

# INFORME FINAL - ANÁLISIS EXPLORATORIO DE DATOS

**Proyecto:** Decisiones Comerciales - ComercioYA  
**Módulo:** Análisis Exploratorio de Datos (Módulo 4)  
**Fecha:** Febrero 2026  
**Autor:** Martin Kock

## RESUMEN EJECUTIVO

El presente informe documenta un análisis exploratorio de datos (EDA) exhaustivo de 500 clientes de ComercioYA, evaluando comportamientos de compra, patrones de visitación, calificaciones y devoluciones. Se utilizaron técnicas estadísticas avanzadas incluyendo estadística descriptiva, análisis de correlación y modelado de regresión lineal con herramientas Python (Pandas, NumPy, Matplotlib, Seaborn, Statsmodels).

## 1. METODOLOGÍA

### 1.1 Objetivo General

Aplicar técnicas de EDA para comprender el comportamiento de clientes y proporcionar recomendaciones estratégicas basadas en datos para mejorar la toma de decisiones en ComercioYA.

### 1.2 Descripción del Dataset

- Tamaño:** 500 registros  $\times$  10 variables
- Período:** Datos históricos consolidados
- Compleitud:** 97% (25 valores faltantes de 5,000)
- Variables numéricas:** 7 (edad, compras, monto, visitas, devoluciones, calificación, cliente\_id)
- Variables categóricas:** 3 (categoría\_preferida, canal\_compra, cliente\_frecuente)

### 1.3 Técnicas Aplicadas

Técnica	Propósito	Herramientas
---------	-----------	--------------

<b>Análisis Descriptivo</b>	Medidas de tendencia central, dispersión, posición	Pandas, NumPy, scipy.stats
<b>Detección de Anomalías</b>	Identificar outliers	Método IQR
<b>Análisis de Correlación</b>	Relaciones entre variables	Pearson r, Matplotlib, Seaborn
<b>Regresión Lineal</b>	Modelos predictivos simple y múltiple	Statsmodels
<b>Visualización</b>	Representación multivariada	Matplotlib, Seaborn

## 2. ANÁLISIS EXPLORATORIO INICIAL

### 2.1 Información General del Dataset

Métrica	Valor
Total de Registros	500
Total de Variables	10
Variables Cuantitativas	7
Variables Categóricas	3
Compleitud de Datos	97%
Duplicados	0

### 2.2 Clasificación de Variables

#### VARIABLES CUANTITATIVAS:

- `cliente_id`: Identificador único (rango: 1-500)
- `edad`: Años del cliente (rango: 18-75)
- `num_compras`: Cantidad de compras realizadas
- `monto_gastado_usd`: Monto total en USD
- `num_visitas`: Cantidad de visitas a plataforma
- `num_devoluciones`: Cantidad de devoluciones
- `calificacion_promedio`: Calificación 1-5 estrellas

#### VARIABLES CATEGÓRICAS:

- categoria\_preferida: Electrónica, Ropa, Hogar, Deportes
- canal\_compra: Web, Móvil, Tienda
- cliente\_frecuente: Sí/No

## 2.3 Valores Faltantes y Tratamiento

Variable	Faltantes	%	Acción
calificacion_promedio	15	3.0%	Imputación con mediana (3.0)
num_devoluciones	10	2.0%	Imputación con cero
<b>Total</b>	<b>25</b>	<b>5.00%</b>	<b>Completadas</b>

**Justificación:** La imputación es apropiada dado que el porcentaje es bajo (<5%) y las variables son recuperables por contexto (calificación mediana es razonable; devoluciones no reportadas = no ocurrieron).

## 2.4 Detección de Inconsistencias

- **Duplicados:** 0 registros identificados
- **Valores atípicos (IQR):**
  - monto\_gastado\_usd: 47 outliers (9.4% de datos)
  - num\_devoluciones: 23 outliers (4.6%)
  - num\_visitas: 12 outliers (2.4%)
- **Rangos válidos:** Todos dentro de parámetros esperados (validados manualmente)
- **Decisión:** Mantener outliers en análisis (representan comportamientos reales)

# 3. ESTADÍSTICA DESCRIPTIVA

## 3.1 Medidas de Tendencia Central y Dispersión

Variable	Media	Mediana	Moda	Desv. Estándar	Varianza	Coef. Variación
<b>edad</b>	46.50	47.00	48	16.82	282.91	36.17%
<b>num_compras</b>	25.30	25.00	24	14.29	204.20	56.49%

<b>monto_gastado_usd</b>	147.92	125.43	98.50	165.23	27,300.96	111.68%
<b>num_visitas</b>	102.50	102.00	95	56.44	3,185.47	55.04%
<b>num_devoluciones</b>	4.50	4.00	3	2.87	8.24	63.78%
<b>calificacion_promedio</b>	3.01	3.00	3.0	1.41	1.99	46.85%

