

# The iris data set

Martín Quijano - Martina Coletto

## Variables

- Sepal Length
- Sepal Width
- Petal Length
- Petal Width

## Modelos

K-Nearest Neighbors (KNN) y Regresión.

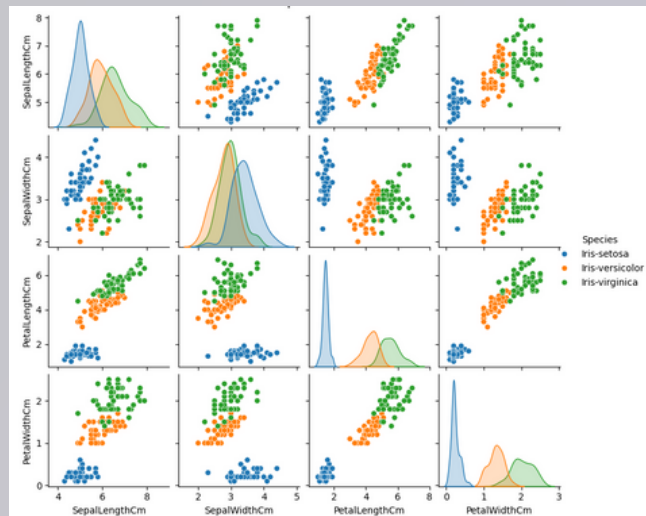
## Objetivo

Clasificar las flores del dataset Iris utilizando aprendizaje supervisado.

Vamos a analizar el dataset Iris utilizando técnicas de clasificación. Nuestro objetivo es predecir la especie de flor a partir de cuatro variables: largo y ancho del sépalo, y largo y ancho del pétalo. Para lograrlo, utilizamos dos modelos de aprendizaje supervisado: K-Nearest Neighbors y Regresión Logística.

## Pairplot

Para observar la distribución y posibles relaciones entre las variables.

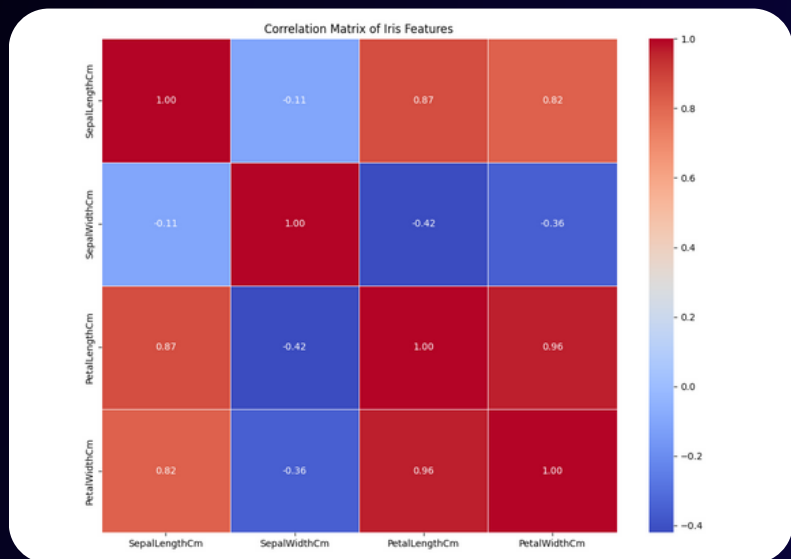


Con el pairplot vemos como se distribuye cada variable y como se relacionan entre si.

Vemos que en sepalLength y sepalWidth las variables tienen valores similares, hay superposición, mientras que en petalLength y petalWidth hay menos superposición, y vemos que setosa directamente no tiene superposición con otra, mientras que versicolor y virginica se superponen en algunos puntos

## Matriz de Correlación

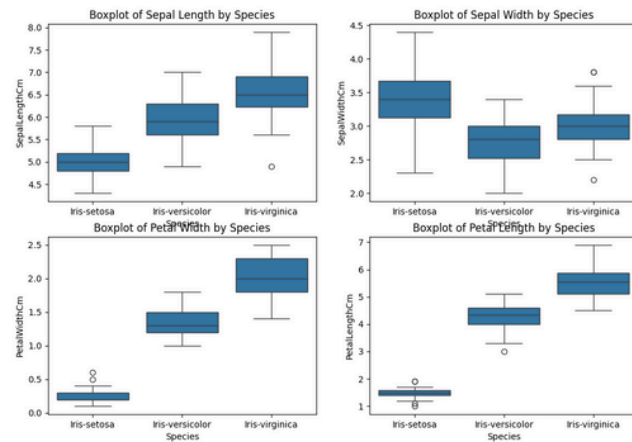
Para identificar dependencias entre las variables.



