

# Taller de Análisis Exploratorio de Datos (EDA)

## con R Studio

Predicción del Rendimiento Académico Estudiantil

---

Colegio Universitario de Cartago

Curso: Introducción a la Estadística / Probabilidad y Estadística

Profesor: David Martínez Salazar  
III Cuatrimestre 2025

## 1 Introducción al Taller

### 1.1 Objetivo General

En este taller práctico, se trabajará con un conjunto de datos real que contiene registros académicos de **2000 estudiantes**. El estudiante aprenderá a realizar un **Análisis Exploratorio de Datos (EDA)** utilizando R Studio, desde la importación de datos hasta la construcción de un modelo de regresión lineal simple.

### 1.2 ¿Qué es el EDA?

#### ♦ Concepto Teórico

El **Análisis Exploratorio de Datos (EDA)** es el proceso de examinar y resumir las características principales de un conjunto de datos. Sus objetivos son:

- Comprender la estructura y contenido de los datos
- Detectar errores, valores faltantes o atípicos
- Calcular estadísticas descriptivas (media, mediana, moda, etc.)
- Visualizar patrones mediante gráficos
- Verificar supuestos estadísticos (como la normalidad)
- Identificar relaciones entre variables

### 1.3 Objetivos Específicos

Al finalizar este taller, el estudiante será capaz de:

1. Importar y explorar datos en R Studio
2. Calcular medidas de tendencia central y dispersión
3. Crear gráficos estadísticos (pastel, barras, histogramas, boxplots)
4. Evaluar la normalidad de una variable
5. Analizar correlación entre dos variables
6. Construir un modelo de regresión lineal simple

### 1.4 Modalidad de Trabajo

- **Modalidad:** Individual o en parejas
- **Herramienta:** R Studio
- **Duración estimada:** 2-3 horas
- **Entregable:** Script de R (.R) con todos los análisis

## 2 Descripción del Conjunto de Datos

### 2.1 Contexto del Dataset

El dataset “**Student Performance Prediction**” fue obtenido de la plataforma Kaggle y contiene información académica de **2000 estudiantes**. Este conjunto de datos permite analizar y predecir las **calificaciones del examen final** basándose en diversos indicadores académicos.

#### ♦ Concepto Teórico

##### **¿Por qué es importante este análisis?**

En el ámbito educativo, identificar los factores que influyen en el rendimiento académico permite:

- Diseñar intervenciones tempranas para estudiantes en riesgo
- Optimizar estrategias de enseñanza-aprendizaje
- Comprender qué hábitos conducen al éxito estudiantil
- Desarrollar sistemas de alerta temprana basados en datos

### 2.2 Variables del Dataset

El dataset contiene las siguientes variables:

Variable Original	Tipo	Descripción
Student_ID	Identificador	Código único del estudiante
Attendance_Rate	Numérica (0-100)	Porcentaje de asistencia a clases
Internal_Test_1	Numérica (0-100)	Calificación del primer examen interno
Internal_Test_2	Numérica (0-100)	Calificación del segundo examen interno
Assignment_Avg	Numérica (0-100)	Promedio de calificaciones en tareas
Study_Hours_Per_Day	Numérica	Horas diarias dedicadas al estudio
Final_Exam_Mark	Numérica (0-100)	Calificación del examen final

Table 1: Descripción de las variables del dataset

### 2.3 Flujo del Análisis EDA

En este taller se seguirá el siguiente orden:

1. **Importación de datos:** Cargar el dataset en R
2. **Exploración inicial:** Ver estructura y primeros registros
3. **Limpieza de datos:** Renombrar columnas, verificar valores faltantes

4. **Estadística descriptiva:** Calcular media, mediana, moda, etc.
5. **Visualización:** Gráficos de pastel, barras y boxplots
6. **Análisis de normalidad:** Histograma, Q-Q plot, pruebas estadísticas
7. **Correlación:** Relación entre asistencia y nota final
8. **Regresión lineal:** Modelo predictivo simple

### 3 Configuración del Entorno de Trabajo

#### ✓ Instrucciones

Antes de comenzar, se deben instalar y cargar las librerías necesarias. Es importante seguir cada paso cuidadosamente.

#### 3.1 Paso 1: Instalación de Paquetes

##### ★ Nota Importante

Solo es necesario instalar los paquetes UNA VEZ. Si ya se encuentran instalados, se puede omitir este paso.

```
1 # Instalar paquetes necesarios
2 install.packages("ggplot2")      # Para graficos
3 install.packages("dplyr")        # Para manipular datos
4 install.packages("corrplot")     # Para mapa de calor
5 install.packages("reshape2")      # Para reestructurar datos
```

Listing 1: Instalación de paquetes (ejecutar solo una vez)

#### 3.2 Paso 2: Cargar las Librerías

```
1 # Cargar librerias
2 library(ggplot2)
3 library(dplyr)
4 library(corrplot)
5 library(reshape2)
6
7 # Mostrar numeros sin notacion cientifica
8 options(scipen = 999)
```

Listing 2: Cargar librerías al inicio de cada sesión

#### 3.3 Paso 3: Importar el Dataset

```
1 # Importar el dataset
2 # NOTA: Se debe ajustar la ruta segun la ubicacion del archivo
3 datos <- read.csv("student_performance.csv")
4
5 # Verificar que se cargo correctamente
6 cat("El dataset tiene", nrow(datos), "filas y", ncol(datos),
    "columnas\n")
```

Listing 3: Importar el archivo CSV

#### 📎 Actividad

Se debe ejecutar el código anterior y anotar cuántas filas y columnas tiene el dataset.