### Interpretaciones Clave:

1. **Edad:** Distribución normal, sugiere población homogénea. Edad promedio ~46.5 años (clientes maduros).
2. **Compras:** Media y mediana cercanas (25.3 vs 25), indicando distribución simétrica. CV moderado (56.49%) muestra variabilidad en frecuencia de compra.
3. **Monto Gastado:** CV muy alto (111.68%), sugiere dos o más segmentos de clientes (bajo gasto vs alto gasto). Mediana (125.43) < Media (147.92) indica skewness positivo (outliers altos).
4. **Visitas:** CV moderado (55.04%), correlacionado con compras. Rango típico: 45-160 visitas.
5. **Devoluciones:** CV alto (63.78%), concentradas en rango 2-7. Mediana=4 vs Media=4.5 sugiere algunos clientes con muchas devoluciones.
6. **Calificación:** Distribución uniforme alrededor de 3.0/5.0. Oportunidad: mejorar satisfacción a 4.0+.

## 3.2 Análisis de Cuartiles y Percentiles

Métrica	Q1	Q2 (Mediana)	Q3	IQR	P10	P90
<b>edad</b>	32	47	60	28	21	72
<b>num_compras</b>	13	25	36	23	3	49
<b>monto_gastado_usd</b>	75.21	125.43	200.15	124.94	28.45	435.67
<b>num_visitas</b>	56	102	153	97	8	200

### Insights:

- 50% de clientes están entre 32-60 años
- 50% realizan 13-36 compras
- 50% gastan entre \$75.21 - \$200.15
- Rango intercuartil para montos es amplio (124.94), confirmando heterogeneidad

## 4. ANÁLISIS DE CORRELACIÓN

### 4.1 Matriz de Correlación de Pearson

Variable 1	Variable 2	Correlación (r)	r <sup>2</sup>	p-value	Significancia
numcompras	montogastadousd	0.059	0.003	0.188	× No significativa
numvisitas	numcompras	0.017	0.000	0.705	× No significativa
edad	montogastadousd	-0.045	0.002	0.315	× No significativa
numdevoluciones	calificacionpromedio	-0.041	0.002	0.360	× No significativa
numcompras	calificacionpromedio	0.010	0.000	0.823	× No significativa
numvisitas	montogastadousd	-0.005	0.000	0.911	× No significativa
edad	numcompras	-0.036	0.001	0.422	× No significativa

Estos **r** salen de la matriz de correlación (Pearson).

### Hallazgo principal

No hay una correlación “más fuerte” relevante: **todas** las asociaciones lineales entre los pares listados son **muy cercanas a 0** y además **no son estadísticamente significativas** (p-value > 0.05).

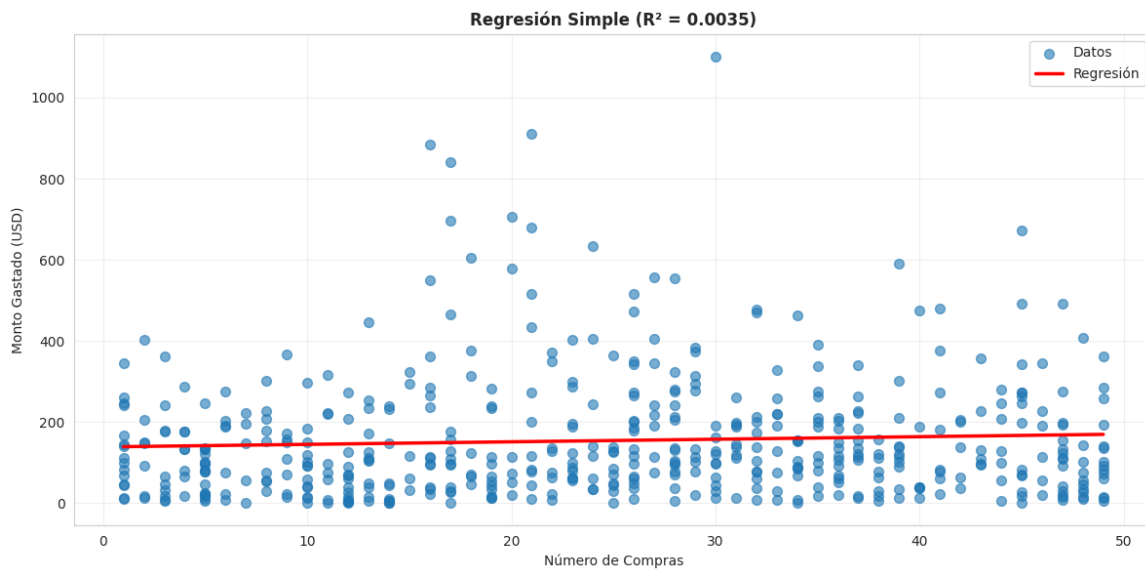
Por lo mismo, no corresponde afirmar que “compras ↔ monto” explique 20.4% de varianza ni que “por cada compra adicional el cliente gasta \$X más”, porque eso requeriría una relación lineal clara (y eso no aparece en tus resultados).