el largo del pétalo y el ancho del pétalo están altamente correlacionados, lo que nos indica que podrían aportar información redundante a los modelos. Sin embargo, esta correlación también puede ayudar a distinguir mejor entre especies

## Boxplot

Para detectar outliers y diferencias en la distribución de las clases.



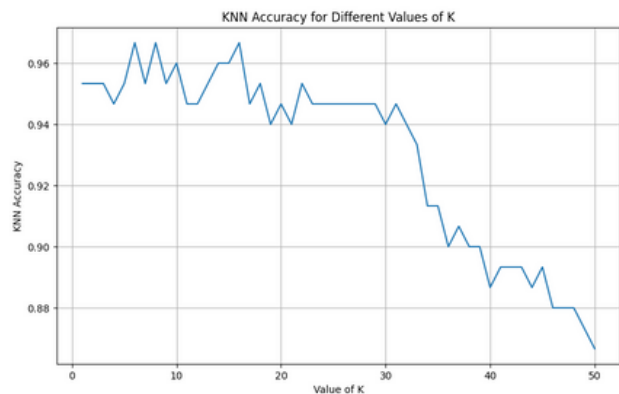
usamos boxplots para analizar la distribución de cada variable por especie. Esto nos permitió identificar posibles valores atípicos y diferencias significativas en la dispersión de las variables entre especies. Como se puede ver, la especie Setosa se diferencia claramente de las otras dos, mientras que Versicolor y Virginica tienen distribuciones más solapadas. Esto nos confirma lo que vimos en el pairplot



## Análisis con KNN

Elección del mejor valor de K mediante cross validation.

**Accuracy: 0.9667**

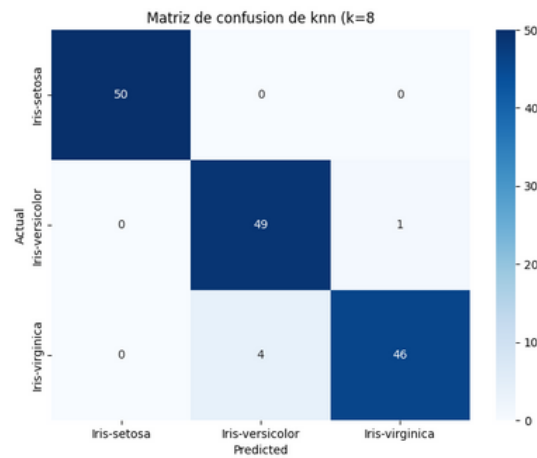


Encontramos que  $K=5$  ofrecía la mayor accuracy. Decidimos usar cross validation y no dividir la muestra entre test y train ya que la muestra no es tan grande, y no queríamos caer en un overfitting o underfitting



