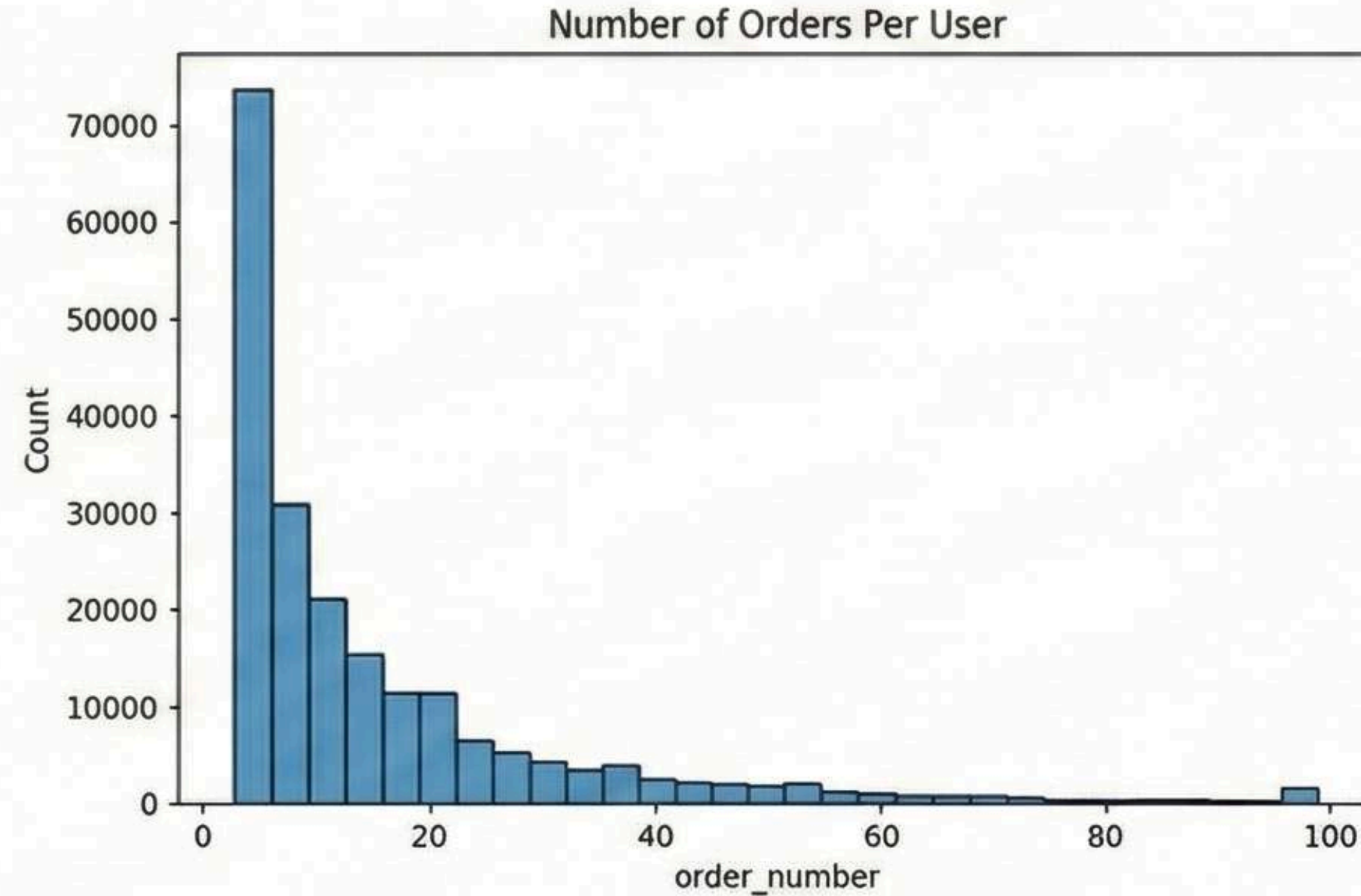


تحليل عدد الطلبات لكل مستخدم

توزيع تكرار الشراء بين العملاء



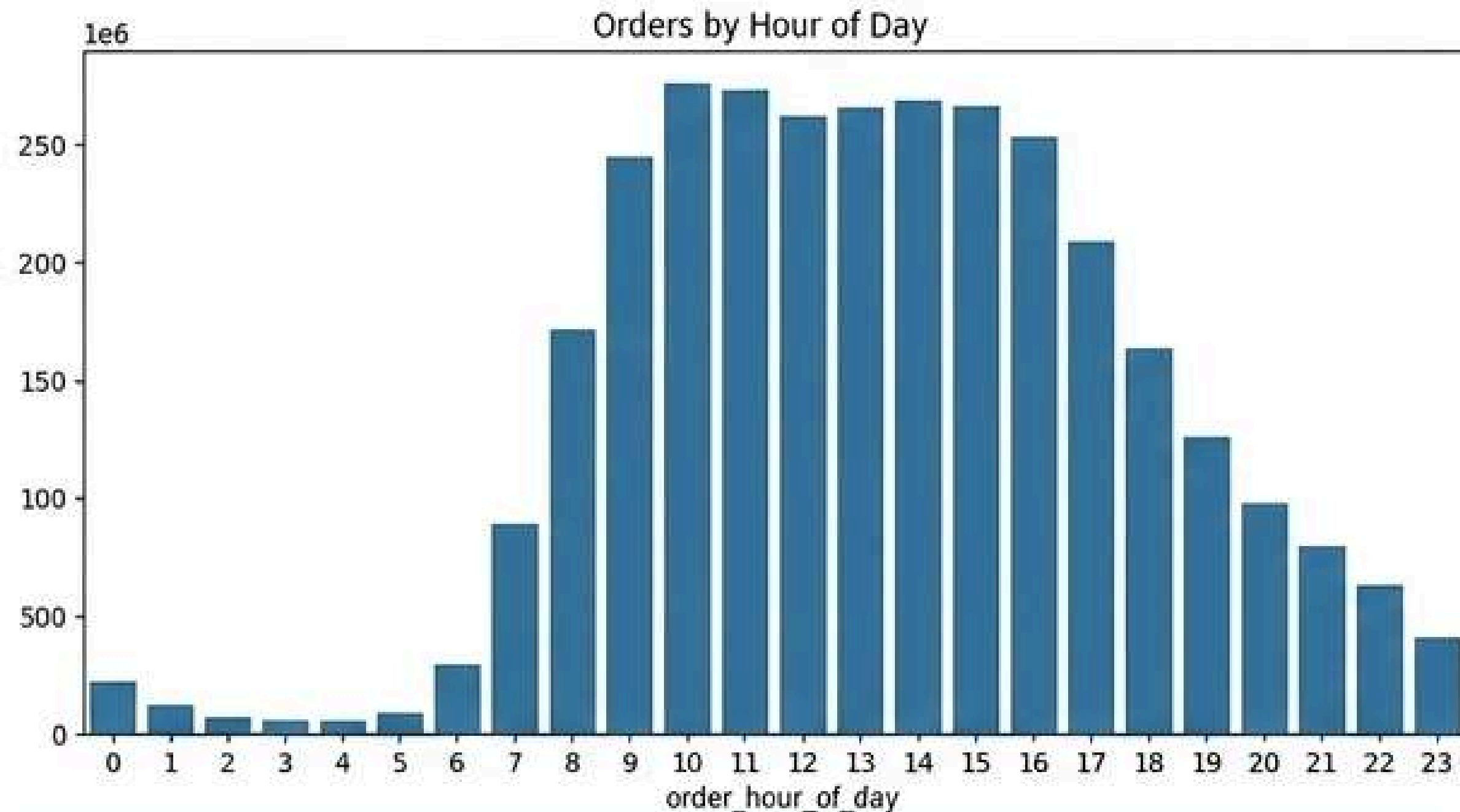
ماذا تخبرنا البيانات؟ 🔍

يوضح الرسم البياني "هستوجرام" توزيع عدد الطلبات التي قام بها كل مستخدم. نلاحظ الأغلبية العظمى من المستخدمين قاموا بعدد قليل من الطلبات. القمة (أعلى عمود) تظهر عند عدد طلبات منخفض جداً (حوالي 4-5 طلبات)، حيث يتجاوز عدد هؤلاء المستخدمين 70,000 مستخدم. ينخفض عدد المستخدمين بشكل حاد وسريع كلما زاد عدد الطلبات. هناك "ذيل طويل" يمثل شريحة من العملاء الأوفياء