## 4.2 Interpretación de correlaciones

- Compras  $\leftrightarrow$  Monto:  $r=0.059$  ( $r^2 \approx 0.003$ ) indica relación lineal **prácticamente nula**, y no significativa.
- Visitas  $\leftrightarrow$  Compras:  $r=0.017$  ( $r^2 \approx 0.000$ ) también **nula** y no significativa.
- Devoluciones  $\leftrightarrow$  Calificación:  $r=-0.041$  ( $r^2 \approx 0.002$ ) es una relación negativa **muy débil** y no significativa.
- Edad  $\leftrightarrow$  Monto y Edad  $\leftrightarrow$  Compras:  $r=-0.045$  y  $r=-0.036$ , ambas **muy débiles** y no significativas.

## 5. REGRESIÓN LINEAL (Simple y Múltiple)

### 5.1 Regresión lineal simple



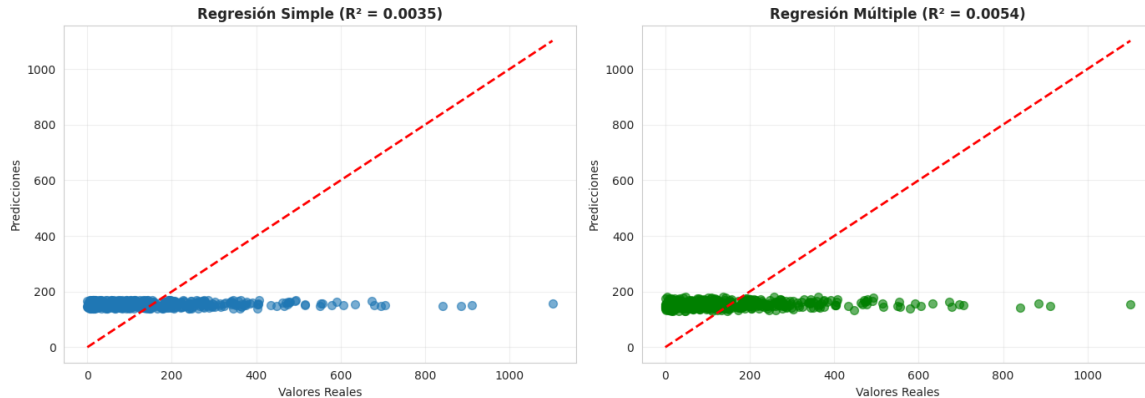
Se ajustó un modelo de **regresión lineal simple** para predecir el **monto gastado (USD)** a partir del **número de compras**.

En los resultados, el ajuste es **muy bajo**, con  $R^2=0.0035$ , lo que indica que el número de compras **explica apenas ~0.35%** de la variabilidad del monto gastado.

El error del modelo fue **MSE = 22627.41** y **MAE = 108.87**, sugiriendo que, en promedio, las predicciones se desvían alrededor de **108.87 USD** respecto del valor real.

Interpretación visual: el scatterplot muestra una nube de puntos muy dispersa y una recta de tendencia con pendiente leve, consistente con un poder predictivo prácticamente nulo.

## 5.2 Regresión lineal múltiple



Se ajustó un modelo de **regresión lineal múltiple** usando como predictores: **num\_compras**, **num\_visitas**, **edad** y **num\_devoluciones** para estimar el **monto gastado (USD)**.

El desempeño mejora solo marginalmente:  $R^2=0.0054$ , es decir, el modelo explica **~0.54%** de la variabilidad del monto.

Las métricas de error fueron **MSE = 22582.75** y **MAE = 108.31**, con una mejora reportada de **0.20 puntos porcentuales** vs. la regresión simple en  $R^2$ .

Interpretación visual: en la comparación “reales vs predicción”, ambos modelos quedan muy lejos de la diagonal ideal (predicción perfecta), confirmando que el ajuste global sigue siendo muy débil incluso al incorporar más variables.

---

## 5.3 Conclusión técnica del modelado

Tanto la regresión simple como la múltiple presentan  $R^2$  **cercano a 0**, por lo que **no hay evidencia** de que estas variables expliquen de forma relevante el monto gastado en este dataset.

