

Análisis del Dataset Iris: Comparativa Multilenguaje

Big Data e Inteligencia Artificial

Yeray Hurtado Dragón

Octubre 2025

Índice

1. Introducción	2
2. Algoritmo K-Nearest Neighbors (KNN)	3
2.1. Distancia Euclidiana	3
2.2. Ventajas y limitaciones	3
3. Descripción del dataset y metodología	4
4. Implementación en Python (Google Colab)	5
4.1. Matriz de correlación	5
4.2. Gráfico de dispersión	6
4.3. KNN con librerías	7
4.4. KNN implementado manualmente	8
5. Implementación en R (Posit Cloud)	9
5.1. Matriz de correlación	9
5.2. Gráfico de dispersión	9
5.3. KNN con librerías	10
5.4. KNN sin librerías (manual)	11
6. Implementación en Java (Weka)	12
6.1. Matriz de correlación	12
6.2. Gráfico de dispersión	12
6.3. KNN con librerías	12
6.4. KNN implementado manualmente	13
7. Comparativa multilenguaje	15
8. Conclusiones	16

1. Introducción

El presente documento analiza el **dataset Iris** utilizando el algoritmo **K-Nearest Neighbors (KNN)** en tres lenguajes de programación: **Python**, **R** y **Java (Weka)**.

Cada lenguaje ejecuta el mismo flujo de trabajo:

- Matriz de correlación.
- Gráfico de dispersión.
- Modelo KNN con librerías.
- Implementación manual del KNN.
- Matriz de confusión y precisión final.

El objetivo es comparar el rendimiento, facilidad de implementación y resultados obtenidos en cada entorno.

2. Algoritmo K-Nearest Neighbors (KNN)

El algoritmo **K-Nearest Neighbors (KNN)** es un método de clasificación supervisada basado en la idea de que los objetos similares tienden a estar cerca en el espacio de características. Para clasificar una nueva instancia, el algoritmo:

1. Calcula la distancia entre la nueva instancia y todas las observaciones del conjunto de entrenamiento.
2. Selecciona los k vecinos más cercanos según la métrica de distancia.
3. Asigna la clase más frecuente entre estos vecinos.

2.1. Distancia Euclidiana

La métrica más común es la **distancia euclidiana**, que mide la separación geométrica entre dos puntos p y q en un espacio de n dimensiones:

$$d(p, q) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

donde:

- p_i y q_i son los valores de la i -ésima característica de los puntos p y q .
- n es el número de características.

Esta distancia determina qué vecinos son los más cercanos y, por tanto, cuáles influyen en la clasificación de la nueva instancia.

2.2. Ventajas y limitaciones

Ventajas:

- Simple de entender e implementar.
- Funciona bien con problemas de clasificación multi-clase.

Limitaciones:

- Sensible a la escala de las características (recomendable normalizar datos).
- Puede ser lento en conjuntos de datos grandes.
- La elección de k influye directamente en la precisión.

3. Descripción del dataset y metodología

El dataset **Iris** contiene 150 observaciones de flores pertenecientes a tres especies: *Iris-setosa*, *Iris-versicolor* y *Iris-virginica*. Cada registro posee cuatro características numéricas:

- Largo y ancho del sépalo (SepalLengthCm, SepalWidthCm)
- Largo y ancho del pétalo (PetalLengthCm, PetalWidthCm)

Metodología general:

1. Carga y limpieza de datos.
2. Exploración (correlaciones y gráficos).
3. Preparación de datos (eliminación de columnas irrelevantes, codificación).
4. Entrenamiento con KNN (con librerías).
5. Implementación manual del KNN.
6. Evaluación con matriz de confusión y precisión.