✍ Espacio para Respuesta

Número de filas: \_\_\_\_\_

Número de columnas: \_\_\_\_\_

## 4 Exploración Inicial de los Datos

### 4.1 Ver las Primeras y Últimas Filas

```

1 # Ver las primeras 6 filas
2 head(datos)
3
4 # Ver las ultimas 6 filas
5 tail(datos)

```

Listing 4: Explorar las primeras filas del dataset

### 4.2 Estructura del Dataset

```

1 # Ver estructura: tipos de variables
2 str(datos)
3
4 # Ver nombres de las columnas
5 names(datos)

```

Listing 5: Ver la estructura de los datos

#### ☞ Actividad

Se debe ejecutar el código `str(datos)` y responder:

1. ¿Cuántas variables numéricas existen?
2. ¿Existe alguna variable de tipo texto (character)?

#### ✍ Espacio para Respuesta

1. Variables numéricas: \_\_\_\_\_

2. Variables de texto: \_\_\_\_\_

### 4.3 Resumen Estadístico Rápido

```

1 # Resumen de todas las variables
2 summary(datos)

```

Listing 6: Resumen estadístico general

#### ☞ Actividad

Con el resultado de `summary(datos)`, se deben anotar los valores para la variable `Final_Exam_Mark`:

✍ Espacio para Respuesta

Mínimo: \_\_\_\_\_ Máximo: \_\_\_\_\_

Mediana: \_\_\_\_\_ Media: \_\_\_\_\_

## 5 Limpieza y Preparación de Datos

### 5.1 Renombrar las Columnas

Para facilitar el trabajo, se renombrarán las columnas a español.

```

1 # Ver nombres actuales
2 names(datos)
3
4 # Renombrar las columnas
5 colnames(datos) <- c(
6   "ID",                      # Student_ID
7   "Asistencia",              # Attendance_Rate
8   "Examen_Interno_1",        # Internal_Test_1
9   "Examen_Interno_2",        # Internal_Test_2
10  "Promedio_Tareas",         # Assignment_Avg
11  "Horas_Estudio",          # Study_Hours_Per_Day
12  "Nota_Final"             # Final_Exam_Mark
13 )
14
15 # Verificar los nuevos nombres
16 names(datos)

```

Listing 7: Renombrar columnas a español

### 5.2 Verificar Valores Faltantes

```

1 # Contar valores NA en cada columna
2 colSums(is.na(datos))
3
4 # Total de valores faltantes
5 cat("Total de valores faltantes:", sum(is.na(datos)), "\n")

```

Listing 8: Buscar valores faltantes (NA)

#### 💡 Actividad

¿Existen valores faltantes en el dataset? ¿En cuáles variables?

#### ✍ Espacio para Respuesta

¿Existen valores faltantes? (Sí / No): \_\_\_\_\_

¿En cuáles variables? \_\_\_\_\_

## 6 Estadística Descriptiva: Nota del Examen Final

El análisis se enfocará en la variable **Nota\_Final** (calificación del examen final).

### ◆ Concepto Teórico

#### Medidas de Tendencia Central:

- **Media ( $\bar{x}$ )**: Promedio de todos los valores
- **Mediana**: Valor que divide los datos a la mitad
- **Moda**: Valor que más se repite

#### Medidas de Dispersión:

- **Varianza ( $s^2$ )**: Indica qué tan dispersos están los datos
- **Desviación Estándar ( $s$ )**: Raíz cuadrada de la varianza
- **Cuartiles**: Dividen los datos en 4 partes iguales

### 6.1 Medidas de Tendencia Central

```

1 # MEDIA (promedio)
2 media <- mean(datos$Nota_Final)
3 cat("Media:", round(media, 2), "\n")
4
5 # MEDIANA (valor central)
6 mediana <- median(datos$Nota_Final)
7 cat("Mediana:", round(mediana, 2), "\n")
8
9 # MODA (valor mas frecuente)
10 # Se crea una tabla de frecuencias y se busca el maximo
11 tabla_freq <- table(round(datos$Nota_Final, 0))
12 moda <- as.numeric(names(tabla_freq)[which.max(tabla_freq)])
13 cat("Moda:", moda, "\n")

```

Listing 9: Calcular media, mediana y moda de Nota Final

### ✉ Actividad

Se deben calcular las medidas de tendencia central para **Nota\_Final** y anotar los resultados:

### ✍ Espacio para Respuesta

**Media:** \_\_\_\_\_

**Mediana:** \_\_\_\_\_

**Moda:** \_\_\_\_\_

## 6.2 Medidas de Dispersión

```

1 # VARIANZA
2 varianza <- var(datos$Nota_Final)
3 cat("Varianza:", round(varianza, 2), "\n")
4
5 # DESVIACION ESTANDAR
6 desv_est <- sd(datos$Nota_Final)
7 cat("Desviacion Estandar:", round(desv_est, 2), "\n")
8
9 # CUARTILES
10 cuartiles <- quantile(datos$Nota_Final, probs = c(0.25, 0.50, 0.75))
11 cat("Q1 (25%):", round(cuartiles[1], 2), "\n")
12 cat("Q2 (50%):", round(cuartiles[2], 2), "\n")
13 cat("Q3 (75%):", round(cuartiles[3], 2), "\n")
14
15 # RANGO
16 rango <- max(datos$Nota_Final) - min(datos$Nota_Final)
17 cat("Rango:", round(rango, 2), "\n")

```

Listing 10: Calcular varianza, desviación estándar y cuartiles

### ☞ Actividad

Se deben calcular las medidas de dispersión para **Nota\_Final**:

### ✍ Espacio para Respuesta

Varianza: \_\_\_\_\_ Desviación Estándar: \_\_\_\_\_

Q1: \_\_\_\_\_ Q2: \_\_\_\_\_ Q3: \_\_\_\_\_

Rango: \_\_\_\_\_

## 7 Visualización de Datos

### 7.1 Gráfico de Pastel: Categorías de Rendimiento

Primero, se crearán categorías basadas en la Nota Final.

```

1 # Crear categorias de rendimiento
2 datos$Rendimiento <- cut(datos$Nota_Final ,
3   breaks = c(0, 60, 70, 80, 90, 100),
4   labels = c("Reprobado", "Suficiente", "Bueno", "Muy Bueno",
5   "Excelente"),
6   include.lowest = TRUE
7 )
8
8 # Contar estudiantes por categoria
9 frecuencias <- table(datos$Rendimiento)
10 porcentajes <- round(prop.table(frecuencias) * 100, 1)
11
12 # Crear etiquetas con porcentajes
13 etiquetas <- paste(names(frecuencias), "\n", porcentajes, "%", sep = "")
14
15 # Definir colores
16 colores <- c("#FF6B6B", "#FFA94D", "#FFD93D", "#6BCB77", "#4D96FF")
17
18 # Crear grafico de pastel
19 pie(frecuencias,
20   labels = etiquetas,
21   col = colores,
22   main = "Distribucion del Rendimiento Academicoo")
```

Listing 11: Crear categorías y gráfico de pastel

#### 🔗 Actividad

Se debe observar el gráfico de pastel y responder:

1. ¿Cuál es la categoría con más estudiantes?
2. ¿Qué porcentaje de estudiantes reprobó?

#### ✍ Espacio para Respuesta

1. Categoría más frecuente: \_\_\_\_\_
2. Porcentaje de reprobados: \_\_\_\_\_

## 7.2 Diagrama de Barras Horizontales: Nivel de Asistencia

```

1 # Crear categorias de asistencia
2 datos$Nivel_Asistencia <- cut(datos$Asistencia,
3   breaks = c(0, 60, 75, 90, 100),
4   labels = c("Bajo", "Regular", "Bueno", "Excelente"),
5   include.lowest = TRUE
6 )
7
8 # Contar frecuencias
9 freq_asistencia <- table(datos$Nivel_Asistencia)
10
11 # Grafico de barras horizontales
12 barplot(freq_asistencia,
13           horiz = TRUE,
14           col = c("#FF6B6B", "#FFA94D", "#6BCB77", "#4D96FF"),
15           main = "Estudiantes por Nivel de Asistencia",
16           xlab = "Numero de Estudiantes",
17           ylab = "Nivel de Asistencia",
18           las = 1)

```

Listing 12: Crear categorías de asistencia y gráfico de barras

### 支线任务

Se debe observar el gráfico de barras y responder:

1. ¿Cuántos estudiantes tienen asistencia “Excelente”?
2. ¿Cuál nivel de asistencia tiene menos estudiantes?

### 支线任务

1. **Estudiantes con asistencia Excelente:** \_\_\_\_\_
2. **Nivel con menos estudiantes:** \_\_\_\_\_

## 7.3 Diagrama de Cajas: Exámenes Internos 1 y 2

### ◆ Concepto Teórico

El **diagrama de cajas** (boxplot) muestra:

- La **mediana** (línea central de la caja)
- Los **cuartiles Q1 y Q3** (bordes de la caja)
- Los **valores mínimo y máximo** (bigotes)
- Los **valores atípicos** (puntos fuera de los bigotes)

```

1 # Crear diagrama de cajas para ambos exámenes internos
2 boxplot(datos$Examen_Interno_1, datos$Examen_Interno_2,
3           names = c("Examen Interno 1", "Examen Interno 2"),
4           col = c("#4D96FF", "#6BCB77"),
5           main = "Comparacion de Examenes Internos",
6           ylab = "Calificacion (0-100)",
7           border = "darkblue")
8
9 # Agregar linea de la media
10 abline(h = mean(datos$Examen_Interno_1), col = "red", lty = 2)