الذين قاموا بعدد كبير جداً من الطلبات (يصل إلى 100)، ولكن عددهم قليل مقارنة بإجمالي المستخدمين.

تحليل توقيت الطلبات (يومي وساعي)

نظرة على فترات الذروة ونمط الشراء الزمني للعملاء

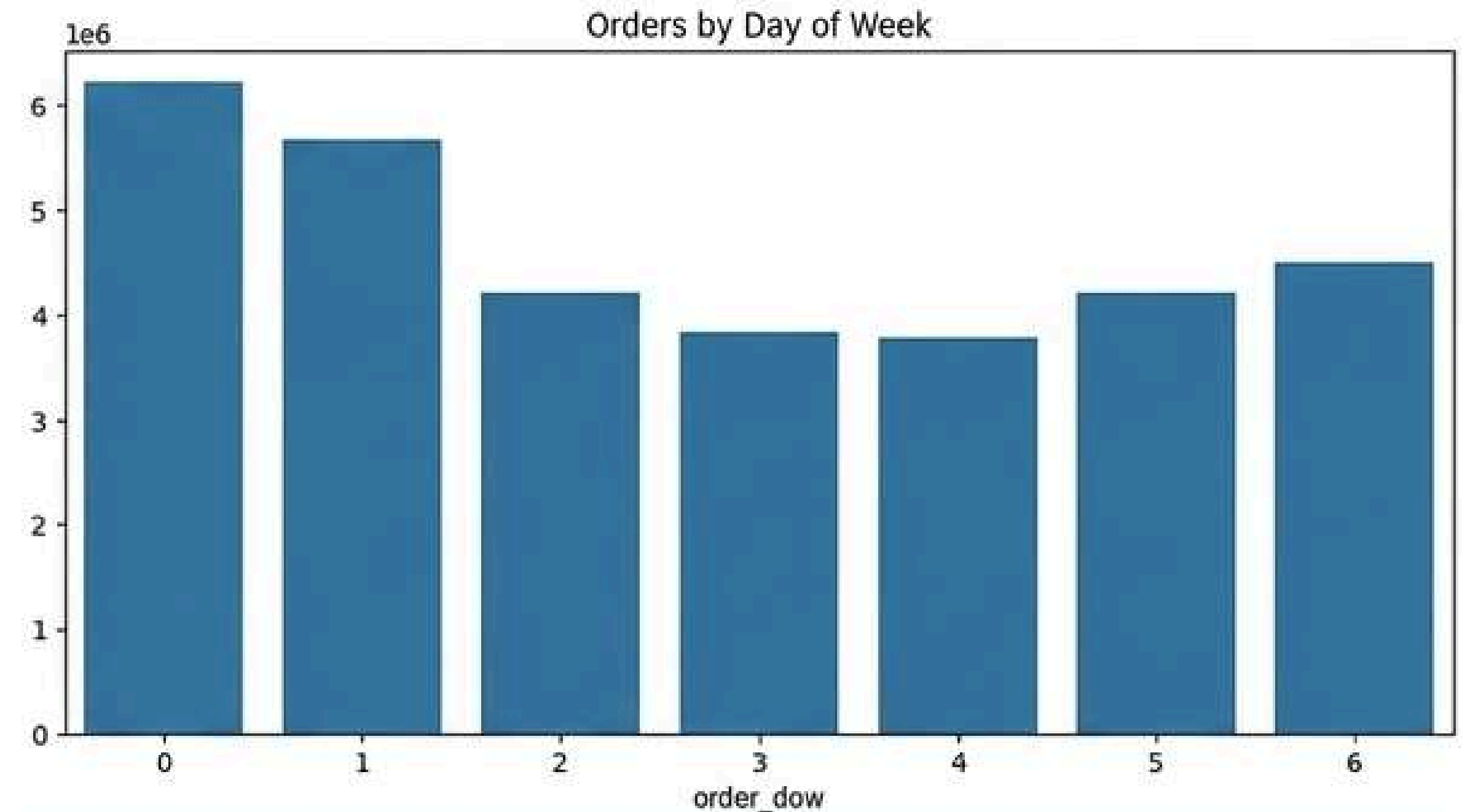
تحليل الطلبات الساعي (Orders by Hour of Day)



ماذا تخبرنا البيانات؟

تبين البيانات أن ذروة نشاط الشراء تكون خلال ساعات النهار، وتحديداً بين الساعة 9 و 4 عصراً (الساعات 9-16)، مع أعلى معدلات طلب حوالي الساعة 10 صباحاً و 2 ظهراً. ينخفض النشاط بشكل كبير في ساعات الصباح الباكر (من منتصف الليل حتى 6 صباحاً).