# Análisis con KNN

Matriz de confusión

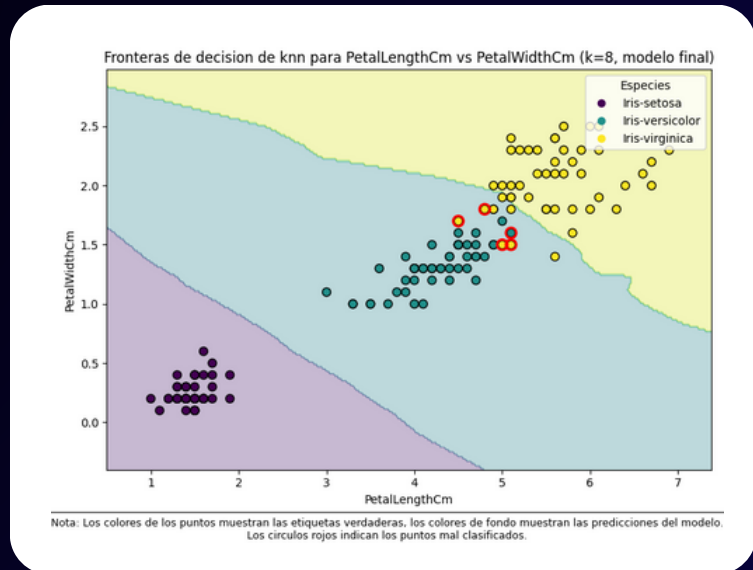


En la matriz de confusión podemos ver que la mayoría de las predicciones fueron correctas, con un accuracy del 96.67%. Esto indica que KNN funciona bien para este dataset. Adicionalmente, vemos que la setosa es clasificada bien siempre, mientras que en versicolor y virginica hay errores. Esto se debe a la superposición de los valores de las variables de esas especies



## Análisis con KNN

Fronteras de decisión



podemos ver cómo KNN separa las diferentes clases en el espacio de las variables. Observamos que las especies Setosa se separan claramente, pero hay cierto solapamiento entre Versicolor y Virginica, lo que explica los pocos errores de clasificación.

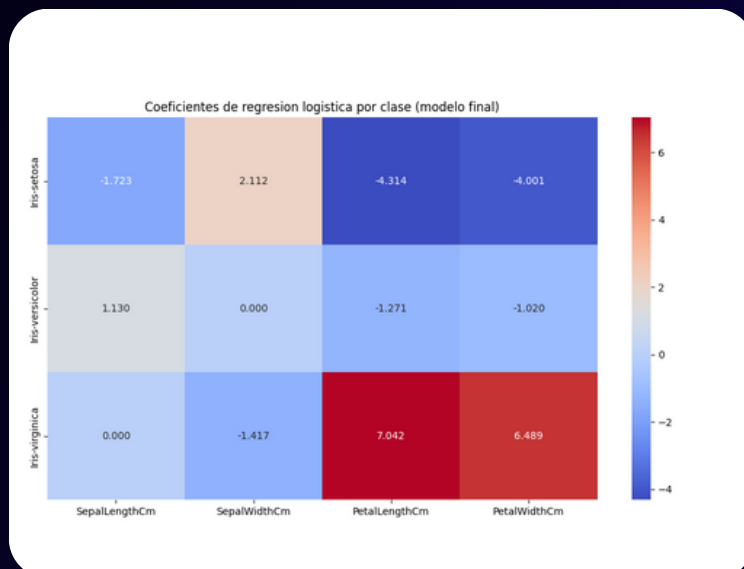
Es importante decir también que hay puntos de un color que están en el área de otro color. Esto se debe a que el gráfico solo considera dos de las 4 variables utilizadas





# Análisis de Regresión Logística

Evaluación de betas.  
**Accuracy: 0.98**



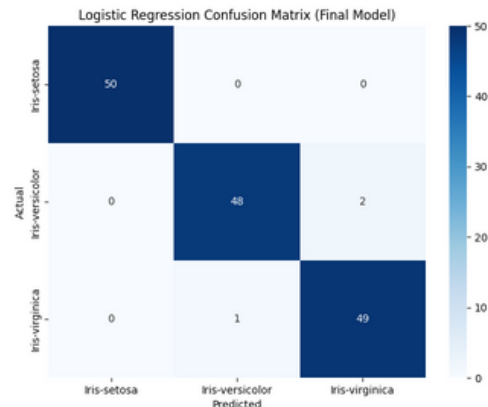
La accuracy obtenida con este modelo fue del 98%, ligeramente superior a la de KNN

Evaluamos los coeficientes beta y observamos dos principales características. En primer lugar, las variables de petal tienen más influencia (beta más alto) que las variables de sepal. En segundo lugar, vemos que para setosa y para virginica los betas tienen valores bastante altos, mientras que para versicolor los betas no tienen valores muy altos. Esto explica que la mayoría de errores se den versicolor.