```

Listing 13: Diagrama de cajas comparando Examen Interno 1 y 2

```

1 # Estadisticas del Examen Interno 1
2 cat("==== EXAMEN INTERNO 1 ====\n")
3 cat("Media:", round(mean(datos$Examen_Interno_1), 2), "\n")
4 cat("Mediana:", round(median(datos$Examen_Interno_1), 2), "\n")
5 cat("Desv. Est.:", round(sd(datos$Examen_Interno_1), 2), "\n\n")
6
7 # Estadisticas del Examen Interno 2
8 cat("==== EXAMEN INTERNO 2 ====\n")
9 cat("Media:", round(mean(datos$Examen_Interno_2), 2), "\n")
10 cat("Mediana:", round(median(datos$Examen_Interno_2), 2), "\n")
11 cat("Desv. Est.:", round(sd(datos$Examen_Interno_2), 2), "\n")

```

Listing 14: Estadísticas de ambos exámenes para comparar

### 📎 Actividad

Se deben comparar los dos exámenes internos usando el boxplot y las estadísticas:

1. ¿Cuál examen tiene la mediana más alta?
2. ¿Cuál examen tiene mayor dispersión (desviación estándar)?
3. ¿Se observan valores atípicos en alguno de los exámenes?

✍ Espacio para Respuesta

1. Examen con mediana más alta: \_\_\_\_\_
2. Examen con mayor dispersión: \_\_\_\_\_
3. ¿Existen valores atípicos? \_\_\_\_\_

## 7.4 Gráfico de Violín: Nota Final por Nivel de Asistencia

### ◆ Concepto Teórico

El **gráfico de violín** combina un boxplot con una estimación de densidad. Muestra:

- La **forma de la distribución** (ancho del violín)
- La **mediana y cuartiles** (boxplot interior)
- La **densidad de los datos** en diferentes valores

Es útil para comparar distribuciones entre grupos.

```

1 # Crear grafico de violin
2 ggplot(datos, aes(x = Nivel_Asistencia, y = Nota_Final,
3                     fill = Nivel_Asistencia)) +
4   geom_violin(trim = FALSE, alpha = 0.7) +
5   geom_boxplot(width = 0.15, fill = "white") +
6   scale_fill_manual(values = c("#FF6B6B", "#FFA94D",
7                      "#6BCB77", "#4D96FF")) +
8   labs(
9     title = "Distribucion de Nota Final por Nivel de Asistencia",
10    x = "Nivel de Asistencia",
11    y = "Nota del Examen Final"
12  ) +
13  theme_minimal() +
14  theme(legend.position = "none")

```

Listing 15: Gráfico de violín: Nota Final por Nivel de Asistencia

### ✎ Actividad

Se debe observar el gráfico de violín y responder:

1. ¿Cuál nivel de asistencia muestra las notas más altas?
2. ¿Se observa alguna diferencia en la forma de las distribuciones?

### ✎ Espacio para Respuesta

1. Nivel con notas más altas: \_\_\_\_\_

2. Diferencias en las distribuciones: \_\_\_\_\_

## 8 Análisis de Normalidad: Nota del Examen Final

### ◆ Concepto Teórico

La **distribución normal** tiene forma de campana simétrica. Muchas pruebas estadísticas (como la correlación de Pearson) asumen que los datos son normales.

#### ¿Cómo se evalúa la normalidad?

1. **Histograma:** Verificar si tiene forma de campana
2. **Q-Q Plot:** Los puntos deben seguir la línea diagonal
3. **Prueba de Shapiro-Wilk:** Si  $p > 0.05$ , los datos son normales
4. **Prueba de Kolmogorov-Smirnov:** Si  $p > 0.05$ , los datos son normales

### 8.1 Histograma de la Nota Final

```

1 # Crear histograma
2 hist(datos$Nota_Final ,
3       breaks = 20 ,
4       col = "#4D96FF" ,
5       border = "white" ,
6       main = "Histograma de Nota Final" ,
7       xlab = "Nota del Examen Final" ,
8       ylab = "Frecuencia" ,
9       freq = FALSE) # Usar densidad en lugar de frecuencia
10
11 # Agregar curva normal teorica
12 curve(dnorm(x, mean = mean(datos$Nota_Final), sd =
13             sd(datos$Nota_Final)),
14         add = TRUE ,
15         col = "red" ,
16         lwd = 2)
17
18 # Agregar leyenda
19 legend("topright",
20        legend = "Curva Normal" ,
21        col = "red" ,
22        lwd = 2)
```

Listing 16: Histograma con curva normal superpuesta

### ⌚ Actividad

Se debe observar el histograma y responder:

1. ¿El histograma tiene forma de campana?
2. ¿La distribución parece simétrica o sesgada?