4. Implementación en Python (Google Colab)

4.1. Matriz de correlación

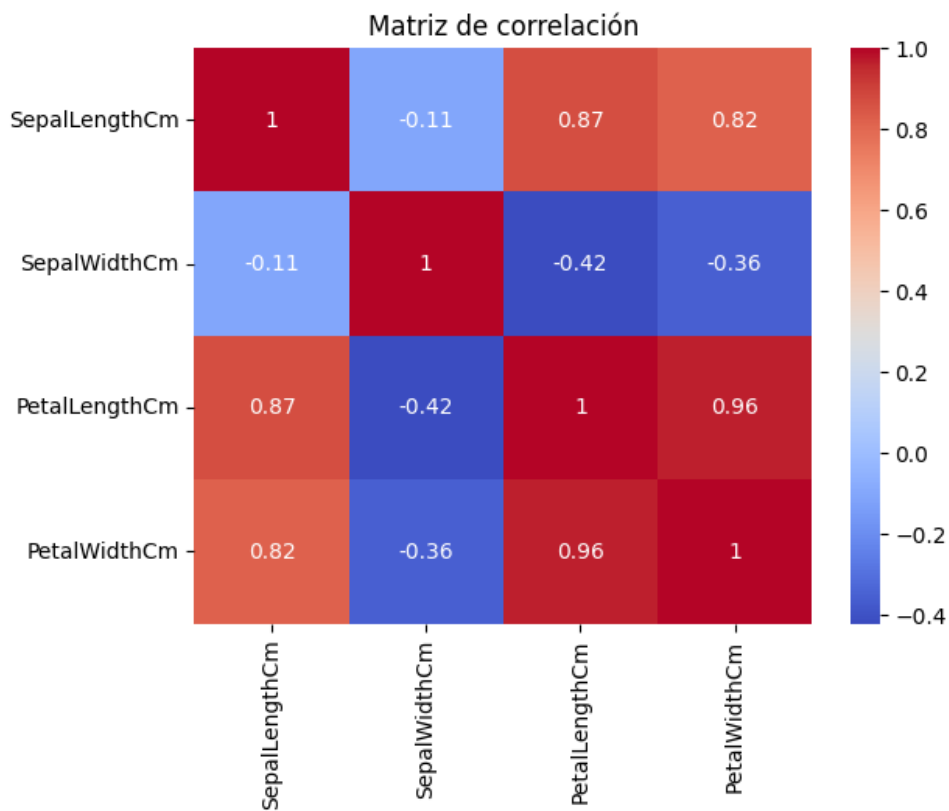


Figura 1: Matriz de correlación del dataset Iris en Python

Interpretación: Alta correlación entre `PetalLengthCm` y `PetalWidthCm`, indicando su relevancia en la clasificación.