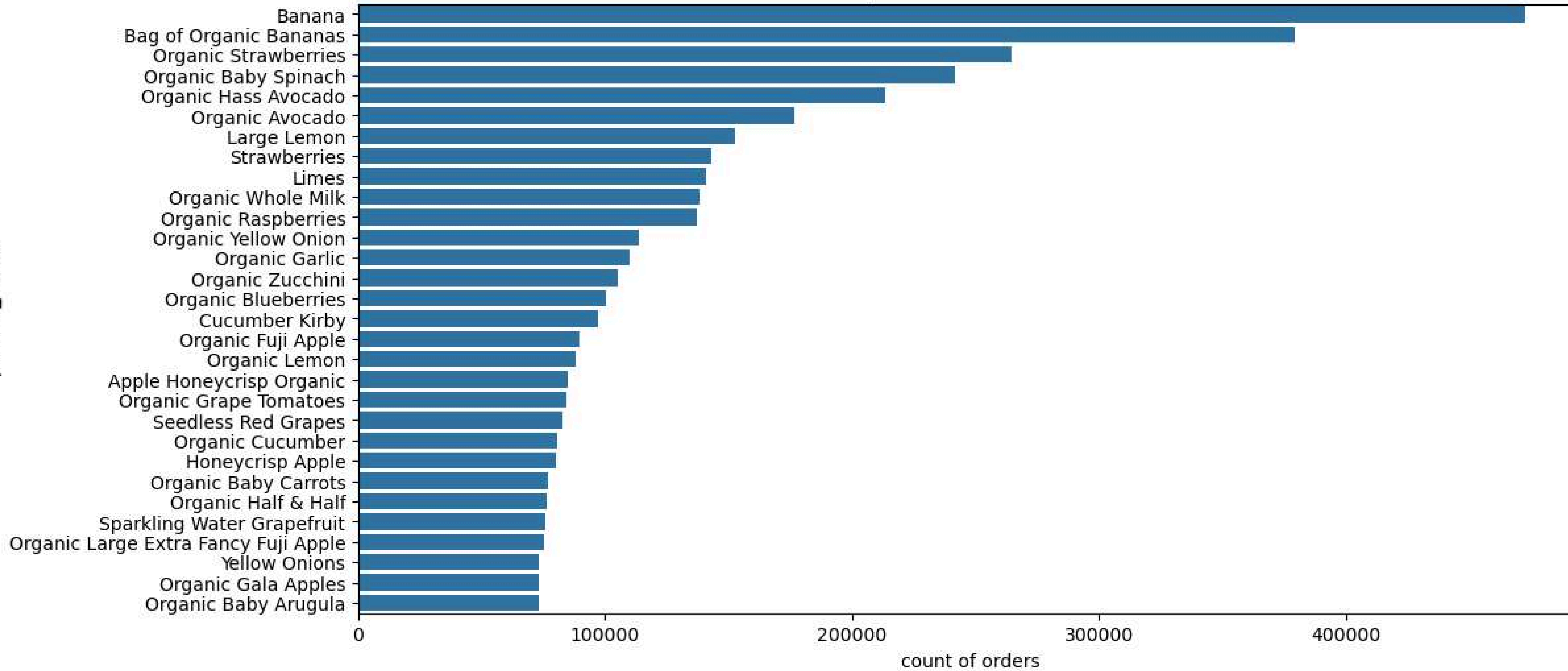
تحليل الطلبات اليومي (Orders by Day of Week)



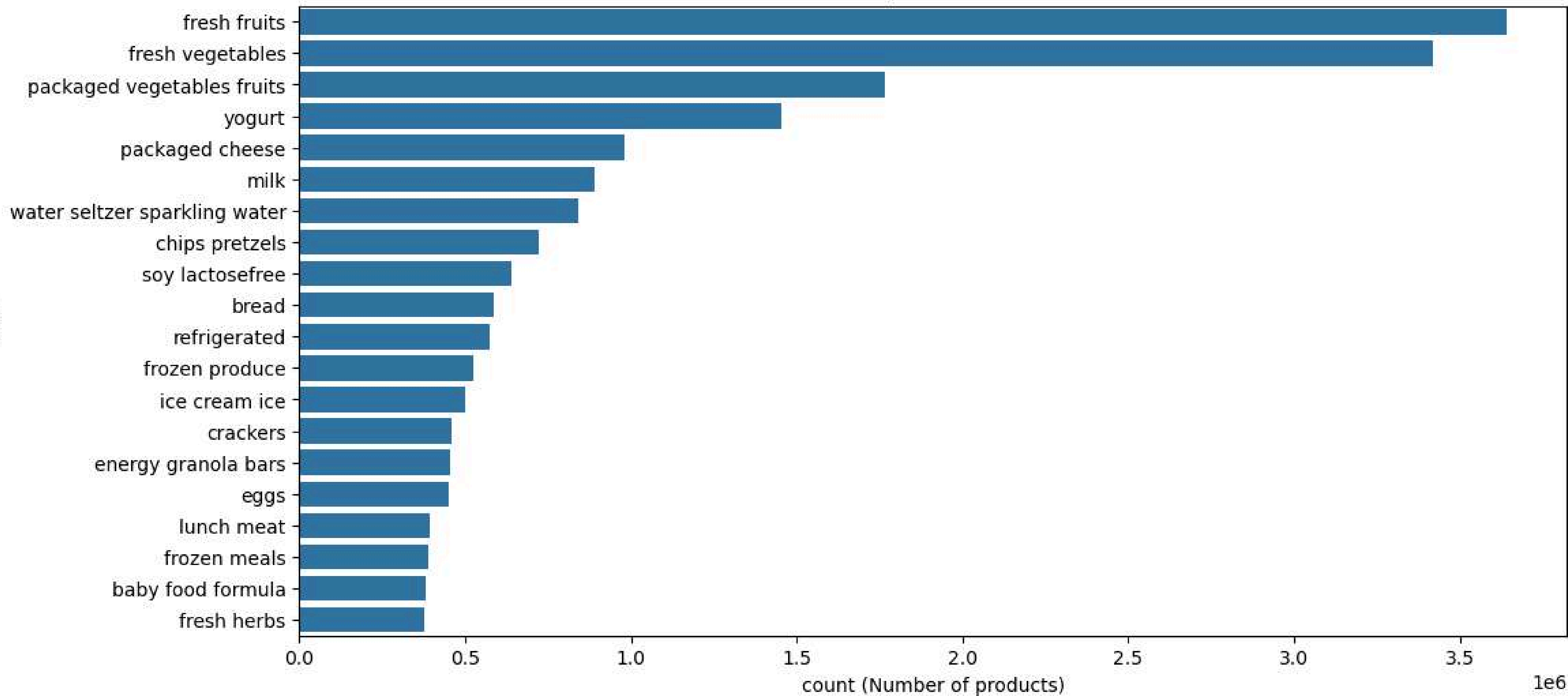
ماذا تخبرنا البيانات؟

يظهر الرسم البياني أن بداية الأسبوع (الأيام 0 و 1) هي الأنشطة في عمليات عمليات الشراء، مع انخفاض ملحوظ في منتصف الأسبوع (الأيام 3 و 4)، ثم ارتفاع طفيف في نهايته (الأيام 5 و 6).