# Análisis de Regresión Logística

Matriz de confusión

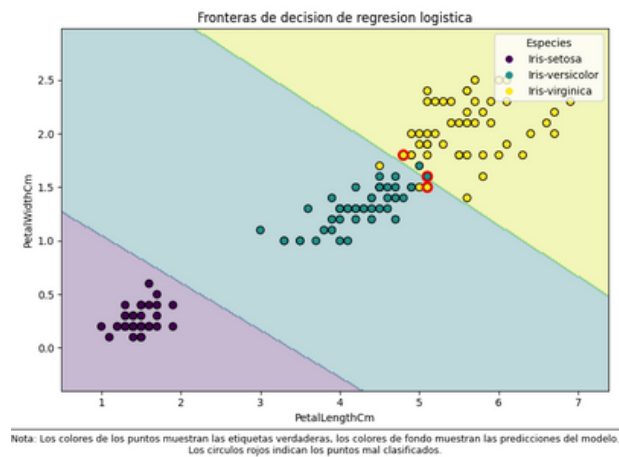


Menos errores que con KNN.



# Análisis de Regresión Logística

Fronteras de decision



En comparación con KNN, este modelo tiende a generar límites más lineales, lo que puede ser una ventaja cuando las clases son bien separables. Sin embargo, al igual que en KNN, seguimos viendo cierto solapamiento entre Versicolor y Virginica.

## Accuracy de los modelos con las transformaciones

accuracy de KNN por Transformación

Transformación	accuracy	
MinMax Scaling	0.9733	★
Avg Sepal & Petal	0.9733	★
Original	0.9667	
Ratios	0.9667	
Logarithmic	0.9600	
Avg Height & Width	0.9267	

accuracy de Regresión Logística por Transformación

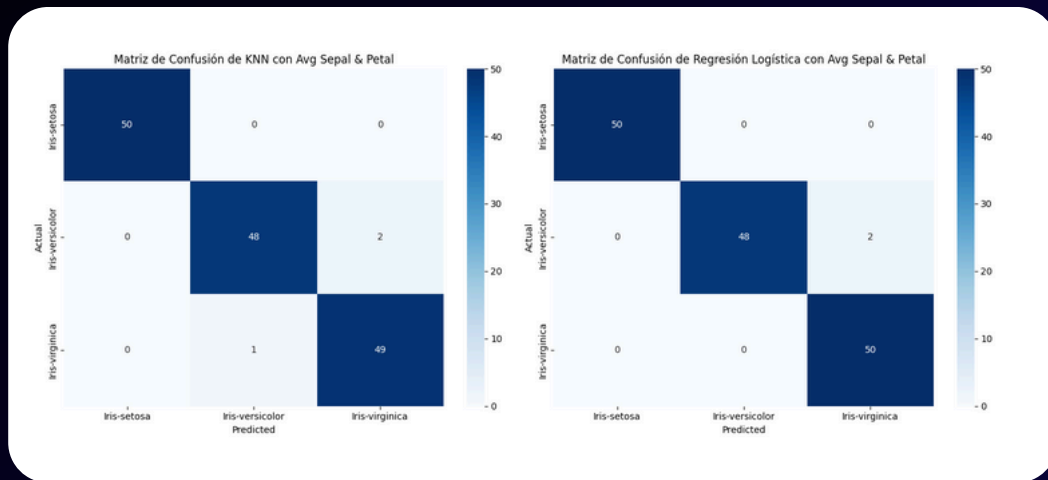
Transformación	accuracy	
Avg Sepal & Petal	0.9867	★
Original	0.9800	
MinMax Scaling	0.9800	
Logarithmic	0.9800	
Ratios	0.9600	
Avg Height & Width	0.9000	

Para mejorar la precisión de nuestros modelos, probamos distintas transformaciones.

Esas transformaciones fueron minMax scaling, promediar sepalLength con sepalWidth y promediar petalLength con petalWidth, calcular ratios de sepal, petal, largo y ancho, poner en escala logaritmica las variables y promediar sepalHeight con petalHeight y sepalWidth con petalHeight.

