

# Análisis exploratorios

Santos G

## Tabla de contenidos

<b>1</b>	<b>Contexto del proyecto</b>	<b>1</b>
<b>2</b>	<b>Carga y verificación inicial de datos</b>	<b>1</b>
<b>3</b>	<b>Matriz de correlaciones y distribuciones entre variables numéricas</b>	<b>3</b>
<b>4</b>	<b>Distribución de variables morfométricas entre especies</b>	<b>5</b>

```
1 # Librerías
2 library(tidyverse) # Manipulación de datos: dplyr, tidyr, readr
3 library(janitor)   # Limpieza: clean_names(), tabyl()
4 library(ggplot2)   # Gráficos profesionales
5 library(skimr)      # EDA rápido y completo (skim())
6 library(GGally)     # Matriz de gráficos para variables múltiples
7 library(knitr)      # Tablas en Quarto
8 library(kableExtra) # Tablas formateadas para informes
```

## 1 Contexto del proyecto

Se realizó una exploración y control de calidad de los datos de entrada para identificar variables relevantes, evaluar supuestos básicos y priorizar rutas analíticas. El objetivo es generar una guía reproducible que permita a futuros analistas (o a un equipo de consultoría) replicar y ampliar los análisis según objetivos específicos (p. ej. comparar tratamientos, modelar abundancias o construir índices de condición).

## 2 Carga y verificación inicial de datos

```
1 #|label: data-load
2 # Carga de datos (ejemplo iris) y limpieza mínima
3 data("iris")
4 df <- as_tibble(iris) %>%
5   janitor::clean_names() # convierte a snake_case: sepal_length, etc.
6
7 # Información básica
```

```

8 n_rows <- nrow(df); n_cols <- ncol(df)
9 glimpse(df)

```

```

Rows: 150
Columns: 5
$ sepal_length <dbl> 5.1, 4.9, 4.7, 4.6, 5.0, 5.4, 4.6, 5.0, 4.4, 4.9, 5.4, 4.~
$ sepal_width  <dbl> 3.5, 3.0, 3.2, 3.1, 3.6, 3.9, 3.4, 3.4, 2.9, 3.1, 3.7, 3.~
$ petal_length <dbl> 1.4, 1.4, 1.3, 1.5, 1.4, 1.7, 1.4, 1.5, 1.4, 1.5, 1.5, 1.~
$ petal_width  <dbl> 0.2, 0.2, 0.2, 0.2, 0.2, 0.4, 0.3, 0.2, 0.2, 0.1, 0.2, 0.~
$ species      <fct> setosa, setosa, setosa, setosa, setosa, setosa, setosa, s~

```

```

1 skim(df) # Resumen compacto por variable

```

Tabla 1: Data summary

Name	df
Number of rows	150
Number of columns	5
Column type frequency:	
factor	1
numeric	4
Group variables	None

#### Variable type: factor

skim_variable	n_missing	complete_rate	ordered	n_unique	top_counts
species	0	1	FALSE	3	set: 50, ver: 50, vir: 50

#### Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100	hist
sepal_length	0	1	5.84	0.83	4.3	5.1	5.80	6.4	7.9	
sepal_width	0	1	3.06	0.44	2.0	2.8	3.00	3.3	4.4	
petal_length	0	1	3.76	1.77	1.0	1.6	4.35	5.1	6.9	
petal_width	0	1	1.20	0.76	0.1	0.3	1.30	1.8	2.5	

El dataset contiene **N = 150** observaciones y **5** variables. Las variables cuantitativas son: **sepal\_length**, **sepal\_width**, **petal\_length**, **petal\_width** (continuas, en cm). La variable categórica **species** indica tres grupos balanceados (n = 50 por grupo). No se detectaron valores faltantes ni duplicados tras una inspección inicial. Esta estructura (muestras balanceadas y variables continuas sin NA) permite aplicar análisis univariados, comparativos y multivariados con mínima preprocesamiento.

La variable **sepal\_length** presenta una media **5.84 cm**, desviación estándar (SD) **0.83 cm**, rango 4.3–7.9.

La variable **sepal\_width** presenta una media **3.06 cm**, SD **0.44 cm**, rango 2.0–4.4.