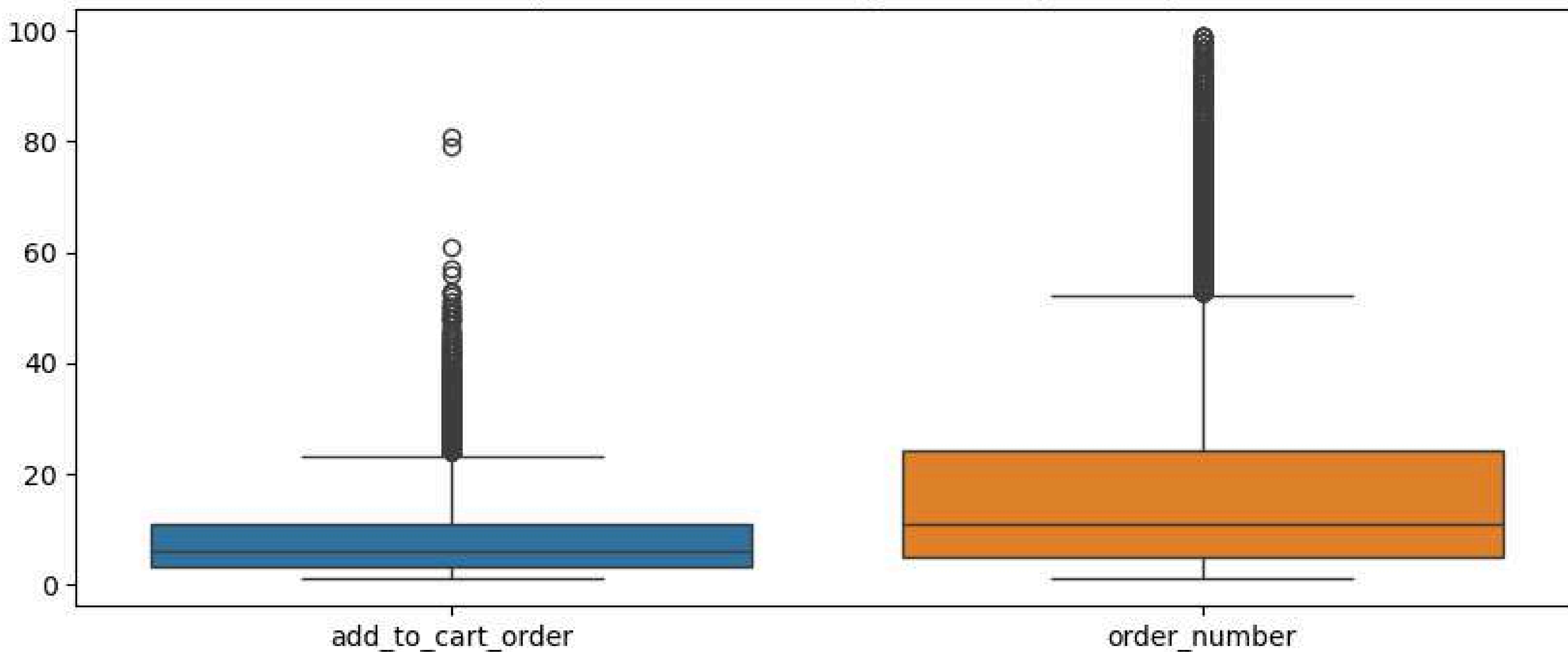
Top 30 Product Ordred



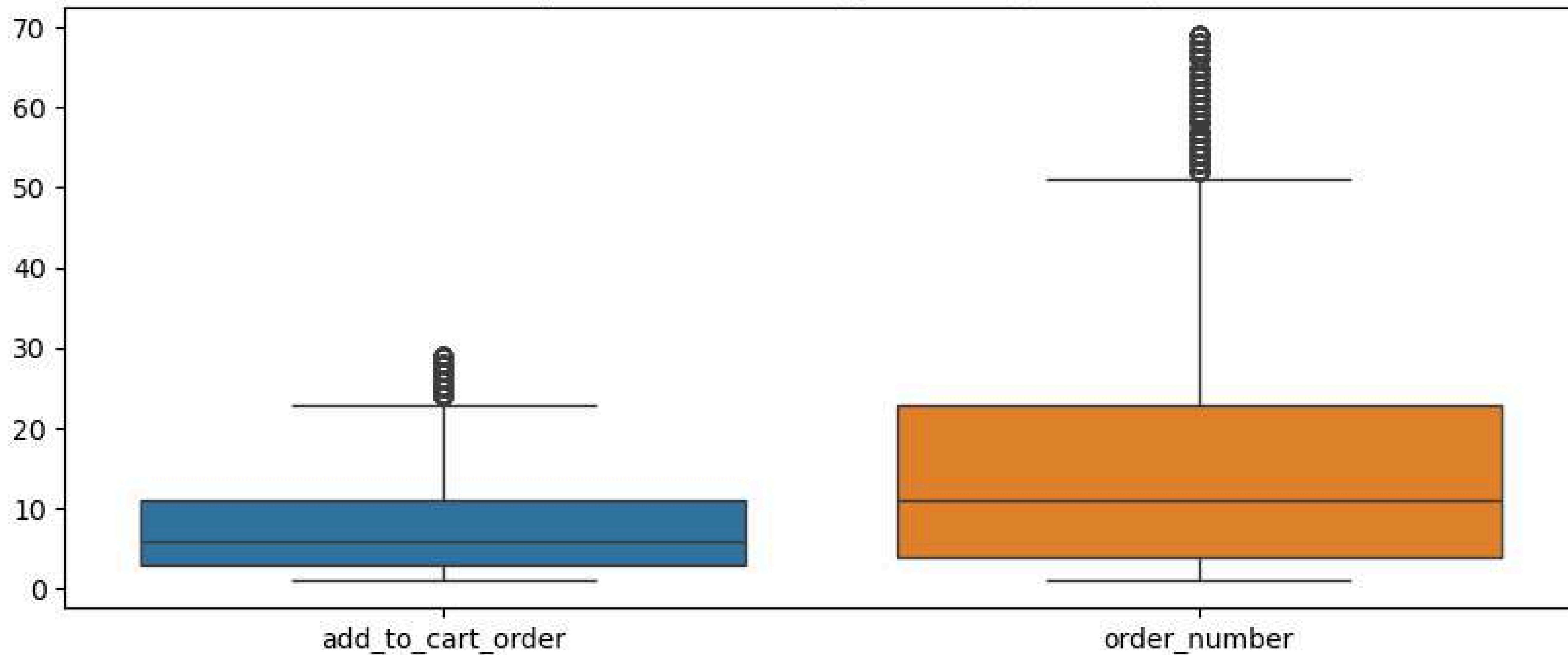
Top 20 aisle



Boxplot Before Removing Outliers (Z-score)



Boxplot After Removing Outliers (Z-score)