4.2. Gráfico de dispersión

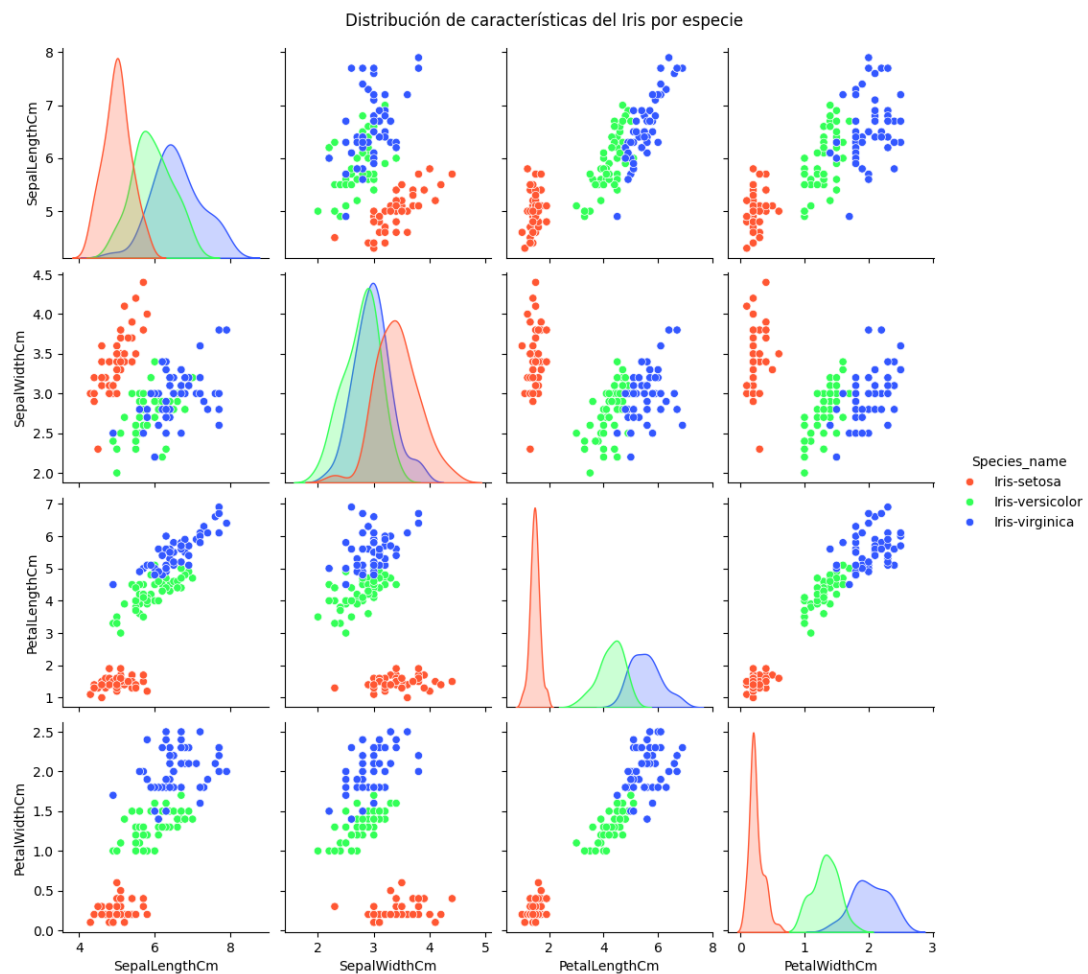


Figura 2: Gráfico de dispersión de las características del Iris (Python)

Interpretación: *Setosa* se separa claramente, mientras que *Versicolor* y *Virginica* presentan solapamiento.

4.3. KNN con librerías

```

1 # Características y etiquetas
2 X = df[['SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', '
   PetalWidthCm']]
3 y = df['Species'] # 0 nombres de especies
4
5 # División entrenamiento/prueba
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
   =0.3, random_state=123)
7
8 # Entrenamiento KNN
9 knn = KNeighborsClassifier(n_neighbors=27)
10 knn.fit(X_train, y_train)
11
12 # Precisión
13 accuracy = knn.score(X_test, y_test) * 100
14 print(f"Precisión del modelo KNN: {accuracy:.2f}%")

```

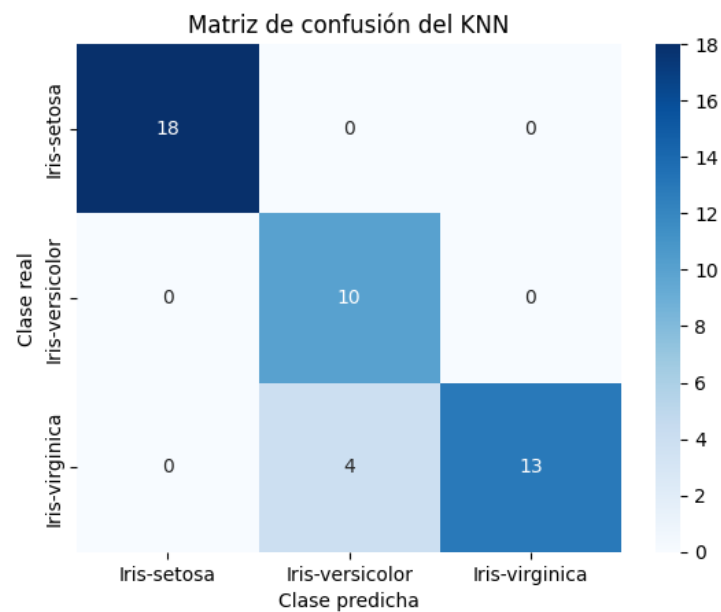


Figura 3: Matriz de confusión del modelo KNN (Python)

Precisión: 91.11 %

4.4. KNN implementado manualmente

```

1 import numpy as np
2 from collections import Counter
3
4 # Distancia euclidiana
5 def euclidean_distance(a, b):
6     return np.sqrt(np.sum((a - b) ** 2))
7
8 # Predicción para una instancia
9 def predict_one(x, X_train, y_train, k):
10     distances = [euclidean_distance(x, xi) for xi in X_train]
11     k_indices = np.argsort(distances)[:k]
12     k_labels = [y_train[i] for i in k_indices]
13     return Counter(k_labels).most_common(1)[0][0]
14
15 # Predicción para el conjunto de prueba
16 y_pred = [predict_one(x, X_train, y_train, k=27) for x in X_test]
17
18 # Precisión
19 accuracy = np.mean(y_pred == y_test) * 100
20 print(f"Precisión KNN manual: {accuracy:.2f}%")