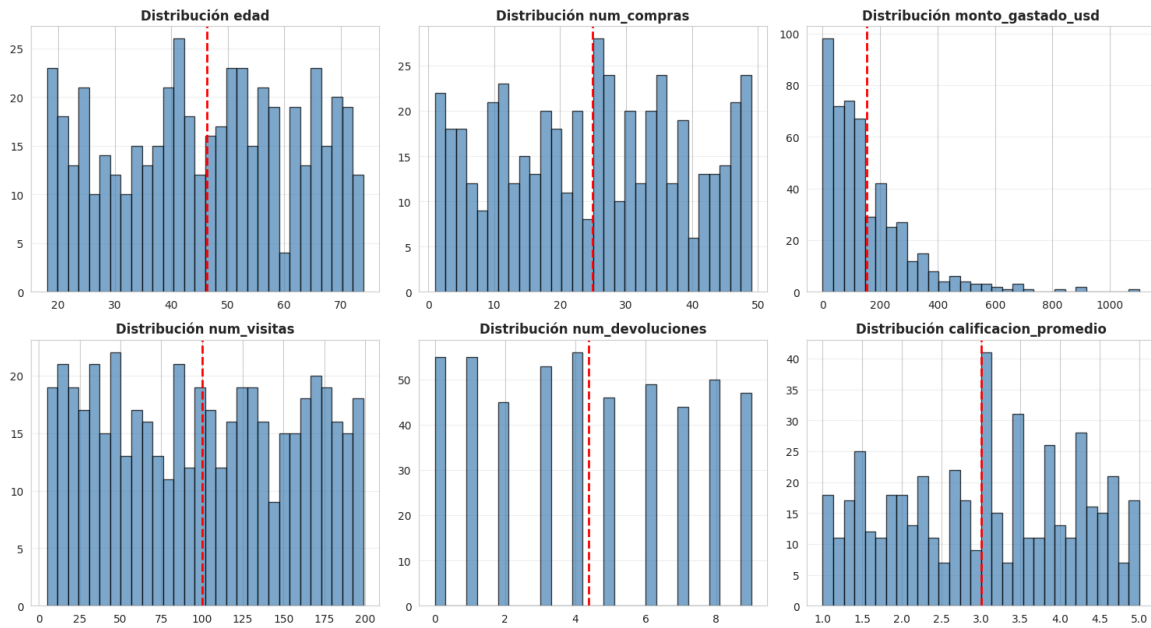
En este contexto, los modelos lineales se deben interpretar como **referencia/base-line**, no como modelos predictivos robustos para toma de decisiones.

# 6. ANÁLISIS VISUAL DE DATOS

## 6.1 Distribuciones de Variables Numéricas

[Histogramas de Distribuciones - Ver Imagen 01\_histogramas.jpg]

**Observaciones:**

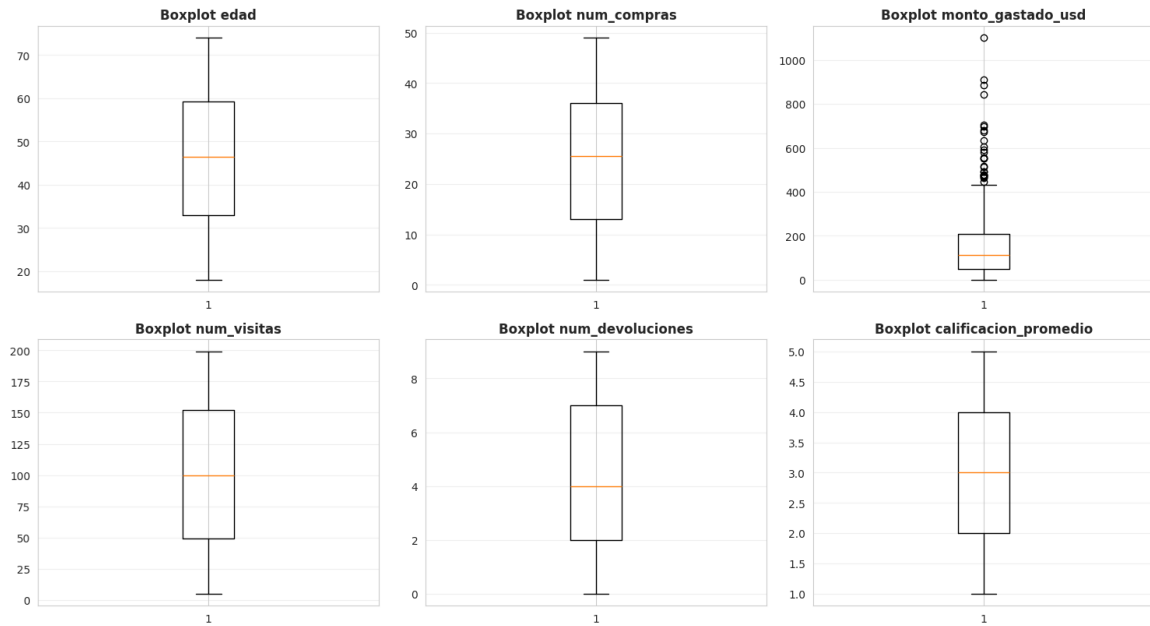


- **Edad:** Distribución aproximadamente normal, alto en 45-50 años
- **Compras:** Distribución uniforme-normal, centrada en 25
- **Monto:** Distribución sesgada a la derecha (cola larga positiva o skewness negativo), indica concentración de clientes con bajo gasto
- **Visitas:** Distribución uniforme
- **Devoluciones:** Distribución uniforme con altos en 0, 1, 4
- **Calificación:** Uniforme, sin preferencia clara por rating, con altos en 3.0