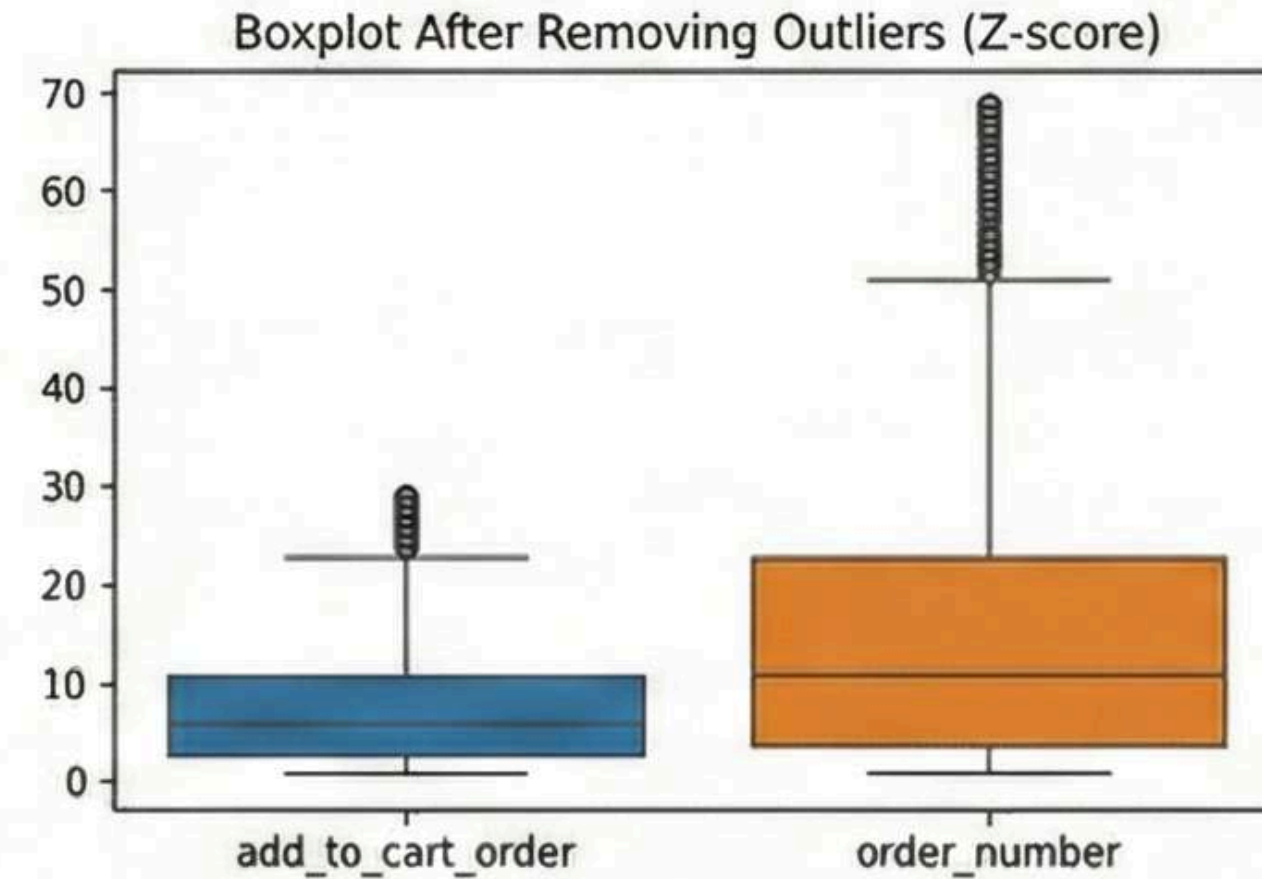
معالجة القيم المتطرفة (Outlier Treatment)

تأثير إزالة القيم الشاذة باستخدام طريقة Z-score على توزيع المتغيرات

بعد إزالة القيم المتطرفة (Z-score) (After Removing Outliers - Z-score)