✍ Espacio para Respuesta

1. ¿Tiene forma de campana? (Sí / No / Aproximadamente): \_\_\_\_\_

2. ¿Es simétrica o sesgada? \_\_\_\_\_

## 8.2 Gráfico Q-Q (Quantile-Quantile)

```

1 # Crear grafico Q-Q
2 qqnorm(datos$Nota_Final,
3         main = "Grafico Q-Q: Nota Final",
4         xlab = "Cuantiles Teoricos",
5         ylab = "Cuantiles de los Datos",
6         col = "#4D96FF",
7         pch = 19,
8         cex = 0.6)
9
10 # Agregar linea de referencia
11 qqline(datos$Nota_Final, col = "red", lwd = 2)

```

Listing 17: Gráfico Q-Q para evaluar normalidad

### ★ Nota Importante

En un Q-Q Plot:

- Si los puntos siguen la línea roja → los datos son normales
- Si los puntos se desvían en los extremos → existen colas pesadas o livianas
- Si los puntos forman una curva → los datos son asimétricos

### 📎 Actividad

Se debe observar el gráfico Q-Q y describir lo observado:

### ✍ Espacio para Respuesta

¿Los puntos siguen la línea roja? \_\_\_\_\_

¿Existen desviaciones en los extremos? \_\_\_\_\_

## 8.3 Pruebas Estadísticas de Normalidad

```

1 # Prueba de Shapiro-Wilk
2 # NOTA: Esta prueba funciona mejor con menos de 5000 datos
3 shapiro_test <- shapiro.test(datos$Nota_Final)
4
5 cat("==== PRUEBA DE SHAPIRO-WILK ===\n")
6 cat("Estadistico W:", round(shapiro_test$statistic, 4), "\n")
7 cat("Valor p:", shapiro_test$p.value, "\n")
8
9 # Interpretar resultado
10 if(shapiro_test$p.value > 0.05) {
11   cat("Conclusion: Los datos SON NORMALES (p > 0.05)\n")
12 } else {
13   cat("Conclusion: Los datos NO SON NORMALES (p <= 0.05)\n")

```

14 }

Listing 18: Prueba de Shapiro-Wilk

```

1 # Prueba de Kolmogorov-Smirnov
2 # Primero se estandarizan los datos
3 datos_std <- (datos$Nota_Final - mean(datos$Nota_Final)) /
  sd(datos$Nota_Final)
4 ks_test <- ks.test(datos_std, "pnorm")
5
6 cat("==== PRUEBA DE KOLMOGOROV-SMIRNOV ===\\n")
7 cat("Estadistico D:", round(ks_test$statistic, 4), "\\n")
8 cat("Valor p:", ks_test$p.value, "\\n")
9
10 # Interpretar resultado
11 if(ks_test$p.value > 0.05) {
12   cat("Conclusion: Los datos SON NORMALES (p > 0.05)\\n")
13 } else {
14   cat("Conclusion: Los datos NO SON NORMALES (p <= 0.05)\\n")
15 }

```

Listing 19: Prueba de Kolmogorov-Smirnov

**📎 Actividad**

Se deben ejecutar las pruebas de normalidad y completar:

**✍ Espacio para Respuesta****Prueba de Shapiro-Wilk:**

Valor p = \_\_\_\_\_ Conclusión: \_\_\_\_\_

**Prueba de Kolmogorov-Smirnov:**

Valor p = \_\_\_\_\_ Conclusión: \_\_\_\_\_

**¿Los datos de Nota\_Final son normales?** \_\_\_\_\_

## 9 Visualizaciones Avanzadas

### 9.1 Mapa de Calor (Heatmap) de Correlaciones

#### ♦ Concepto Teórico

El **mapa de calor** permite visualizar la matriz de correlaciones entre todas las variables numéricas. Los colores indican la fuerza y dirección de la correlación:

- **Colores cálidos (rojos):** Correlación positiva
- **Colores fríos (azules):** Correlación negativa
- **Colores neutros (blancos):** Correlación cercana a cero

```

1 # Seleccionar solo las variables numericas (sin ID)
2 vars_numericas <- datos[, c("Asistencia", "Examen_Interno_1",
3                               "Examen_Interno_2", "Promedio_Tareas",
4                               "Horas_Estudio", "Nota_Final")]
5
6 # Calcular la matriz de correlacion
7 matriz_cor <- cor(vars_numericas)
8
9 # Ver la matriz de correlacion
10 round(matriz_cor, 2)
```

Listing 20: Crear matriz de correlación

```

1 # Crear mapa de calor basico
2 heatmap(matriz_cor,
3          col = colorRampPalette(c("#4D96FF", "white", "#FF6B6B"))(20),
4          symm = TRUE,
5          main = "Mapa de Calor - Correlaciones")
```