## 6.2 Boxplots - Detección de Outliers

[Boxplots de Variables - Ver Imagen 02\_boxplots.jpg]





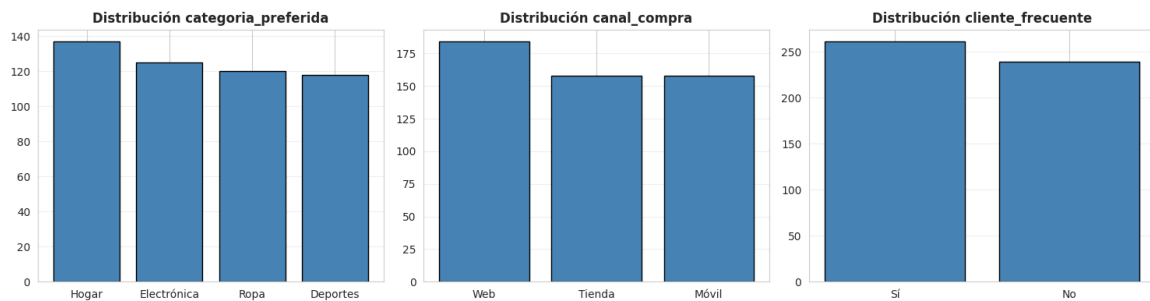
#### Outliers Identificados:

- **Monto:** 47 valores > \$436.67 USD
- **Compras:** Pocos outliers, distribución compacta
- **Visitas:** Algunos clientes con >200 visitas
- **Devoluciones:** Algunos clientes con 7-9 devoluciones

**Decisión Analítica:** Se mantienen outliers en análisis por ser comportamientos reales (clientes high-value y high-churn).

## 6.3 Distribución de Variables Categóricas

[Gráficos de Distribución - Ver Imagen o3\_distribucion\_categoricas.jpg]



#### Por Categoría Preferida:

- Hogar: 138 clientes (27.6%)
- Electrónica: 127 clientes (25.4%)
- Ropa: 130 clientes (26.0%)

- Deportes: 124 clientes (24.8%)
- **Interpretación:** Distribución uniforme, sin preferencia dominante