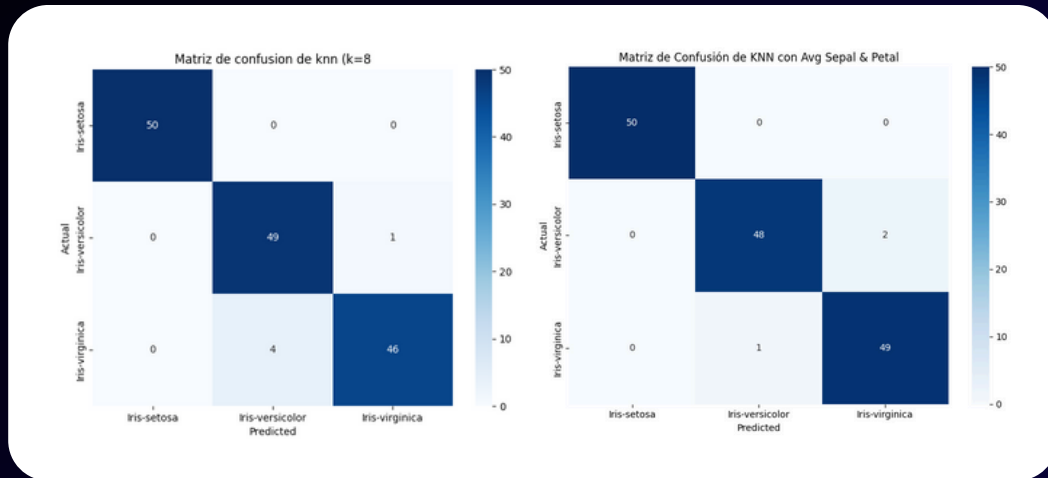
La mejor transformacion fue MinMaxScaler. Notamos que el accuracy en ambos modelos mejora un 0.0067 al ser evaluados con los datos transformados

## Matriz de confusión de los modelos con mejor transformación



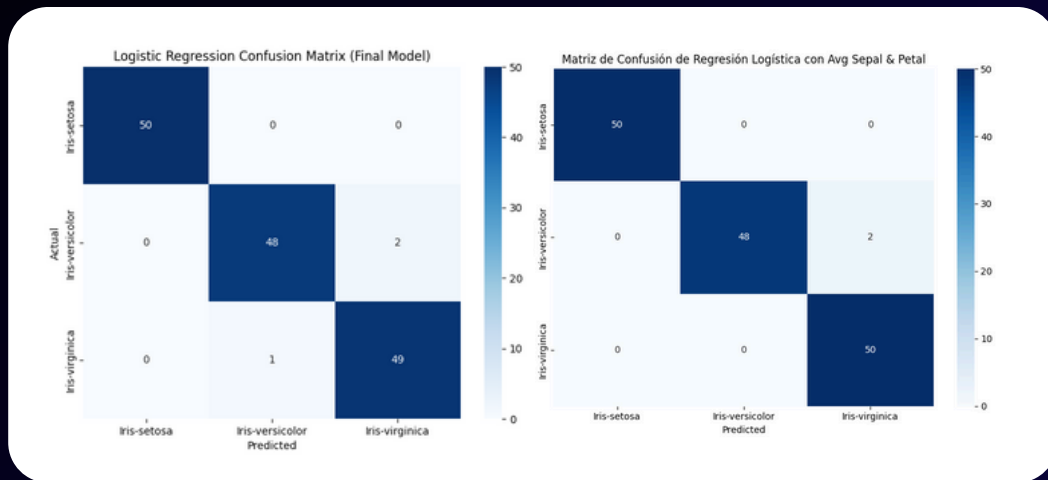
Vemos que con los datos transformados en knn hay 3 errores y en regresion solo 2

## KNN: normal vs. transformación



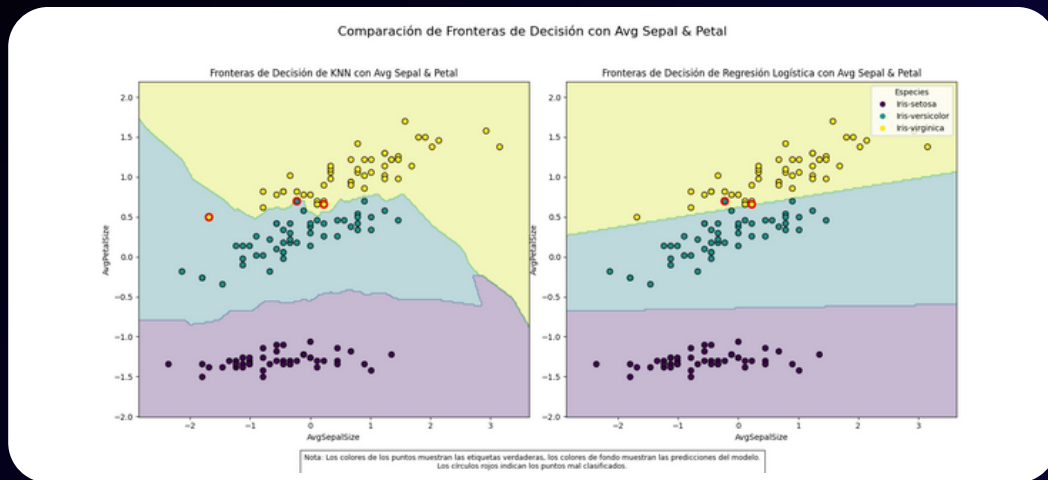
Vemos que con la transformación hay menos errores, pasamos de 5 errores a 3 errores en KNN