Listing 21: Crear mapa de calor con heatmap básico

```

1 # Instalar si es necesario: install.packages("reshape2")
2 library(reshape2)
3
4 # Convertir matriz a formato largo
5 matriz_largo <- melt(matriz_cor)
6 colnames(matriz_largo) <- c("Variable1", "Variable2", "Correlacion")
7
8 # Crear mapa de calor con ggplot2
9 ggplot(matriz_largo, aes(x = Variable1, y = Variable2, fill =
10   Correlacion)) +
11   geom_tile(color = "white") +
12   geom_text(aes(label = round(Correlacion, 2)), size = 3) +
13   scale_fill_gradient2(low = "#4D96FF", mid = "white", high = "#FF6B6B",
14                         midpoint = 0, limit = c(-1, 1)) +
15   labs(title = "Mapa de Calor - Matriz de Correlaciones",
16        x = "", y = "") +
17   theme_minimal() +
18   theme(axis.text.x = element_text(angle = 45, hjust = 1))
```

---

Listing 22: Mapa de calor mejorado con ggplot2

### 💡 Actividad

Se debe observar el mapa de calor y responder:

1. ¿Cuáles dos variables tienen la correlación más fuerte (más cercana a 1)?
2. ¿Existe alguna correlación negativa entre las variables?
3. ¿Cuál variable tiene la mayor correlación con Nota\_Final?

### ✍ Espacio para Respuesta

1. **Variables con correlación más fuerte:** \_\_\_\_\_
2. **¿Existe correlación negativa?** \_\_\_\_\_
3. **Variable más correlacionada con Nota\_Final:** \_\_\_\_\_

## 9.2 Gráfico de Violín

### ♦ Concepto Teórico

El **gráfico de violín** combina un diagrama de cajas con una estimación de densidad. Permite visualizar:

- La **distribución completa** de los datos (forma del violín)
- La **mediana y cuartiles** (caja interior)
- La **concentración** de datos en diferentes valores

Es útil para comparar distribuciones entre grupos.

```

1 # Crear grafico de violin
2 ggplot(datos, aes(x = Nivel_Asistencia, y = Nota_Final,
3                     fill = Nivel_Asistencia)) +
4   geom_violin(trim = FALSE, alpha = 0.7) +
5   geom_boxplot(width = 0.15, fill = "white") +
6   scale_fill_manual(values = c("#FF6B6B", "#FFA94D", "#6BCB77",
7                      "#4D96FF")) +
8   labs(title = "Distribucion de Nota Final por Nivel de Asistencia",
9        x = "Nivel de Asistencia",
10       y = "Nota del Examen Final") +
11  theme_minimal() +
12  theme(legend.position = "none")

```

Listing 23: Gráfico de violín: Nota Final por Nivel de Asistencia

```

1 # Grafico de violin por categoria de rendimiento
2 ggplot(datos, aes(x = Rendimiento, y = Examen_Interno_1,
3                     fill = Rendimiento)) +
4   geom_violin(trim = FALSE, alpha = 0.7) +
5   geom_boxplot(width = 0.1, fill = "white") +
6   scale_fill_manual(values = c("#FF6B6B", "#FFA94D", "#FFD93D",
7                      "#6BCB77", "#4D96FF")) +
8   labs(title = "Distribucion de Examen Interno 1 por Rendimiento Final",
9        x = "Categoria de Rendimiento",
10       y = "Nota Examen Interno 1") +
11  theme_minimal() +
12  theme(legend.position = "none",
13        axis.text.x = element_text(angle = 45, hjust = 1))

```

Listing 24: Gráfico de violín: Nota Final por Rendimiento

### 💡 Actividad

Se debe observar el gráfico de violín de Nota Final por Nivel de Asistencia y responder:

1. ¿Cuál nivel de asistencia muestra la distribución más concentrada (violín más angosto)?
2. ¿En cuál nivel de asistencia se observan las notas más altas?
3. ¿Qué se puede concluir sobre la relación entre asistencia y rendimiento?

### ✍ Espacio para Respuesta

1. Nivel con distribución más concentrada: \_\_\_\_\_

2. Nivel con notas más altas: \_\_\_\_\_

3. Conclusión sobre asistencia y rendimiento: \_\_\_\_\_  
\_\_\_\_\_

## 10 Análisis de Correlación

Primero se analizará la correlación entre todas las variables numéricas mediante un mapa de calor, y luego se profundizará en la relación entre **Asistencia** y **Nota Final**.

### ◆ Concepto Teórico

El **coeficiente de correlación de Pearson** ( $r$ ) mide la fuerza y dirección de la relación lineal entre dos variables.