ماذا تخبرنا البيانات؟

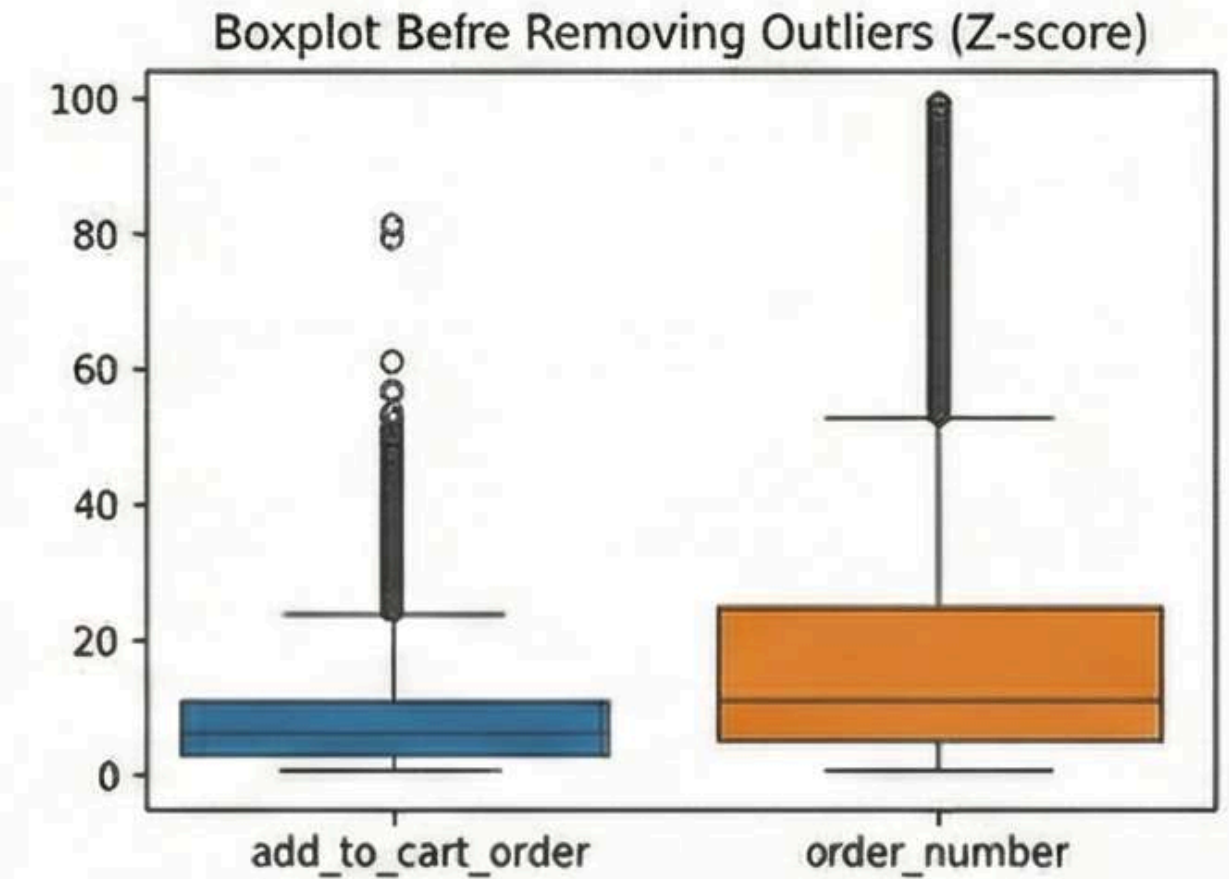
بعد تطبيق طريقة Z-score، تم
تقليص عدد القيم المتطرفة
بشكل كبير. أصبحت القيم
القصوى لـ `add_to_cart_order`
عند 30 و `order_number` عند 70
تقريباً. هذا يؤدي إلى توزيع
بيانات أكثر نظافة وموثوقية
لبناء النماذج.



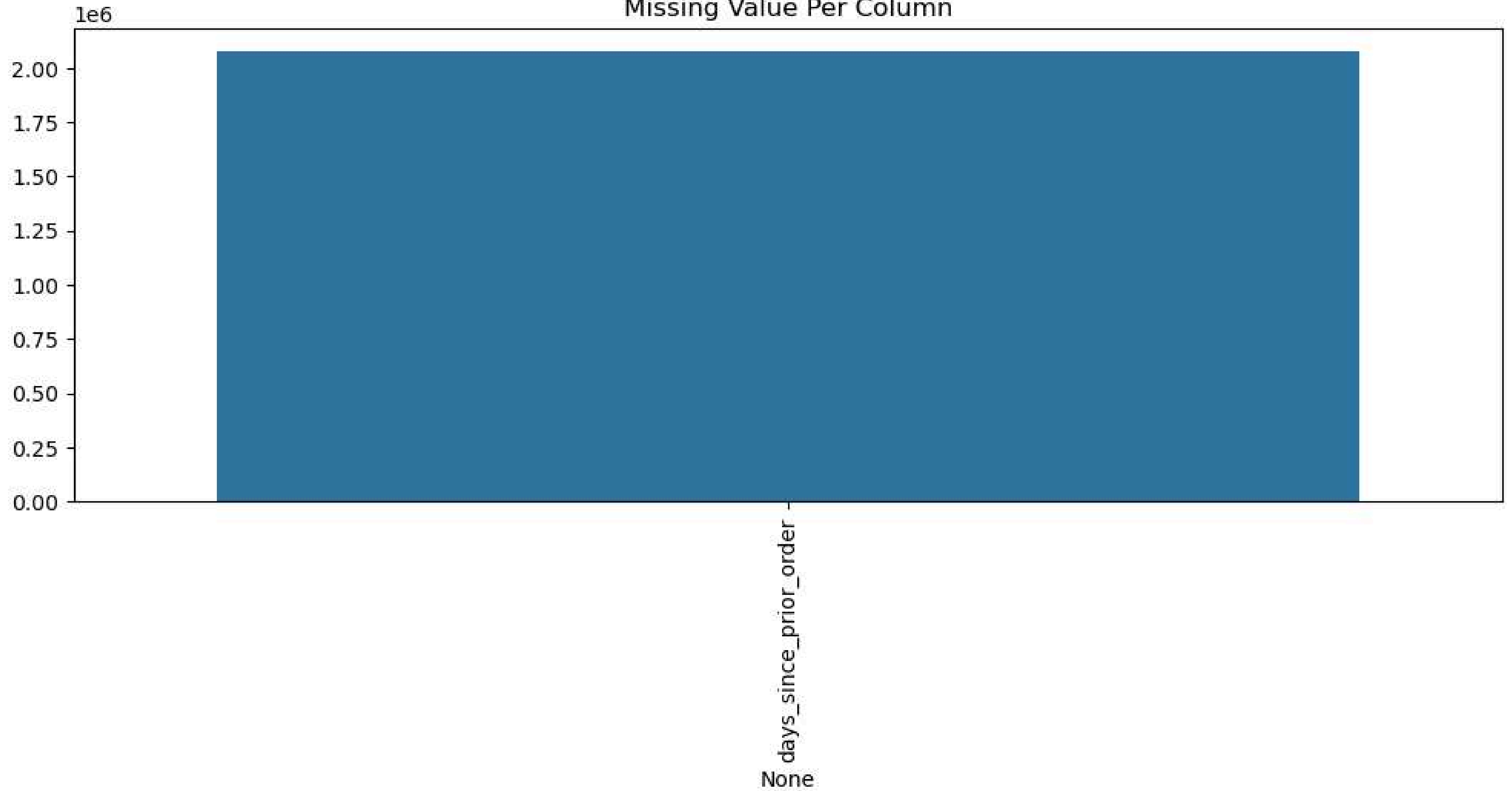
قبل إزالة القيم المتطرفة (Z-score) (Before Removing Outliers - Z-score)

ماذا تخبرنا البيانات؟

يوضح المخطط الصندوقي قبل
المعالجة وجود عدد كبير من
القيم المتطرفة (النقاط
الفردية خارج "الشوارب"). في
`add_to_cart_order`، تصل
القيم إلى 80، بينما في
`order_number`، تصل إلى 100،
مما يشير إلى وجود ملاحظات
شاذة قد تؤثر على التحليل
الإحصائي.



