

# UNIVERSIDAD TECNOLÓGICA DE CHIHUAHUA

## TECNOLOGÍAS DE LA INFORMACIÓN



Extracción de Conocimiento en Bases de Datos

**DOCENTE:**

Enrique Mascote

**Reporte de Métricas de Evaluación (50%)**

**PRESENTA:**

Iván Eduardo Martínez Martínez

**Grupo:**

IDGS91N

# **Introducción**

Este reporte explica y aplica métricas de evaluación de modelos de clasificación y regresión. Se usa un caso práctico de clasificación con K-Nearest Neighbors (KNN) sobre una matriz de datos que contiene glucosa, edad y una etiqueta binaria etiqueta. Se divide el conjunto en 70% entrenamiento y 30% prueba, se escala, se prueban varios k y se elige el mejor según F1-score. Se muestran matriz de confusión, curva ROC y AUC.

# Investigación de métricas

## Métricas de clasificación

### 1. Accuracy

- Definición: ( $\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$ ).
- Interpretación: Proporción de predicciones correctas.
- Ventajas: Fácil de entender. Limitaciones: engañosa en datasets desbalanceados.

### 2. Precision (Precisión)

- Definición: ( $\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$ ).
- Interpretación: De las predicciones positivas, cuántas son realmente positivas.
- Ventajas: Útil cuando costo de falsos positivos es alto. Limitación: no mide falsos negativos.

### 3. Recall (Sensibilidad)

- Definición: ( $\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$ ).
- Interpretación: Proporción de positivos reales correctamente detectados.
- Ventajas: útil cuando se quiere detectar la mayoría de positivos. Limitación: puede aumentar FP.

### 4. F1-score

- Definición: ( $F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$ ).
- Interpretación: Balance entre precision y recall.
- Ventajas: buena métrica cuando hay desbalance. Limitación: no muestra trade-off en distintos umbrales.

### 5. ROC-AUC

- Definición: Área bajo la curva ROC que grafica TPR vs FPR para distintos umbrales.
- Interpretación: Probabilidad de que el clasificador ordene aleatoriamente una instancia positiva por encima de una negativa.
- Ventajas: independiente del umbral; útil para comparar clasificadores. Limitaciones: puede ser optimista si datos muy desbalanceados.

## Métricas de regresión

### 1. MAE (Mean Absolute Error)

- Fórmula: ( $\text{MAE} = \frac{1}{n} \sum |y_i - \hat{y}_i|$ ).
- Interpretación: Error absoluto promedio. Ventaja: interpretable en unidades. Limitación: no penaliza fuertemente grandes errores.

### 2. RMSE (Root Mean Squared Error)

- Fórmula: ( $\text{RMSE} = \sqrt{\frac{1}{n} \sum (y_i - \hat{y}_i)^2}$ ).
- Interpretación: desviación típica del error; penaliza más errores grandes. Limitación: sensible a outliers.

## Solución con KNN (preprocesamiento, entrenamiento, evaluación)

### Preparación de datos

- Se cargó `Matriz.csv` y se seleccionaron `glucosa` y `edad` como predictoras y `etiqueta` como etiqueta binaria.
- División: 70% entrenamiento / 30% prueba (estratificada por etiqueta).
- Escalado: StandardScaler (media 0, varianza 1) aplicado a variables numéricas.

### Implementación (KNN)

- Se entrenaron clasificadores KNN con  $k = 3, 5$  y  $7$ .
- Para cada  $k$  se calculó: accuracy, precision, recall, F1-score y ROC-AUC.
- Se generó la matriz de confusión y curva ROC (con AUC) para cada  $k$ .

## Análisis de resultados

- Compara la F1 entre  $k=3,5,7$ . Indica cuál obtuvo mayor F1 y por qué (por ejemplo,  $k$  pequeño -> mayor varianza;  $k$  grande -> más suavizado).
- Observa la matriz de confusión para ver si hay muchos falsos negativos o falsos positivos (importante según el contexto clínico).
- ROC-AUC te dice si el modelo separa bien clases independientemente del umbral.