La variable **petal\_length** presenta una media media **3.76 cm**, SD **1.77 cm**, rango 1.0–6.9.

La variable **petal\_width** presenta una media **1.20 cm**, SD **0.76 cm**, rango 0.1–2.5.

Además, el conjunto de datos presenta las siguientes características:

- **Escala homogénea de medidas:** todas las variables cuantitativas están en la misma unidad (centímetros), lo que facilita comparaciones directas y análisis multivariados sin necesidad inmediata de reescalado.
- **Posible colinealidad:** inspecciones preliminares sugieren que **petal\_length** y **petal\_width** están altamente correlacionadas, lo que puede condicionar modelos de regresión o técnicas de reducción de dimensionalidad (ej. PCA).
- **Grupos biológicos claros:** las especies representan categorías naturales y balanceadas, condición poco común en datos ecológicos reales donde suele haber desbalance → este dataset es ideal como caso didáctico y de referencia.
- **Potencial de discriminación:** dado el balance de clases y la separación conocida entre especies de Iris, el dataset es adecuado para ilustrar desde comparaciones descriptivas hasta modelos predictivos (GLM, LDA, clustering).

### 3 Matriz de correlaciones y distribuciones entre variables numéricas

```
1 num_df <- df %>% select(where(is.numeric))
2 GGally::ggpairs(num_df, upper = list(continuous = wrap("cor", size = 3)))
```

