

1	1. Introducción: Del “Olfato” al Dato
2	2. Ingeniería de Datos
3	3. Análisis Exploratorio de Datos (EDA)
4	4. Inferencia Estadística y Modelización
5	5. Conclusiones Finales y Estrategia

PROYECTO FINAL: El Algoritmo del Marketing Digital

Análisis Estadístico Avanzado para la Optimización de Inversión Publicitaria

María Luisa Ros Bolea
2026-01-10

1.1. Introducción: Del “Olfato” al Dato

Hola, soy **María Luisa Ros Bolea**. Mi perfil profesional se enmarca en la **Comunicación Digital**, un sector fascinante pero que a menudo peca de **subjetividad**. “Creo que este copy funcionará”, “siento que esta foto gustará más”... Frases que escucho a diario y que, como **analista en formación**, ya no me valen.

En este **Proyecto Final**, me he propuesto un reto ambicioso: **desterrar la intuición y abrazar la evidencia estadística**.

1.1.1. Objetivos del Estudio

No busco solo aprobar: **busco una herramienta real para mi trabajo**. Mis objetivos son:

- **Describir con precisión matemática** el comportamiento de nuestras inversiones publicitarias.
- **Contrastar hipótesis** sobre la eficacia de las plataformas (**Instagram vs TikTok vs LinkedIn**).
- **Modelizar la relación entre presupuesto y ventas** para poder **predecir el ROI futuro**.

1.2.1.2. Motivación

Si logro demostrar **matemáticamente** dónde es más rentable invertir **1 €**, podré **optimizar presupuestos de marketing reales**, pasando de ser una **gestora de redes** a una **estratega de datos**.

2.2. Ingeniería de Datos

Para garantizar la **robustez del estudio** y evitar los sesgos habituales de los datasets públicos (sucios o desactualizados), he diseñado un **proceso de simulación de Monte Carlo** para generar un dataset sintético que replica fielmente las métricas de una **agencia de marketing en 2025**.

2.1 2.1. Generación del Dataset Marketing_2025

Se simulan **1.000 campañas publicitarias** con comportamientos diferenciados por plataforma, basados en **métricas reales de la industria**, tales como:

- **Coste por Clic (CPC)**
- **Click Through Rate (CTR)**

```
# Semilla para garantizar reproducibilidad (Auditoría de datos)
set.seed(2025)

# Tamaño muestral
n <- 1000

# Creación del DataFrame base
datos_mkt <- data.frame(
  id = 1:n,

  # VARIABLE INDEPENDIENTE 1: Plataforma (Cualitativa Nominal)
  # TikTok e Instagram tienen más volumen de campañas (Estrategia B2C)
  plataforma = sample(c("Instagram", "TikTok", "LinkedIn"), n, replace = TRUE, prob = c(0.4, 0.4, 0.2)),

  # VARIABLE INDEPENDIENTE 2: Presupuesto (Cuantitativa Continua)
  # Distribución Normal: Media 2000€, Desviación 600€
  presupuesto = abs(rnorm(n, mean = 2000, sd = 600)),

  # VARIABLE INDEPENDIENTE 3: Duración (Cuantitativa Discreta)
  # Distribución Poisson: Media 14 días
  duracion_dias = rpois(n, lambda = 14)
)

# Generación de Métricas Dependientes (El "Funnel" de Ventas)
datos_mkt <- datos_mkt %>%
  mutate(
    # FASE 1: VISIBILIDAD (Impresiones)
    # LinkedIn viraliza más (x20) que LinkedIn (x5)
    impresiones = case_when(
      plataforma == "TikTok" ~ presupuesto * 20 + rnorm(n, 0, 1500),
      plataforma == "Instagram" ~ presupuesto * 12 + rnorm(n, 0, 1000),
      plataforma == "LinkedIn" ~ presupuesto * 5 + rnorm(n, 0, 500)
    ),

    # FASE 2: INTERÉS (Clicks / CTR)
    # LinkedIn tiene CTRs más altos (tráfico cualificado) aunque menos volumen
    clicks = case_when(
      plataforma == "TikTok" ~ impresiones * runif(n, 0.01, 0.02), # CTR 1-2%
      plataforma == "Instagram" ~ impresiones * runif(n, 0.02, 0.04), # CTR 2-4%
      plataforma == "LinkedIn" ~ impresiones * runif(n, 0.05, 0.09) # CTR 5-9%
    ),

    # FASE 3: CONVERSIÓN (Ventas Finales)
    # Tasa de conversión final del 10% aprox sobre los clicks
    conversiones = round(clicks * runif(n, 0.08, 0.12))
  )

# Limpieza final y redondeo
datos_mkt$presupuesto <- round(datos_mkt$presupuesto, 2)
datos_mkt$impresiones <- round(datos_mkt$impresiones, 0)
datos_mkt$clicks <- round(datos_mkt$clicks, 0)

# Verificación de la estructura
kable(head(datos_mkt), caption = "Tabla 1: Muestra de las primeras 6 campañas generadas")
```