### Mejoras sugeridas:

- Probar más  $k$  (ej. 1–20) y usar validación cruzada para elegir  $k$ .
- Probar otros algoritmos (Logistic Regression, SVM, Random Forest).
- Si hay desbalance, aplicar sobremuestreo (SMOTE) o cambiar el umbral de decisión.
- Añadir más variables predictoras (si existen) y hacer selección de features.

# Código

```
EXPLORER ... Welcome reporte_metricas_knn.py ... reporte_metricas_knn.py > ...
reporte_metricas_knn.py > ...
1 # reporte_metricas_knn.py
2 import os
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
8 import matplotlib.pyplot as plt
9
10 # Ruta del CSV (ajusta si es necesario)
11 path = "Matriz.csv"
12 df = pd.read_csv(path)
13
14 # Seleccionar columnas
15 X = df[['glucosa','edad']]
16 y = df['etiqueta']
17
18 # División 70/30 estratificada
19 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42, stratify=y)
20
21 # Escalado
22 scaler = StandardScaler()
23 X_train_s = scaler.fit_transform(X_train)
24 X_test_s = scaler.transform(X_test)
25
26 # Probar varios k
27 ks = [3,5,7]
28 results = []
29 outdir = "knn_results"
30 os.makedirs(outdir, exist_ok=True)
31
32 for k in ks:
33     knn = KNeighborsClassifier(n_neighbors=k)
34     knn.fit(X_train_s, y_train)
35     y_pred = knn.predict(X_test_s)
36     y_proba = knn.predict_proba(X_test_s)[:,1] if hasattr(knn, "predict_proba") else y_pred
37
38     acc = accuracy_score(y_test, y_pred)
39     prec = precision_score(y_test, y_pred, zero_division=0)
40     rec = recall_score(y_test, y_pred, zero_division=0)
41     f1 = f1_score(y_test, y_pred, zero_division=0)
42     try:
43         auc = roc_auc_score(y_test, y_proba)
44     except Exception:
45         auc = float('nan')
46
47     results.append({"k":k, "accuracy":acc, "precision":prec, "recall":rec, "f1":f1, "roc_auc":auc})
48
49 # guardar confusion matrix
50 cm = confusion_matrix(y_test, y_pred)
51 pd.DataFrame(cm, index=[True 0, True 1], columns=[Pred 0, Pred 1]).to_csv(f"{outdir}/confusion_matrix.csv")
52
53 # guardar ROC
54 fpr, tpr, _ = roc_curve(y_test, y_proba)
55 plt.figure()
56 plt.plot(fpr, tpr)
57 plt.plot([0,1],[0,1], linestyle='--')
58 plt.xlabel("False Positive Rate")
59 plt.ylabel("True Positive Rate")
60 plt.title(f"ROC k={k} (AUC: {auc:.3f})")
61 plt.savefig(f"{outdir}/roc_k{k}.png")
62 plt.close()
63
64 pd.DataFrame(results).set_index('k').to_csv(f"{outdir}/knn_metrics.csv")
65 print("Resultados guardados en:", outdir)
66 print(pd.DataFrame(results).set_index('k'))
```

## Conclusión

En conclusión, el análisis realizado permitió comprender la importancia de seleccionar adecuadamente las métricas de evaluación para interpretar correctamente el rendimiento de modelos de clasificación. Las métricas estudiadas —accuracy, precision, recall, F1-score y ROC-AUC— mostraron que cada una aporta una perspectiva distinta del comportamiento del modelo, especialmente en contextos con posibles desbalances en los datos. En la aplicación práctica con KNN, el preprocessamiento mediante escalado y la comparación de diferentes valores de  $k$  demostraron que pequeñas variaciones en los hiperparámetros pueden influir significativamente en los resultados.

## Referencias (APA)

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.

Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning*. Springer.

Fawcett, T. (2006). *An introduction to ROC analysis*. Pattern Recognition Letters, 27(8), 861–874.