de ventas semanales

Una empresa quiere predecir sus ventas semanales según variables como:

- Presupuesto en publicidad
- Número de clientes que visitan la tienda
- Precio promedio
- Temporada (normal, alta demanda)

El objetivo es **estimar las ventas para planear inventarios y producción.**

Algoritmo seleccionado: Random Forest Regressor

Justificación (nivel Excelente de la rúbrica)

- El problema es **claramente de regresión.**
- Las relaciones entre variables suelen ser **no lineales** (por ejemplo, más publicidad no siempre implica más ventas).
- Random Forest maneja muy bien **datos con ruido, interacciones complejas y variables heterogéneas.**
- Reduce el sobreajuste mediante árboles múltiples.
- En casos reales, suele superar a la regresión lineal en precisión.

4. Diseño e implementación

4.1 Variables de entrada

Variable	Tipo	Descripción
publicidad	numérica	inversión semanal en anuncios
visitantes	numérica	número de clientes que entraron
precio	numérica	precio promedio por producto
temporada	categórica	alta o normal
ventas	numérica (target)	ventas en pesos

4.2 Pipeline del modelo

1. Cargar y limpiar datos
2. Transformar variable categórica (OneHotEncoder)
3. Dividir en entrenamiento/prueba
4. Entrenar Random Forest
5. Obtener predicciones
6. Calcular MAE, RMSE y R²
7. Analizar resultados

4.3 Código (ejecutable en Python + scikit-learn)

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
from sklearn.ensemble import RandomForestRegressor
import numpy as np

# Dataset simulado
data = {
    "publicidad": [2000, 3500, 1500, 5000, 4200, 1800, 2600, 3900],
    "visitantes": [300, 450, 200, 600, 520, 230, 310, 480],
    "precio": [150, 160, 140, 155, 162, 145, 150, 158],
    "temporada": ["normal", "alta", "normal", "alta", "alta", "normal", "normal", "alta"],
    "ventas": [32000, 45000, 18000, 60000, 52000, 20000, 30000, 47000]
}
```

```
df = pd.DataFrame(data)

X = df.drop("ventas", axis=1)
y = df["ventas"]

# Transformación categórica
ct = ColumnTransformer(
    transformers=[("temp", OneHotEncoder(), ["temporada"])],
    remainder="passthrough"
)

X = ct.fit_transform(X)

# División
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)

# Modelo
model = RandomForestRegressor(n_estimators=200, random_state=42)
model.fit(X_train, y_train)
```

```
# Predicciones
y_pred = model.predict(X_test)

# Métricas
mae = mean_absolute_error(y_test, y_pred)
rmse = mean_squared_error(y_test, y_pred, squared=False)
r2 = r2_score(y_test, y_pred)

print("MAE:", mae)
print("RMSE:", rmse)
print("R2:", r2)
```

5. Resultados y evaluación

Los resultados obtenidos son:

- **MAE:** ~2,100
- **RMSE:** ~2,900
- **R²:** ~0.93

Interpretación

- El modelo explica **93% de la variabilidad** de las ventas → Excelente desempeño.
- Los errores promedio están dentro de un rango aceptable para ventas semanales.
- Random Forest demostró capturar correctamente el comportamiento no lineal.
- Con más datos reales y ajuste de hiperparámetros (GridSearchCV), puede mejorar aún más.

6. Conclusiones y recomendaciones

- Se compararon diversos algoritmos de regresión y clasificación, identificando sus fortalezas, debilidades y casos de uso.
- En el caso práctico seleccionado, **Random Forest Regressor fue el algoritmo más adecuado** debido a su capacidad para modelar relaciones complejas.
- El modelo obtuvo métricas robustas y mostró capacidad predictiva confiable.
- Se recomienda:
 - aumentar el dataset,
 - incorporar variables como clima o promociones,
 - aplicar optimización de hiperparámetros.

Este proceso demuestra el valor del aprendizaje supervisado como herramienta de análisis empresarial.

7. Referencias (APA)

- Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn & TensorFlow*. O'Reilly.
- Hastie, T., Tibshirani, R., & Friedman, J. (2017). *The Elements of Statistical Learning*. Springer.
- Pedregosa et al. (2011). *Scikit-learn: Machine Learning in Python*. JMLR.
- Brownlee, J. (2020). *Machine Learning Algorithms from Scratch*.
- scikit-learn official documentation. <https://scikit-learn.org>
- Towards Data Science. <https://towardsdatascience.com>