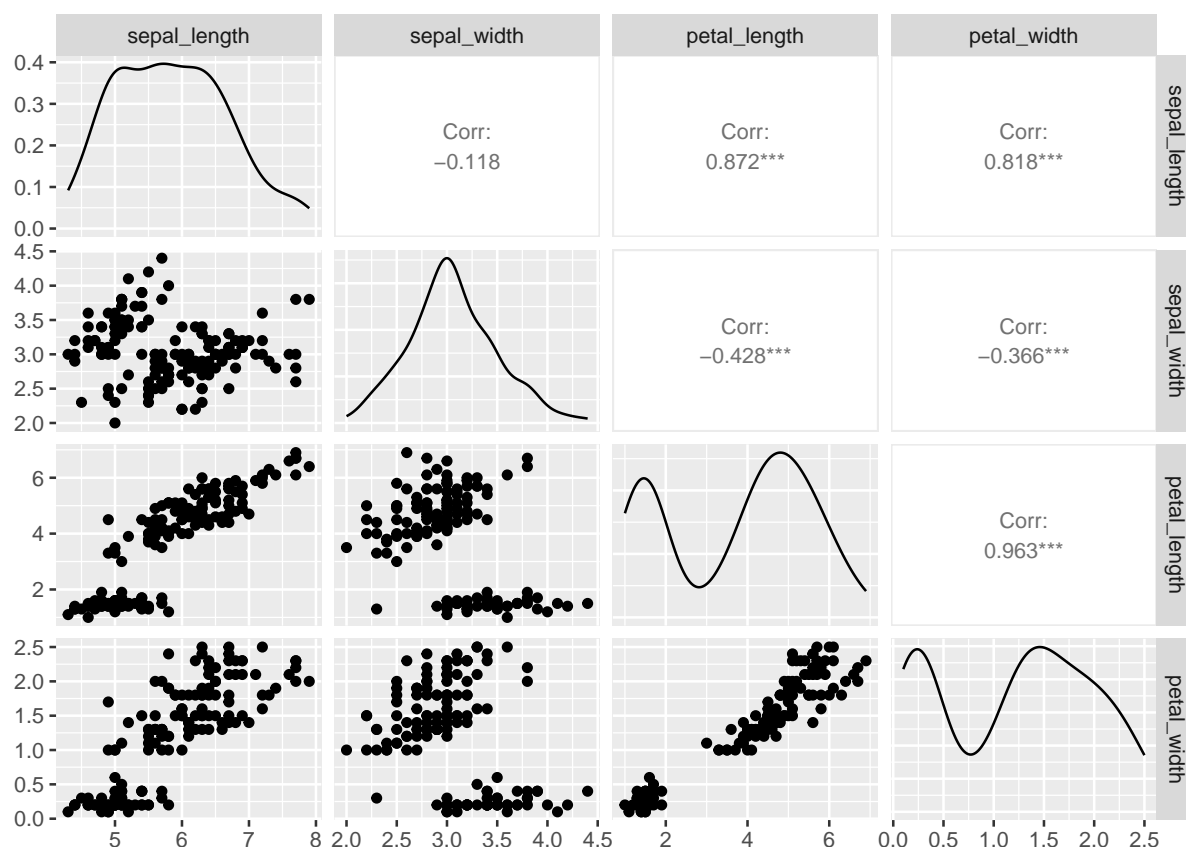


Figura 1: Matriz de dispersión y correlación de las variables cuantitativas.

La matriz de gráficos (scatterplots + histogramas + coeficientes de correlación) permite evaluar de un vistazo correlaciones, distribuciones marginales y relaciones lineales o no lineales entre variables.

- **Colinealidad fuerte:** `petal_length` y `petal_width` muestran correlación muy alta ( $r > 0.9$ ), lo que evidencia redundancia informativa.
- **Correlaciones moderadas:** `sepal_length` correlaciona de forma consistente con `petal_length` ( $r = 0.87$ ), lo que sugiere que ambas crecen conjuntamente y pueden estar asociadas a un **gradiente de tamaño corporal general**.
- **Menor asociación:** `sepal_width` se muestra menos correlacionada con las demás, lo que indica que puede aportar información diferenciada en análisis multivariados.
- **Separación por especie:** en las distribuciones univariadas y bivariadas se observa que las especies tienden a formar agrupamientos diferenciados, especialmente *setosa*, lo que sugiere potencial de clasificación/discriminación.
- **Linealidad vs no linealidad:** las nubes de puntos muestran relaciones aproximadamente lineales, aunque con cierto solapamiento entre *versicolor* y *virginica*, lo cual podría requerir modelos más flexibles en análisis predictivos.
- **Distribución:** Para `sepal_length` es aproximadamente simétrica, con dispersión moderada. No aparecen valores extremos. `sepal_width` presenta una dispersión menor que en `sepal_length`; sin embargo, la distribución presenta cierta asimetría negativa (colas hacia valores pequeños). Finalmente tanto `petal_length` como `petal_width` muestran

una distribución marcadamente bimodal, porque las especies difieren mucho en esta variable. Es la de mayor varianza relativa.

En presencia de **alta colinealidad** ( $r > 0.8$ ), considerar retener solo una variable del par altamente correlacionado si el objetivo es simplificación. Emplear técnicas multivariadas (PCA, análisis discriminante) para reducir dimensionalidad y evitar inestabilidad en modelos de regresión.

## 4 Distribución de variables morfométricas entre especies

```
1 # Pasar el dataset a formato largo
2 iris_long <- df %>%
3   pivot_longer(cols = -species,
4                 names_to = "Variable",
5                 values_to = "Valor")
6
7 # Gráfico unificado
8 ggplot(iris_long, aes(x = species, y = Valor, fill = species)) +
9   geom_boxplot(outlier.shape = 21, alpha = 0.7) +
10  facet_wrap(~ Variable, scales = "free_y") +
11  labs(
12    title = "Comparación de variables morfométricas en especies de Iris",
13    x = "Especies",
14    y = "Valor (cm)"
15  ) +
16  theme_minimal(base_size = 13) +
17  theme(
18    plot.title = element_text(hjust = 0.5, face = "bold"),
19    legend.position = "none",
20    strip.text = element_text(face = "bold")
21  )
```