#### Interpretación:

- $r = 1$ : Correlación positiva perfecta
- $0.7 \leq r < 1$ : Correlación positiva fuerte
- $0.4 \leq r < 0.7$ : Correlación positiva moderada
- $0.2 \leq r < 0.4$ : Correlación positiva débil
- $-0.2 < r < 0.2$ : Correlación muy débil o nula
- $r = -1$ : Correlación negativa perfecta

### 10.1 Mapa de Calor (Heatmap) de Correlaciones

```

1 # Seleccionar solo variables numericas
2 vars_numericas <- datos[, c("Asistencia", "Examen_Interno_1",
3                               "Examen_Interno_2", "Promedio_Tareas",
4                               "Horas_Estudio", "Nota_Final")]
5
6 # Calcular matriz de correlacion
7 matriz_cor <- cor(vars_numericas)
8
9 # Mostrar la matriz
10 round(matriz_cor, 2)

```

Listing 25: Calcular matriz de correlación

```

1 # Mapa de calor con corrplot
2 corrplot(matriz_cor,
3           method = "color",
4           type = "upper",
5           addCoef.col = "black",
6           number.cex = 0.7,
7           tl.col = "black",
8           tl.srt = 45,
9           col = colorRampPalette(c("#FF6B6B", "white", "#4D96FF"))(100),
10          title = "Mapa de Calor - Correlaciones",
11          mar = c(0, 0, 2, 0))

```

Listing 26: Crear mapa de calor con corrplot

```

1 # Convertir matriz a formato largo
2 matriz_largo <- melt(matriz_cor)

```

```

3 colnames(matriz_largo) <- c("Variable1", "Variable2", "Correlacion")
4
5 # Crear heatmap con ggplot2
6 ggplot(matriz_largo, aes(x = Variable1, y = Variable2,
7                           fill = Correlacion)) +
8   geom_tile(color = "white") +
9   geom_text(aes(label = round(Correlacion, 2)), size = 3) +
10  scale_fill_gradient2(low = "#FF6B6B", mid = "white",
11                        high = "#4D96FF", midpoint = 0) +
12  labs(title = "Matriz de Correlacion") +
13  theme_minimal() +
14  theme(axis.text.x = element_text(angle = 45, hjust = 1))

```

Listing 27: Mapa de calor alternativo con ggplot2

### 📎 Actividad

Se debe observar el mapa de calor y responder:

1. ¿Cuáles dos variables tienen la correlación más alta?
2. ¿Existe alguna correlación negativa?
3. ¿Qué variable tiene la correlación más fuerte con Nota\_Final?

### ✍ Espacio para Respuesta

1. **Variables con correlación más alta:** \_\_\_\_\_
2. **¿Existe correlación negativa?** \_\_\_\_\_
3. **Variable más correlacionada con Nota\_Final:** \_\_\_\_\_

## 10.2 Gráfico de Dispersión: Asistencia vs Nota Final

```

1 # Crear grafico de dispersion
2 plot(datos$Asistencia, datos$Nota_Final,
3       main = "Relacion entre Asistencia y Nota Final",
4       xlab = "Porcentaje de Asistencia",
5       ylab = "Nota del Examen Final",
6       col = "#4D96FF",
7       pch = 19,
8       cex = 0.7)
9
10 # Agregar linea de tendencia
11 abline(lm(Nota_Final ~ Asistencia, data = datos),
12         col = "red",
13         lwd = 2)

```

Listing 28: Gráfico de dispersión: Asistencia vs Nota Final

## 10.3 Cálculo del Coeficiente de Correlación

```

1 # Calcular correlacion de Pearson
2 correlacion <- cor(datos$Asistencia, datos$Nota_Final)
3 cat("Coeficiente de correlacion (r):", round(correlacion, 4), "\n")
4
5 # Prueba de significancia
6 prueba_cor <- cor.test(datos$Asistencia, datos$Nota_Final)
7 cat("Valor p:", prueba_cor$p.value, "\n")
8
9 # Interpretar
10 if(prueba_cor$p.value < 0.05) {
11   cat("La correlacion ES estadisticamente significativa\n")
12 } else {
13   cat("La correlacion NO es significativa\n")
14 }

```

Listing 29: Calcular correlación de Pearson

### Actividad

Se debe calcular la correlación entre Asistencia y Nota Final:

✍ Espacio para Respuesta

**Coeficiente r =** \_\_\_\_\_

**Valor p =** \_\_\_\_\_

**Tipo de correlación** (fuerte/moderada/débil): \_\_\_\_\_

**¿Es significativa? (Sí / No):** \_\_\_\_\_

## 11 Regresión Lineal Simple

Se creará un modelo para predecir la **Nota Final** basándose en la **Asistencia**.

### ♦ Concepto Teórico

La **regresión lineal simple** modela la relación entre dos variables con la ecuación:

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

Donde:

- $Y$ : Variable dependiente (Nota Final)
- $X$ : Variable independiente (Asistencia)
- $\beta_0$ : Intercepto (valor de  $Y$  cuando  $X = 0$ )
- $\beta_1$ : Pendiente (cambio en  $Y$  por cada unidad de  $X$ )
- $\varepsilon$ : Error aleatorio

El **coeficiente de determinación ( $R^2$ )** indica qué porcentaje de la variabilidad de  $Y$  es explicado por  $X$ .