```

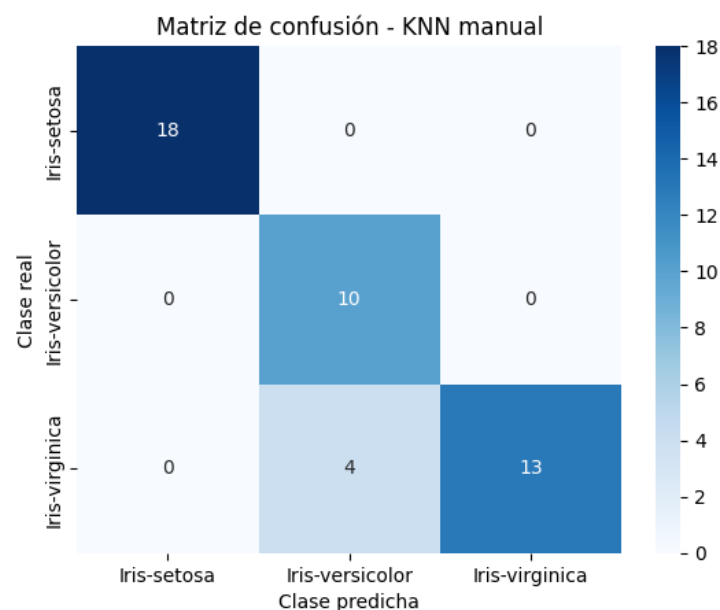


Figura 4: Matriz de confusión del modelo KNN sin librerías (Python)