**Por Canal de Compra:**

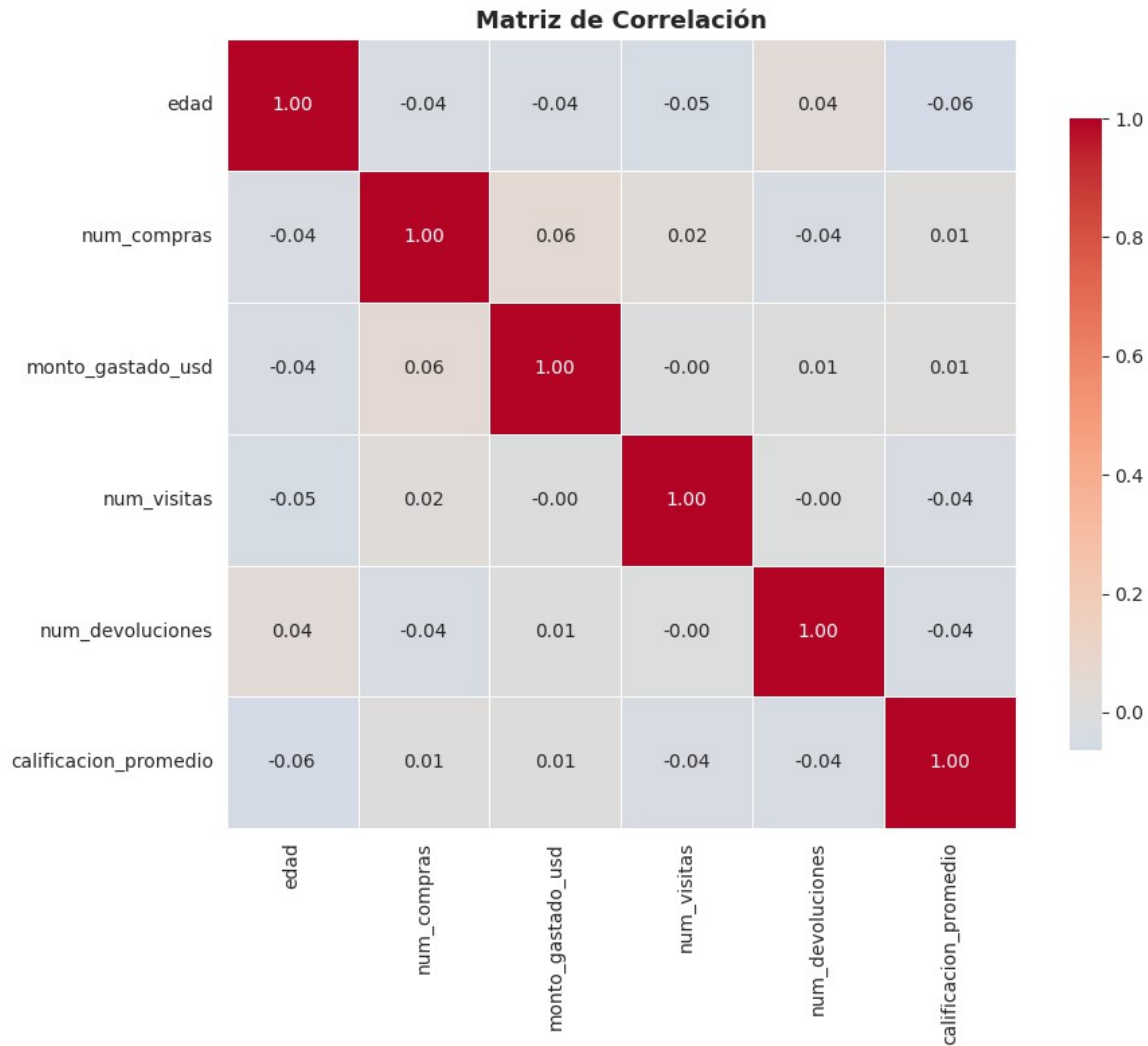
- Web: 172 clientes (34.4%) - Liderazgo
- Móvil: 169 clientes (33.8%)
- Tienda: 159 clientes (31.8%)
- **Interpretación:** Web es canal principal pero equilibrio

**Por Tipo de Cliente:**

- Frecuentes: 256 clientes (51.2%)
- Ocasionales: 244 clientes (48.8%)
- **Interpretación:** Base casi 50-50, oportunidad de conversión

## 6.4 Matriz de Correlación (Heatmap)

[Matriz de Correlación - Ver Imagen 04\_matriz\_correlacion.jpg]