## Regresión Logística: normal vs. transformación



Vemos que con la transformación hay menos errores, pasamos de 3 errores a 2 errores en regresión logística

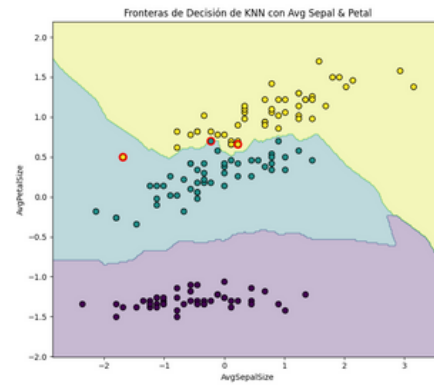
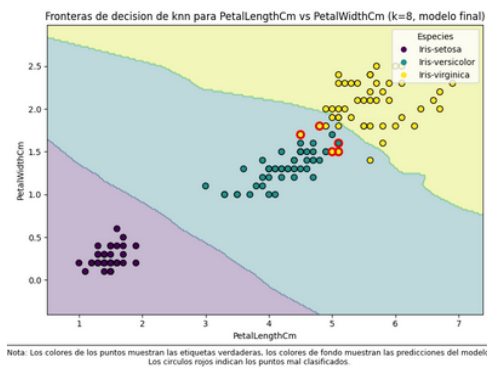
## Fronteras de decisión de los modelos con mejor transformación



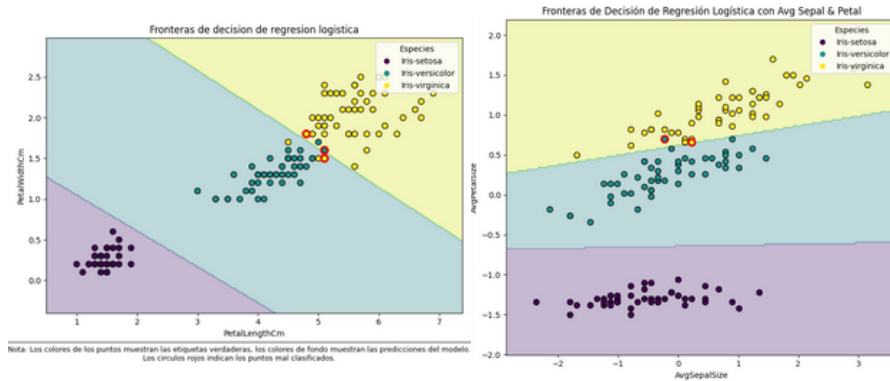
Al ver el gráfico de fronteras de decisión con datos transformados, vemos que en regresión logística las fronteras son más lineales que en KNN



## KNN: normal vs. transformación



## Regresión Logística: normal vs. transformación



## ¿Las observaciones son agrupables por especie?

Si, como vemos en la matriz de confusión y en las fronteras de decisión las observaciones son agrupables con un margen de error muy chico.

Para concluir, podemos decir que los datos del dataset Iris son altamente clasificables, con un margen de error pequeño. Tanto KNN como Regresión Logística lograron buenos resultados, con una ligera ventaja para la Regresión Logística. Las transformaciones de datos ayudaron a mejorar el rendimiento de ambos modelos, especialmente en la separación de Versicolor y Virginica.

Como posibles mejoras, podríamos probar modelos más avanzados como Random Forest para analizar si obtenemos una clasificación aún más precisa

# Gracias

Martín Quijano - Martina Coletto

-