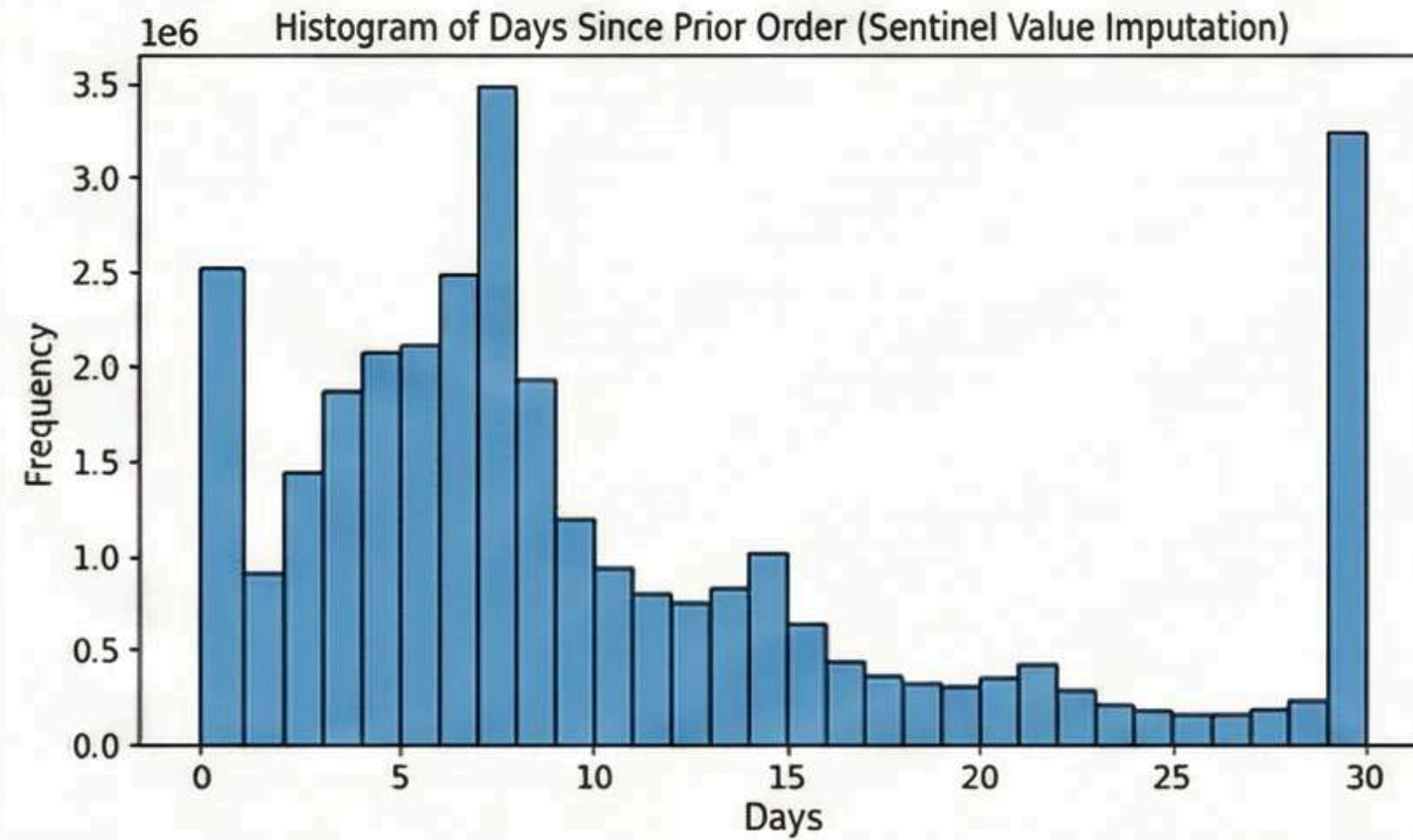
Missing Value Per Column



مقارنة طرق معالجة القيم المفقودة لفترة بين الطلبات

تحليل تأثير الطرق المختلفة على توزيع البيانات واختيار الطريقة الأنسب منطقياً

ملء بالقيمة الثابتة 0 (Sentinel Value - 0) ★ الخيار الأنسب

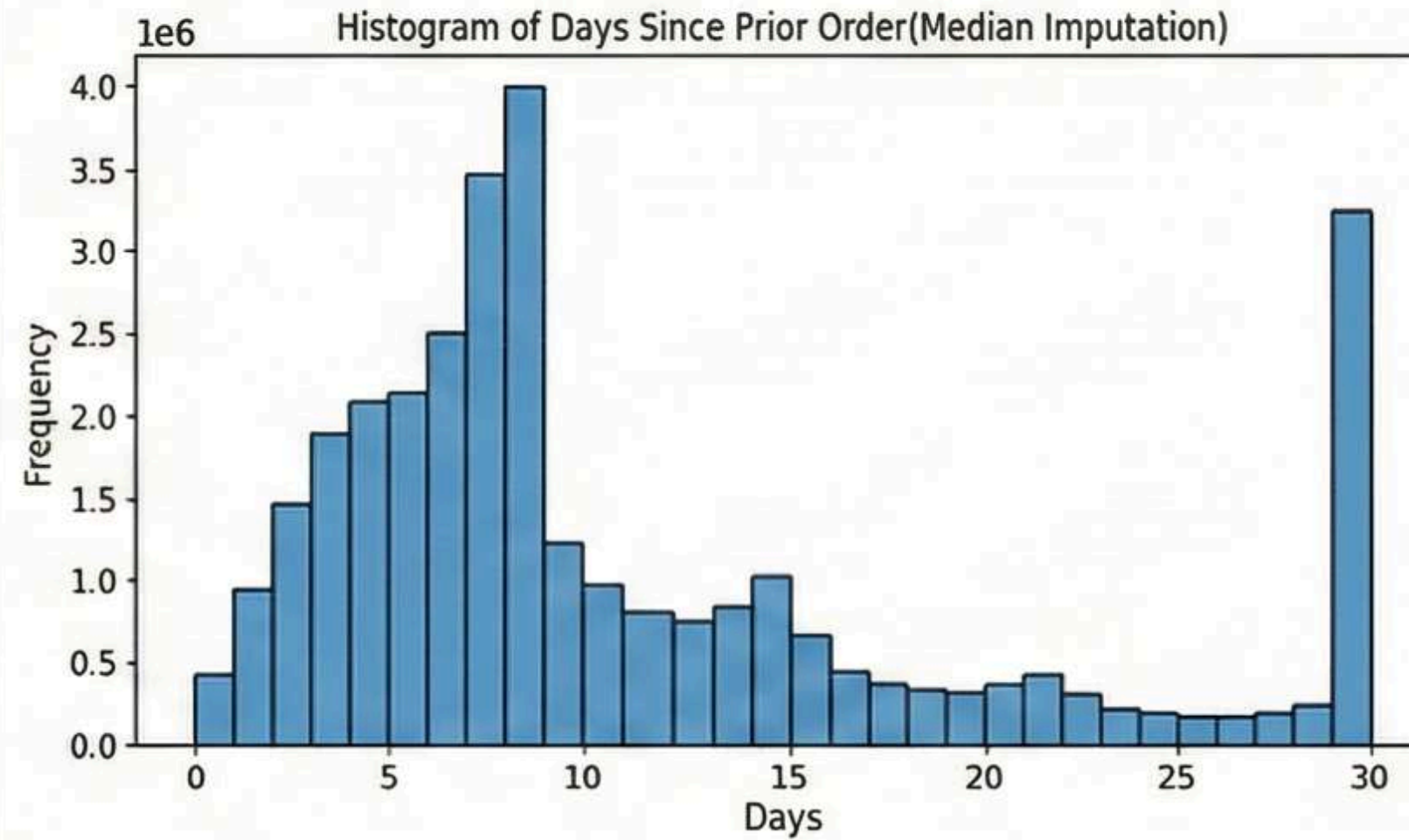


ماذا تخبرنا البيانات؟

اختيار القيمة الثابتة 0 لملء القيم المفقودة هو الأنسب والأكثر منطقية. هذا العمود الكبير عند الصفر يمثل بدقة 'الطلب الأول' للعميل (حيث لا توجد فترة سابقة).