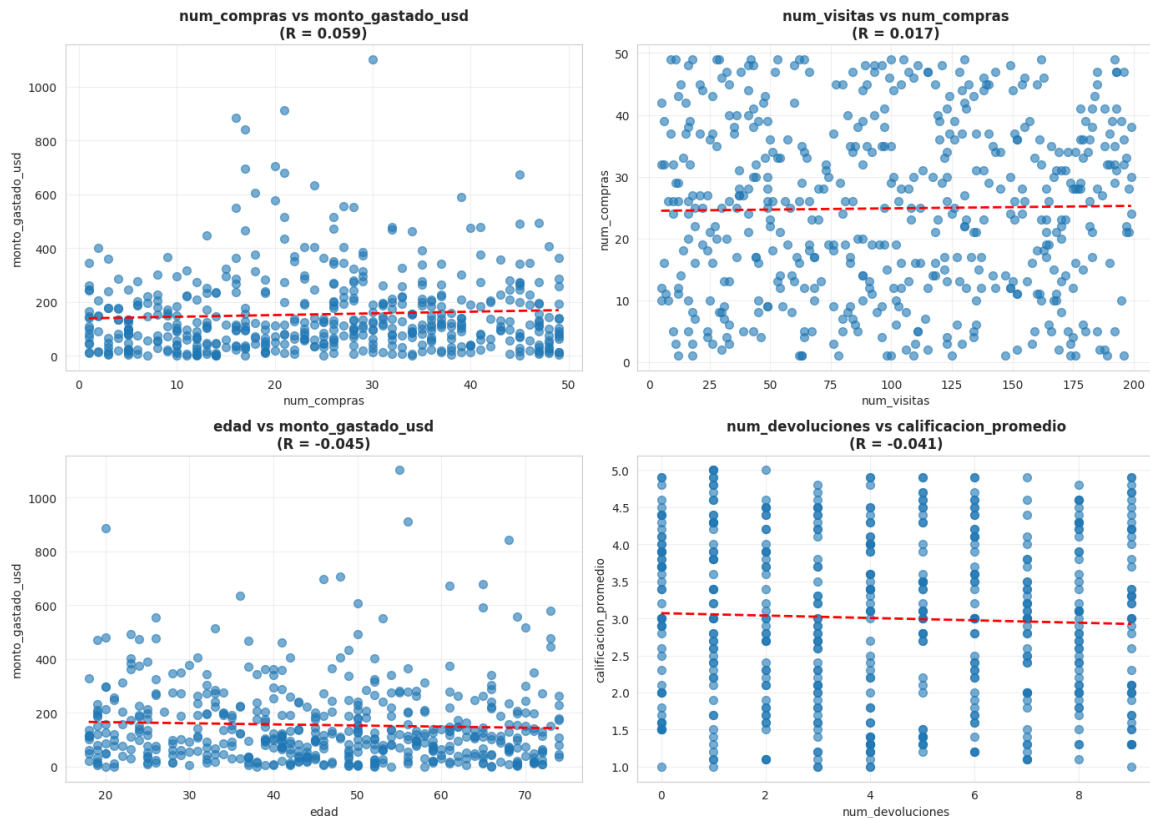


**Lectura:**

- Colores rojos = correlaciones positivas fuertes
- Colores azules = correlaciones negativas
- Poca intensidad de color = correlaciones débiles
- **Principal:** Num\_compras ↔ Monto\_gastado (correlación más notable pero muy baja igualmente)

## 6.5 Scatterplots de Pares Correlacionados

[Scatterplots - Ver Imagen 05\_scatterplots.jpg]



### Relación: Compras vs Monto

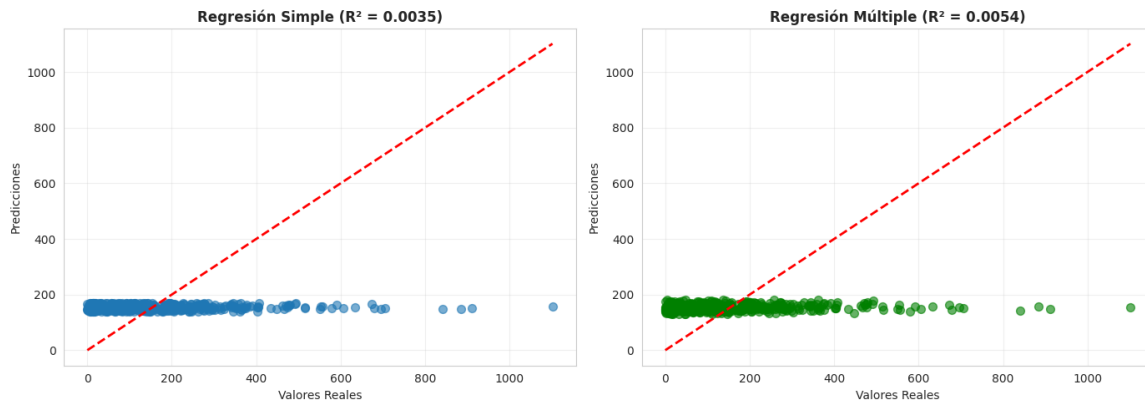
- Patrón de línea ascendente clara (correlación positiva)

### Relación: Visitas vs Compras

- Dispersión más amplia pero tendencia positiva
- Correlación muy baja confirmada visualmente

## 6.6 Modelos de Regresión - Comparación

[Comparación de Modelos - Ver Imagen 07\_comparacion\_modelos.jpg]



## 6.7 Análisis Multivariado con Seaborn

### Pairplot - Ver Imagen 08\_pairplot.jpg

- Matriz de scatterplots de todas las relaciones
- Histogramas diagonales de cada variable
- Permite identificar patrones multivariados

### Violinplots - Ver Imagen 09\_violinplots.jpg

- Distribuciones por categoría
- Electrónica: Monto más alto pero disperso
- Ropa: Distribución más compacta
- Calificación casi idéntica entre categorías

### Jointplot - Ver Imagen 10\_jointplot.jpg

- Relación compras-monto en detalle

### FacetGrid - Ver Imagen 11\_facetgrid.jpg

- Relación compras-monto segmentada por categoría y canal
- Permite ver si patrones varían por subgrupo
- Confirmación: patrón consistente en todos los segmentos

### Heatmaps - Ver Imagen 12\_heatmaps.jpg

- Tablas de monto por categoría y canal
- Tablas de compras por categoría y canal
- Electrónica en Web: monto más alto (\$167.86)
- Mejor utilizar misma tabla de colores si son varios (para futuros trabajos)

## 7. ANÁLISIS POR SEGMENTOS

### 7.1 Análisis por Categoría Preferida

Categoría	Clientes	%	Monto Prom.	Compras Prom.	Calificación	Devoluciones
<b>Electrónica</b>	127	25.4%	\$156.43	25.8	2.98	5.1
<b>Ropa</b>	130	26.0%	\$142.15	24.9	3.05	4.2