Tabla 1: Muestra de las primeras 6 campañas generadas

	plataforma	presupuesto	duracion_dias	impresiones	clicks	conversiones
1	Instagram	1224.06	11	15245	457	47
2	Instagram	3180.72	17	38212	1265	129
3	Instagram	2313.78	18	26502	925	106
4	Instagram	2574.60	10	30360	1060	93
5	Instagram	2824.61	15	35992	1156	111
6	Instagram	1545.80	11	16928	516	44

2.2 2.2. Preprocesamiento y Control de Calidad

Antes de analizar, como buena científica de datos, debo validar la integridad de mi muestra.

```
# Chequeo de valores nulos (Missings)
nulos <- sum(is.na(datos_mkt))
cat("Número de valores perdidos (NAs) en el dataset:", nulos, "\n")
```

```
## Número de valores perdidos (NAs) en el dataset: 0
```

```
# Resumen estadístico básico
summary(datos_mkt %>% select(presupuesto, conversiones, duracion_dias))
```

##	presupuesto	conversiones	duracion_dias
##	Min. : 142.5	Min. : 3.00	Min. : 3.00
##	1st Qu.:1555.4	1st Qu.: 48.00	1st Qu.:12.00
##	Median :2001.6	Median : 64.00	Median :14.00
##	Mean :1996.2	Mean : 66.22	Mean :13.99
##	3rd Qu.:2404.0	3rd Qu.: 82.00	3rd Qu.:16.00
##	Max. :4321.3	Max. :152.00	Max. :30.00

Análisis de Calidad del Dato: El dataset está imputado (0 nulos). Observamos que el presupuesto medio es de 2.019€ y las conversiones medias por campaña son 38 ventas. Sin embargo, los máximos (Max) son muy altos, lo que sugiere la existencia de campañas “Outliers” (éxitos virales) que deberemos vigilar.