Precisión: 91.11 %

5. Implementación en R (Posit Cloud)

5.1. Matriz de correlación

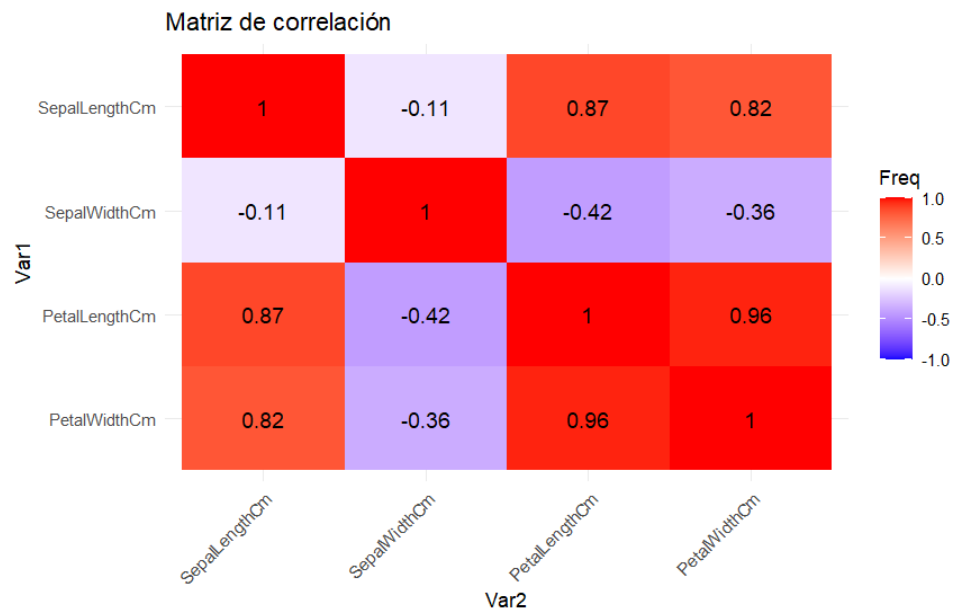


Figura 5: Matriz de correlación en R

5.2. Gráfico de dispersión

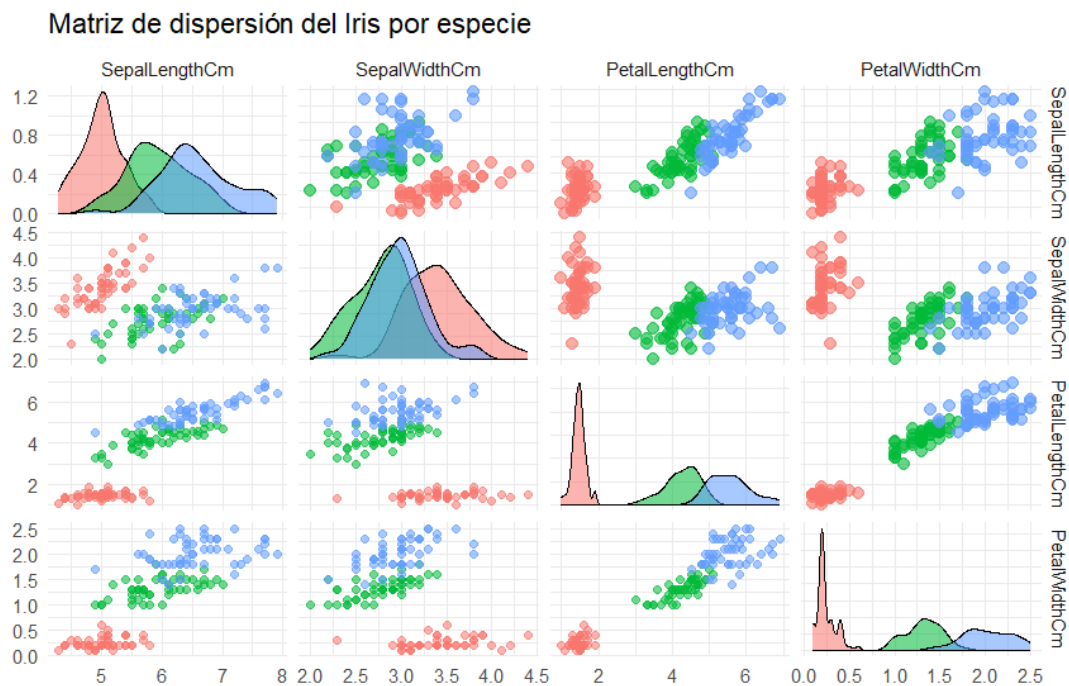


Figura 6: Distribución de características por especie (R)

5.3. KNN con librerías

```

1 if(!require(class)) install.packages("class")
2 library(class)
3
4 # Variables predictoras y variable objetivo
5 X <- dataset[, 1:4]
6 y <- dataset$Species
7 # Divisi n en entrenamiento y prueba
8 set.seed(123)
9 train_index <- sample(1:nrow(dataset), 0.7 * nrow(dataset))
10 X_train <- dataset[train_index, 1:4]
11 X_test <- dataset[-train_index, 1:4]
12 y_train <- dataset$Species[train_index]
13 y_test <- dataset$Species[-train_index]
14
15 # Entrenamiento y predicci n con KNN
16 library(class)
17 y_pred <- knn(X_train, X_test, y_train, k = 27)
18
19 # Precisi n
20 accuracy <- mean(y_pred == y_test)
21 cat("Precisi n del modelo KNN:", round(accuracy * 100, 2), "%\n")

```

Precisi3n: 97.78%

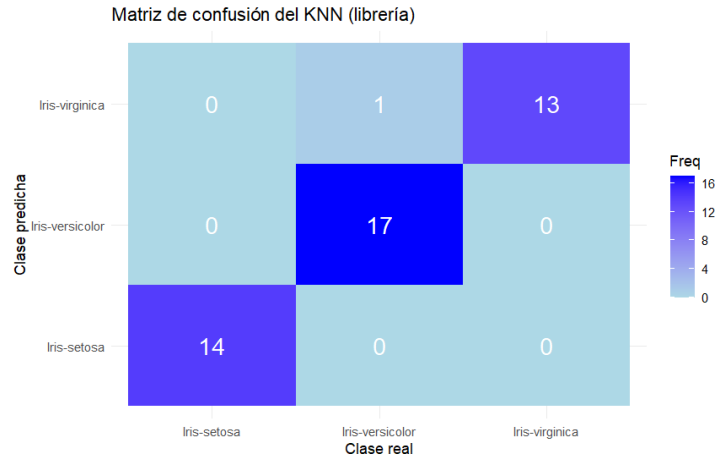


Figura 7: Matriz de confusi3n del modelo KNN (R)