### 11.1 Crear el Modelo de Regresión

```

1 # Crear el modelo de regresion
2 modelo <- lm(Nota_Final ~ Asistencia, data = datos)
3
4 # Ver resumen del modelo
5 summary(modelo)

```

Listing 30: Crear modelo de regresión lineal simple

### 11.2 Interpretar los Resultados

```

1 # Obtener coeficientes
2 intercepto <- coef(modelo)[1]
3 pendiente <- coef(modelo)[2]
4
5 cat("==== ECUACION DEL MODELO ====\n")
6 cat("Nota_Final =", round(intercepto, 2), "+",
7     round(pendiente, 2), "* Asistencia\n\n")
8
9 # Coeficiente de determinacion R^2
10 r_cuadrado <- summary(modelo)$r.squared
11 cat("R cuadrado (R^2):", round(r_cuadrado, 4), "\n")
12 cat("Esto significa que la Asistencia explica el",
13     round(r_cuadrado * 100, 1), "% de la variabilidad en la Nota
Final\n")

```

Listing 31: Extraer e interpretar coeficientes

**💡 Actividad**

Se debe crear el modelo de regresión y completar:

**✍ Espacio para Respuesta**

**Ecuación del modelo:**

Nota\_Final = \_\_\_\_\_ + \_\_\_\_\_  $\times$  Asistencia

**R cuadrado ( $R^2$ ) = \_\_\_\_\_**

**Interpretación:** La Asistencia explica el \_\_\_\_\_ % de la variabilidad en la nota final.

### 11.3 Visualizar el Modelo de Regresión

```

1 # Grafico con ggplot2
2 ggplot(datos, aes(x = Asistencia, y = Nota_Final)) +
3   geom_point(color = "#4D96FF", alpha = 0.5) +
4   geom_smooth(method = "lm", color = "red", se = TRUE) +
5   labs(
6     title = "Regresion Lineal: Nota Final vs Asistencia",
7     subtitle = paste("R^2 =", round(r_cuadrado, 4)),
8     x = "Porcentaje de Asistencia",
9     y = "Nota del Examen Final"
10   ) +
11   theme_minimal()

```

Listing 32: Gráfico del modelo de regresión

### 11.4 Realizar Predicciones

```

1 # Crear datos para predicción
2 nuevos_datos <- data.frame(
3   Asistencia = c(60, 70, 80, 90, 100)
4 )
5
6 # Realizar predicciones
7 predicciones <- predict(modelo, newdata = nuevos_datos)
8
9 # Mostrar resultados
10 cat("==== PREDICCIONES ===\n")
11 for(i in 1:5) {
12   cat("Si la asistencia es", nuevos_datos$Asistencia[i],
13       "% -> Nota predicha:",
14       round(predicciones[i], 1), "\n")
15 }

```

Listing 33: Predecir notas para diferentes niveles de asistencia

#### ↳ Actividad

Se debe usar el modelo para predecir la nota de un estudiante que tiene **85% de asistencia**:

```

1 # Predecir para 85% de asistencia
2 mi_prediccion <- data.frame(Asistencia = 85)
3 resultado <- predict(modelo, newdata = mi_prediccion)
4 cat("Nota predicha para 85% de asistencia:", round(resultado, 1), "\n")

```

Listing 34: Predicción personalizada

#### ↳ Espacio para Respuesta

**Si un estudiante tiene 85% de asistencia:**

Nota Final predicha = \_\_\_\_\_

## 12 Actividades de Cierre

### Actividad

#### Reflexión Final

Con base en todo el análisis realizado, se deben responder las siguientes preguntas:

1. ¿Cuál es el promedio de la Nota Final de todos los estudiantes?
2. ¿Los datos de Nota Final siguen una distribución normal? Se debe justificar la respuesta.
3. ¿Existe una relación significativa entre la Asistencia y la Nota Final?
4. Según el modelo, ¿cuánto aumenta la nota final por cada punto porcentual adicional de asistencia?
5. ¿Qué recomendaciones se le darían a un estudiante que desea mejorar su nota final?

### Espacio para Respuesta

1. Promedio de Nota Final: \_\_\_\_\_

2. ¿Son normales los datos? \_\_\_\_\_

Justificación: \_\_\_\_\_  
\_\_\_\_\_

3. ¿Existe relación significativa? \_\_\_\_\_

4. Aumento por punto de asistencia: \_\_\_\_\_

5. Recomendaciones: \_\_\_\_\_  
\_\_\_\_\_

## 13 Resumen del Taller

### Funciones de R aprendidas

1. Importación de datos: `read.csv()`
2. Exploración: `head()`, `str()`, `summary()`
3. Estadística descriptiva: `mean()`, `median()`, `sd()`, `var()`, `quantile()`
4. Visualización: `pie()`, `barplot()`, `boxplot()`, `hist()`
5. Normalidad: `qqnorm()`, `qqline()`, `shapiro.test()`, `ks.test()`
6. Correlación: `cor()`, `cor.test()`
7. Regresión: `lm()`, `predict()`

## 14 Referencias y Recursos

### 14.1 Dataset

- Fuente: Kaggle - Student Performance Prediction Dataset
- URL: <https://www.kaggle.com/>

### 14.2 Recursos de R

- R Documentation: <https://www.rdocumentation.org/>
- R for Data Science: <https://r4ds.had.co.nz/>
- ggplot2: <https://ggplot2.tidyverse.org/>