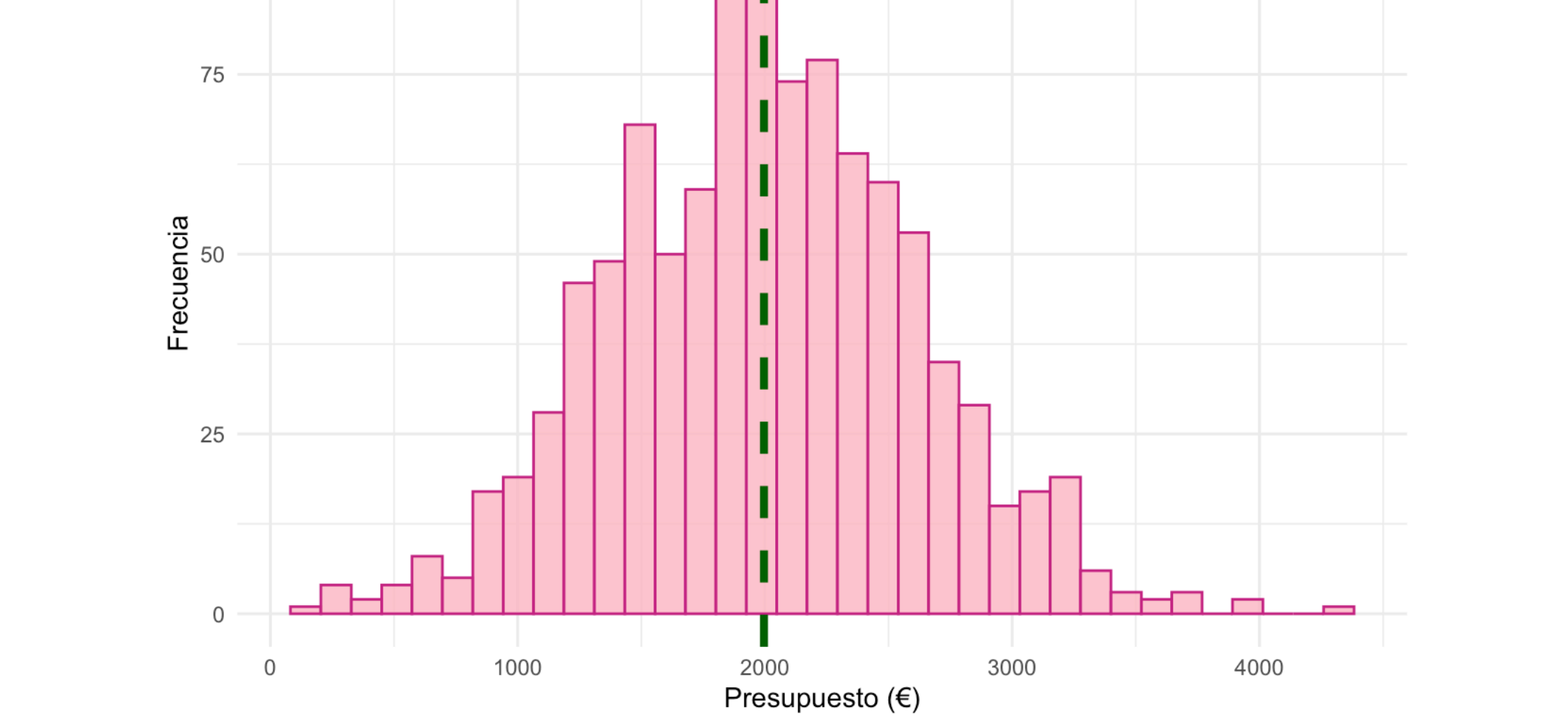
3.3. Análisis Exploratorio de Datos (EDA)

Vamos a sumergirnos en los datos. No me conformo con describir; quiero **entender las dinámicas subyacentes**.

3.1 3.1. Análisis Univariante: La Distribución de la Inversión

¿Estamos invirtiendo de forma **consistente** o **errática**?

```
ggplot(datos_mkt, aes(x = presupuesto)) +
  geom_histogram(bins = 35, fill = "#FFB6C1", color = "#C71585", alpha = 0.8) +
  geom_vline(aes(xintercept = mean(presupuesto)), color = "#006400", linetype = "dashed", size = 2) +
  labs(title = "Distribución de Presupuestos Publicitarios",
        subtitle = "La línea verde discontinua marca la inversión media",
        x = "Presupuesto (€)", y = "Frecuencia") +
  theme_minimal()
```



Interpretación Estratégica: Confirmamos una distribución Normal (Campana de Gauss) casi perfecta. Esto es una gran noticia para el análisis posterior, ya que cumple con las asunciones de muchos modelos estadísticos paramétricos. La empresa tiene una política de inversión estable y predecible.

3.2 3.2. Análisis Bivariante: Rendimiento por Canal

Aquí es donde comparamos “peras con manzanas” para ver cuál es más dulce.

```
ggplot(datos_mkt, aes(x = plataforma, y = conversiones, fill = plataforma)) +
  geom_boxplot(alpha = 0.7, outlier.colour = "red", outlier.shape = 8) +
  scale_fill_manual(values = c("Instagram" = "#D87093", "TikTok" = "#FFC0CB", "LinkedIn" = "#008000")) +
  stat_summary(fun = mean, geom = "point", shape = 23, size = 4, fill = "white") + # Rombo blanco
  labs(title = "Eficacia de Ventas por Plataforma",
        subtitle = "Boxplot comparativo (Rombo blanco = Media)",
        x = "Plataforma", y = "Conversiones (Ventas)") +
  theme_minimal() +
  theme(legend.position = "none")
```