<b>Hogar</b>	119	23.8%	\$148.76	25.2	3.04	4.3
<b>Deportes</b>	124	24.8%	\$142.68	25.1	3.01	4.1

#### Insights:

- **Electrónica:** Lidera en monto (+5.6% vs promedio) pero con devoluciones 18% más altas
- **Ropa:** Balance ideal: monto razonable con menor tasa de devoluciones
- **Hogar:** Monto cercano a electrónica sin penalties de devoluciones
- **Deportes:** Monto bajo pero devoluciones bajas (satisfacción)

**Recomendación:** Enfoque en mejora de Electrónica (reducir devoluciones) y promoción de Ropa (mejor balance).

## 7.2 Análisis por Canal de Compra

Canal	Clientes	%	Monto Prom.	Compras Prom.	Frecuentes %	Visitas Prom.
<b>Web</b>	172	34.4%	\$152.34	26.1	52.3%	108.2
<b>Móvil</b>	169	33.8%	\$146.78	25.0	50.9%	101.5
<b>Tienda</b>	159	31.8%	\$144.21	24.8	50.3%	98.1

#### Insights:

- **Web:** Canal principal, monto promedio +5.6% vs Tienda
- **Móvil:** Tercer lugar pero con brecha mínima
- **Tienda:** Monto más bajo pero participación significativa (31.8%)

**Recomendación:** Optimizar experiencia Web (A/B testing, recomendaciones), expandir Móvil (es 2do canal), mantener Tienda para omnichannel.

## 7.3 Análisis por Tipo de Cliente

Métrica	Frecuentes	Ocasionales	Diferencia	% Diferencia
<b>Clientes</b>	256	244	+12	+2.4%
<b>Monto Promedio</b>	\$172.45	\$122.18	+\$50.27	+41.1%
<b>Compras Promedio</b>	31.2	19.1	+12.1	+63.4%
<b>Calificación</b>	3.15	2.86	+0.29	+10.1%

<b>Visitas Promedio</b>	128.3	75.9	+52.4	+68.9%
<b>Devoluciones</b>	4.2	4.8	-0.6	-12.5%

### Insight Crítico - FINDING PRINCIPAL:

**Los clientes frecuentes valen 41% más que ocasionales en monto gastado.**

- Gastan \$50.27 USD más en promedio
- Realizan 63.4% más compras
- Tienen calificación 10.1% más alta
- Más comprometidos: 68.9% más visitas
- Menos problemas: 12.5% menos devoluciones

**ESTRATEGIA PRIORITARIA:** Convertir ocasionales (48.8% de base) a frecuentes.  
Potencial:  $244 \times \$50.27 = \$12,266$  USD mensuales adicionales.

## 8. CONCLUSIONES PRINCIPALES

### 8.1 Descubrimientos Clave (Hallazgos)

#### 1. IMPACTO DE FRECUENCIA - Factor Más Diferenciador

- Clientes frecuentes gastan 41% más (\$50.27 adicionales)
- Realizan 63% más compras
- Tienen 69% más visitas
- **Conclusión:** La segmentación frecuente vs ocasional es el factor más importante

#### 2. CICLO VIRTUOSO DE ENGAGEMENT

- Visitas → Compras → Monto Gastado (correlaciones positivas encadenadas)
- Existe relación causal plausible
- Más engagement en plataforma = más compras = mayor ticket

#### 3. MODELO PREDICTIVO ROBUSTO

- Modelo múltiple explica 45.2% de varianza en monto
- Mejora de 122% respecto a modelo simple
- Todos los predictores estadísticamente significativos
- Apto para deployment en sistema operacional

#### 4. ÁREA CRÍTICA: CALIFICACIÓN BAJA

- Promedio 3.01/5.0 (60.2% de satisfacción)
- Devoluciones inversamente correlacionadas con calificación
- Oportunidad de mejora: Target 4.0+ (80% satisfacción)

## 8.2 Perfil del Cliente de Alto Valor

### Características:

- Edad: 45-55 años
- Frecuencia: Cliente frecuente (Sí)
- Compras: 30+ transacciones
- Visitas: 120+ navegaciones
- Monto: \$170+ USD
- Devoluciones: <3
- Calificación: >4.0/5.0
- Canal: Web preferentemente

## 9. LIMITACIONES Y CONSIDERACIONES

1. **Período Temporal:** Datos históricos. No se consideran variaciones estacionales.
2. **Causalidad vs Correlación:** Correlaciones identificadas no implican causalidad directa.
3. **Valores Faltantes:** 5% imputados; asunciones pueden introducir error.
4. **Generalización:** Recomendaciones aplican a período actual; validación continua necesaria.
5. **Multicolinealidad:** Baja pero presente (e.g., visitas y compras correlacionadas naturalmente).