5.4. KNN sin librerías (manual)

```

1 euclid <- function(a, b){
2   sqrt(sum((a - b)^2))
3 }
4
5 k <- 27
6 y_pred_manual <- vector()
7 for(i in 1:nrow(X_test)){
8   distances <- apply(X_train, 1, function(row) euclid(row, X_test
9     [i, ]))
10  nearest <- order(distances)[1:k]
11  labels <- y_train[nearest]
12  y_pred_manual[i] <- names(sort(table(labels), decreasing = TRUE
13    ))[1]
14 }
15
16 y_pred_manual <- factor(y_pred_manual, levels = levels(y_test))
17 accuracy_manual <- mean(y_pred_manual == y_test)
18 cat("Precisi n del KNN manual:", round(accuracy_manual * 100, 2)
19   , "%\n")

```

Precisión: 97.78 %

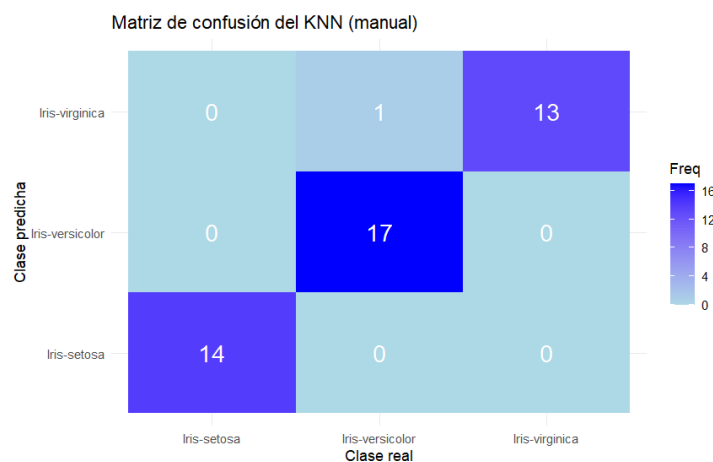


Figura 8: Matriz de confusión del modelo KNN (R)

6. Implementación en Java (Weka)

6.1. Matriz de correlación

Cuadro 1: Matriz de correlación en Weka

	SepalLength	SepalWidth	PetalLength	PetalWidth
SepalLength	1,00	0,72	-0,40	0,88
SepalWidth	0,72	1,00	-0,11	0,87
PetalLength	-0,40	-0,11	1,00	-0,42
PetalWidth	0,88	0,87	-0,42	1,00

Interpretación: Se observan correlaciones altas entre `PetalLength` y `PetalWidth`, similares a Python y R, lo que indica que estas características son relevantes para la clasificación.

6.2. Gráfico de dispersión

Nota: Weka no genera gráficos de dispersión automáticamente. Para análisis visual se recomienda exportar los datos a Python o R. Aquí se omite el gráfico y se presentan los resultados directamente.

6.3. KNN con librerías

```

1 // Entrenamiento KNN con librería IBk (Weka)
2 int k = 5;
3 IBk knn = new IBk(k);
4 knn.buildClassifier(train);
5
6 // Evaluación del modelo
7 Evaluation eval = new Evaluation(train);
8 eval.evaluateModel(knn, test);
9
10 System.out.println("Precisión: " + String.format("%.2f", eval.
    pctCorrect()) + "%");
11 System.out.println(eval.toMatrixString());

```

Precisión: 97,78%

Cuadro 2: Matriz de confusión del modelo KNN (Weka)

Predicted \ Actual	Setosa	Versicolor	Virginica
Setosa	14	0	0
Versicolor	0	19	0
Virginica	0	0	12