## Comparación de variables morfológicas en especies de Iris

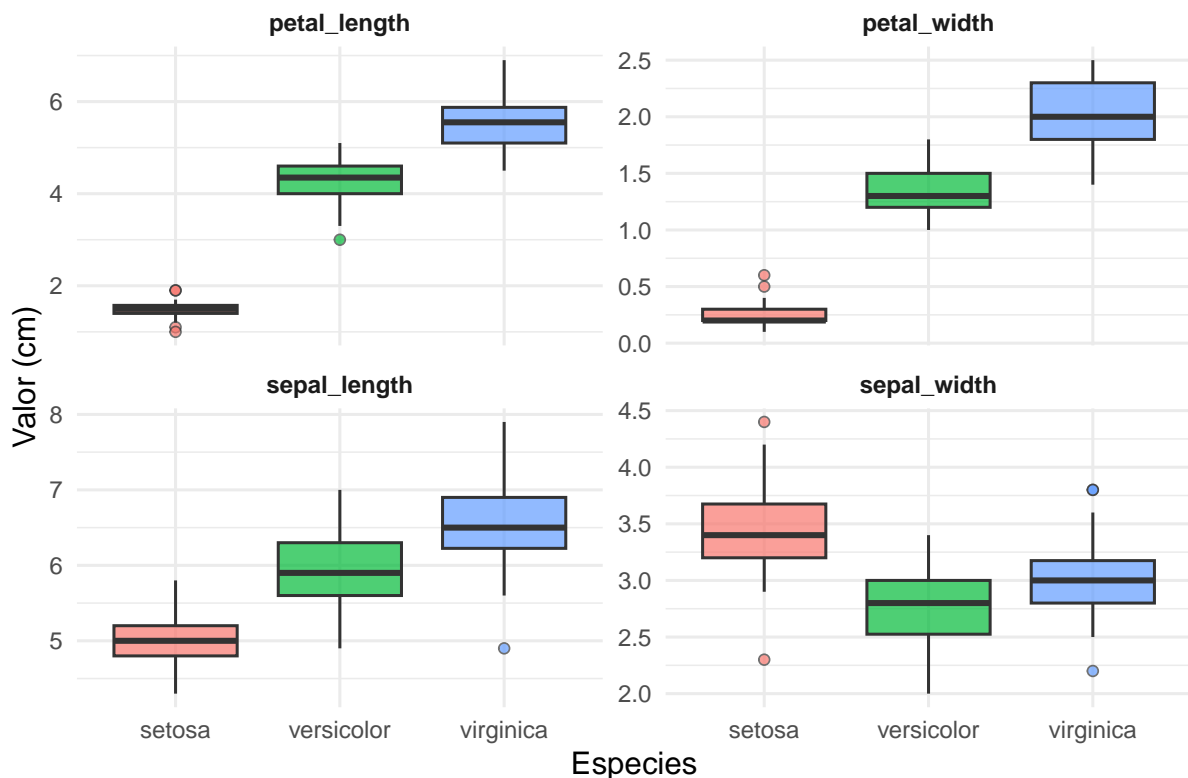


Figura 2: Distribución de variables morfológicas en tres especies de Iris.

En cuanto a la **longitud del sépalo (Sepal.Length)**, se observa que *Iris setosa* presenta los valores más bajos, con una mediana cercana a los 5 cm, claramente diferenciada de las otras especies. Por su parte, *I. versicolor* e *I. virginica* muestran longitudes mayores y con bastante solapamiento, lo que dificulta separarlas únicamente con esta variable. En consecuencia, Sepal.Length resulta útil para discriminar a *I. setosa*, pero limitado para distinguir entre *versicolor* y *virginica*.

La **anchura del sépalo (Sepal.Width)** muestra distribuciones más amplias y con solapamiento importante entre las tres especies. Aunque *setosa* tiende a tener sépalos ligeramente más anchos en promedio, la alta variabilidad reduce su capacidad discriminante. Así, esta variable por sí sola tiene bajo poder para diferenciar especies.

En el caso de la **longitud del pétalo (Petal.Length)**, se aprecia una separación casi perfecta: *setosa* forma un grupo completamente aislado con valores mucho menores, mientras que *versicolor* y *virginica* presentan medianas diferenciadas, aunque con cierto solapamiento en sus rangos intercuartílicos. Esto convierte a Petal.Length en una de las variables más informativas y clave para la clasificación.

Finalmente, la **anchura del pétalo (Petal.Width)** también resulta altamente discriminante. *Setosa* presenta valores notablemente bajos y sin superposición con las otras especies, mientras que *versicolor* y *virginica* se diferencian con claridad, aunque con menor margen que en el caso de Petal.Length. Esta variable se perfila como una de las más robustas para separar especies, especialmente en conjunto con Petal.Length.