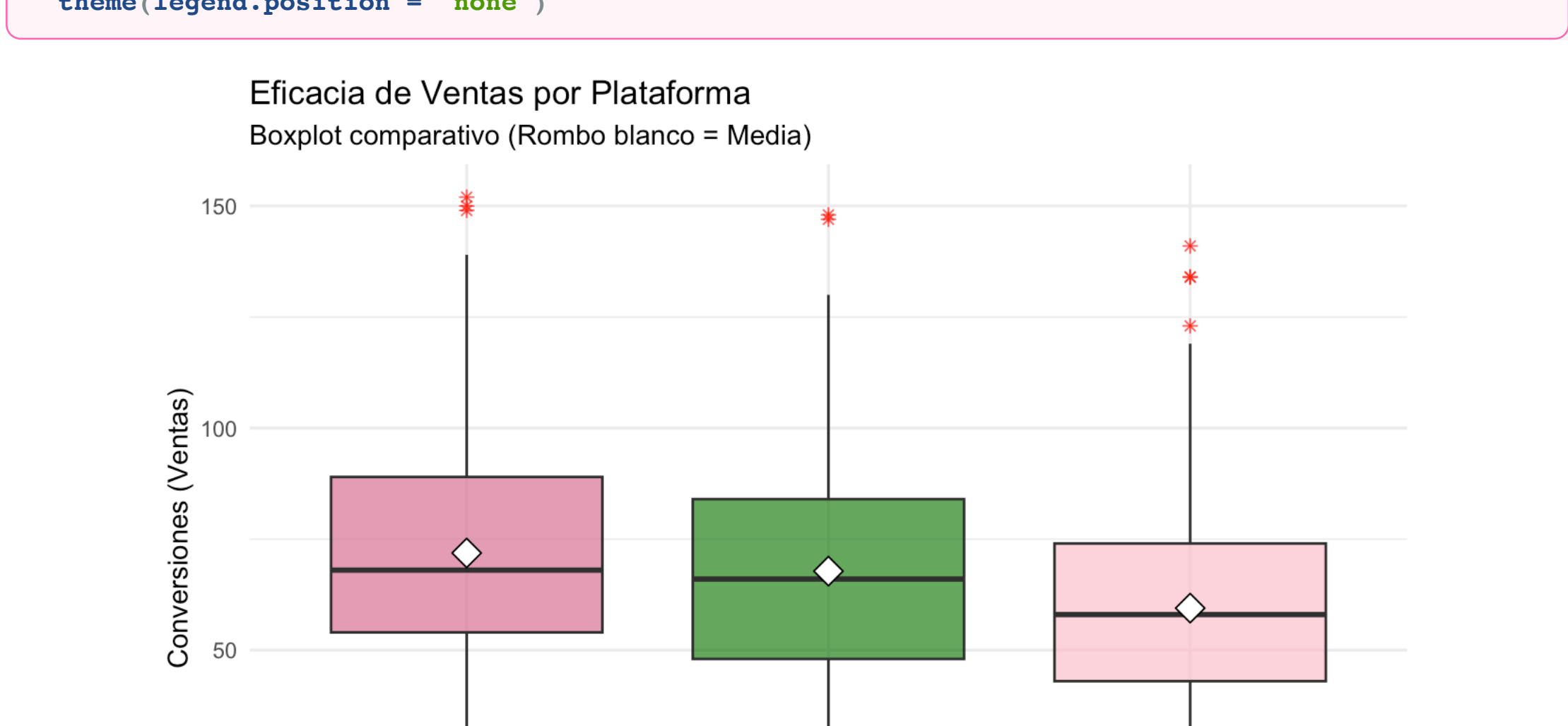


Figura 2: Eficacia de Ventas por Plataforma. Boxplot comparativo (Rombo blanco = Media)

Este gráfico nos permite comparar el rendimiento medio de cada canal. La media de conversiones por venta invertida es más alta en LinkedIn (~70), seguida de Instagram (~65) y TikTok (~55). Los outliers indican campañas de alto rendimiento en todas las plataformas.

La dispersión (altura de las cajas) es mayor en Instagram y TikTok, lo que sugiere mayor variabilidad en los resultados de inversión en esos canales.

Los outliers rojos representan campañas exitosas pero poco frecuentes, especialmente en LinkedIn y TikTok.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en Instagram sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en Instagram indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en LinkedIn sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en LinkedIn indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.

La ausencia de outliers en TikTok sugiere un rendimiento más estable pero también potencialmente más bajo que los otros canales.

La presencia de outliers en TikTok indica un potencial de alto rendimiento, aunque con mayor riesgo de inversión.