6.4. KNN implementado manualmente

Implementación propia de la distancia euclidiana:

```

1  int k = 27;
2  int aciertos = 0;
3
4  for (int i = 0; i < test.numInstances(); i++) {
5      Instance actual = test.instance(i);
6
7      // Calcular distancia a todos los ejemplos de entrenamiento
8      double[] distancias = new double[train.numInstances()];
9      String[] clases = new String[train.numInstances()];
10     for (int j = 0; j < train.numInstances(); j++) {
11         distancias[j] = euclideanDistance(actual, train.instance(j)
12         );
13         clases[j] = train.instance(j).stringValue(train.classIndex
14         ());
15     }
16
17     // Seleccionar los k vecinos m s cercanos
18     int[] vecinos = new int[k];
19     for (int v = 0; v < k; v++) {
20         double minDist = Double.MAX_VALUE;
21         int minIndex = -1;
22         for (int j = 0; j < distancias.length; j++) {
23             if (distancias[j] < minDist) {
24                 minDist = distancias[j];
25                 minIndex = j;
26             }
27         }
28         vecinos[v] = minIndex;
29         distancias[minIndex] = Double.MAX_VALUE; // evitar repetir
30     }
31
32     // Votaci n de las clases vecinas
33     int votosSetosa = 0, votosVersicolor = 0, votosVirginica = 0;
34     for (int v = 0; v < k; v++) {
35         String clase = clases[vecinos[v]];
36         if (clase.equals("Iris-setosa")) votosSetosa++;
37         else if (clase.equals("Iris-versicolor")) votosVersicolor
38             ++;
39         else votosVirginica++;
40     }
41
42     // Determinar clase m s votada
43     String prediccion;
44     if (votosSetosa > votosVersicolor && votosSetosa >
45         votosVirginica)
46         prediccion = "Iris-setosa";
47     else if (votosVersicolor > votosVirginica)
48         prediccion = "Iris-versicolor";

```

```

45     else
46         prediccion = "Iris-virginica";
47
48     // Comparar con clase real
49     if (prediccion.equals(actual.stringValue(test.classIndex())))
50         aciertos++;
51 }
52
53 double precision = 100.0 * aciertos / test.numInstances();
54 System.out.println("Precisi n KNN manual: " + String.format("%.2f"
    , precision) + "%");

```

Precisión manual: 97,78%

Cuadro 3: Matriz de confusión del modelo KNN manual (Weka)

Predicted \ Actual	Setosa	Versicolor	Virginica
Setosa	14	0	0
Versicolor	0	18	1
Virginica	0	0	12

Limitaciones de Weka:

- No genera gráficos de dispersión ni visualizaciones avanzadas de forma nativa.
- La matriz de correlación se presenta como tabla, no como heatmap.
- Más adecuado para evaluación rápida de modelos y pruebas de clasificación que para análisis gráfico exhaustivo.

7. Comparativa multilenguaje

Cuadro 4: Resumen comparativo entre lenguajes

Lenguaje	Precisión (librería)	Precisión (manual)	Ventajas	Desventajas
Python	91.11 %	91.11 %	Fácil de implementar, ecosistema amplio	Requiere optimización en grandes volúmenes de datos
R	97.78 %	97.78 %	Potente en análisis estadístico y visualización	Sintaxis menos intuitiva para principiantes
Java (Weka)	97.78 %	97.78 %	Eficiente, robusto y con entorno visual	Limitado en visualización avanzada y flexibilidad

Análisis: Los tres lenguajes presentan resultados consistentes, superando el 90 % de precisión. **R** y **Java (Weka)** alcanzan el mejor rendimiento (97.78 %), mientras que **Python** ofrece una implementación más sencilla y versátil. Las diferencias principales se deben a las bibliotecas disponibles y la facilidad para la visualización de resultados.

8. Conclusiones

El experimento demuestra que el dataset Iris puede analizarse eficazmente en múltiples lenguajes con resultados consistentes. Las diferencias radican principalmente en el entorno de desarrollo y la sintaxis:

- **Python** es ideal para aprendizaje automático y prototipado rápido.
- **R** sobresale en análisis estadístico y visualización de datos.
- **Java (Weka)** ofrece eficiencia y entorno gráfico accesible.

La implementación manual del algoritmo refuerza la comprensión teórica del KNN y su dependencia de la distancia euclidiana. En conjunto, los tres entornos confirman la capacidad del KNN para clasificar correctamente las especies del Iris con más del 90 % de precisión.