بالإضافة إلى ذلك، نرى الأنماط المعتادة والممثلة بشكل صحيح لإعادة الطلب الأسبوعي (7 أيام) والشهري (30 يوماً).

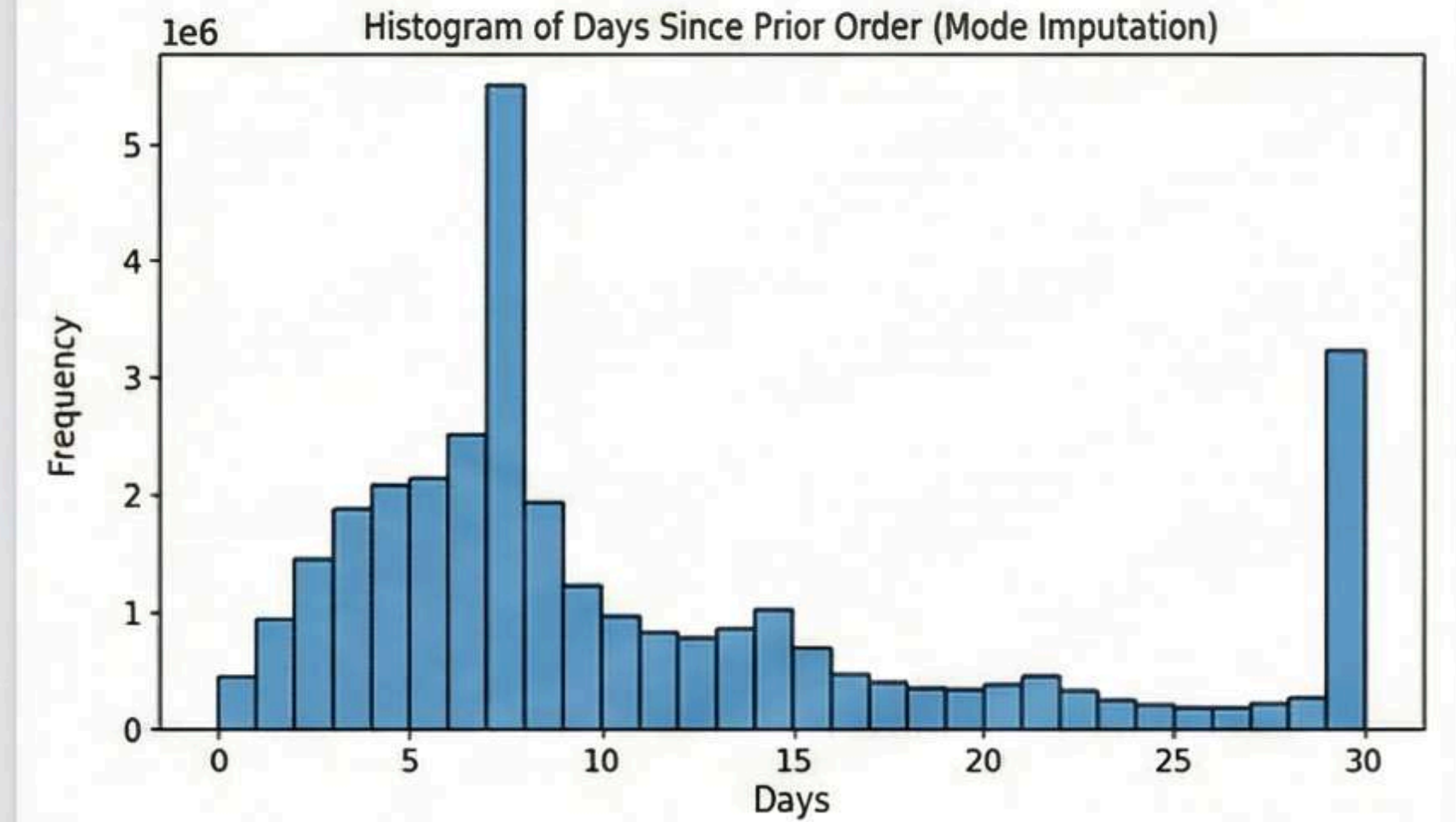
طريقة بالوسيط (Median Imputation)



ماذا تخبرنا البيانات؟

استخدام القيمة المتوسطة (الوسيط) يعطي توزيعاً مشابهاً للمنوال، مع قمم واضحة عند 7 و 30 يوماً، لكن القمة عند 7 أيام تبدو تبدو تدبره أقل حدة قليلاً، مما قد يكون أكثر تمثيلاً للوسط الحقيقي للبيانات مقارنة بالمنوال.

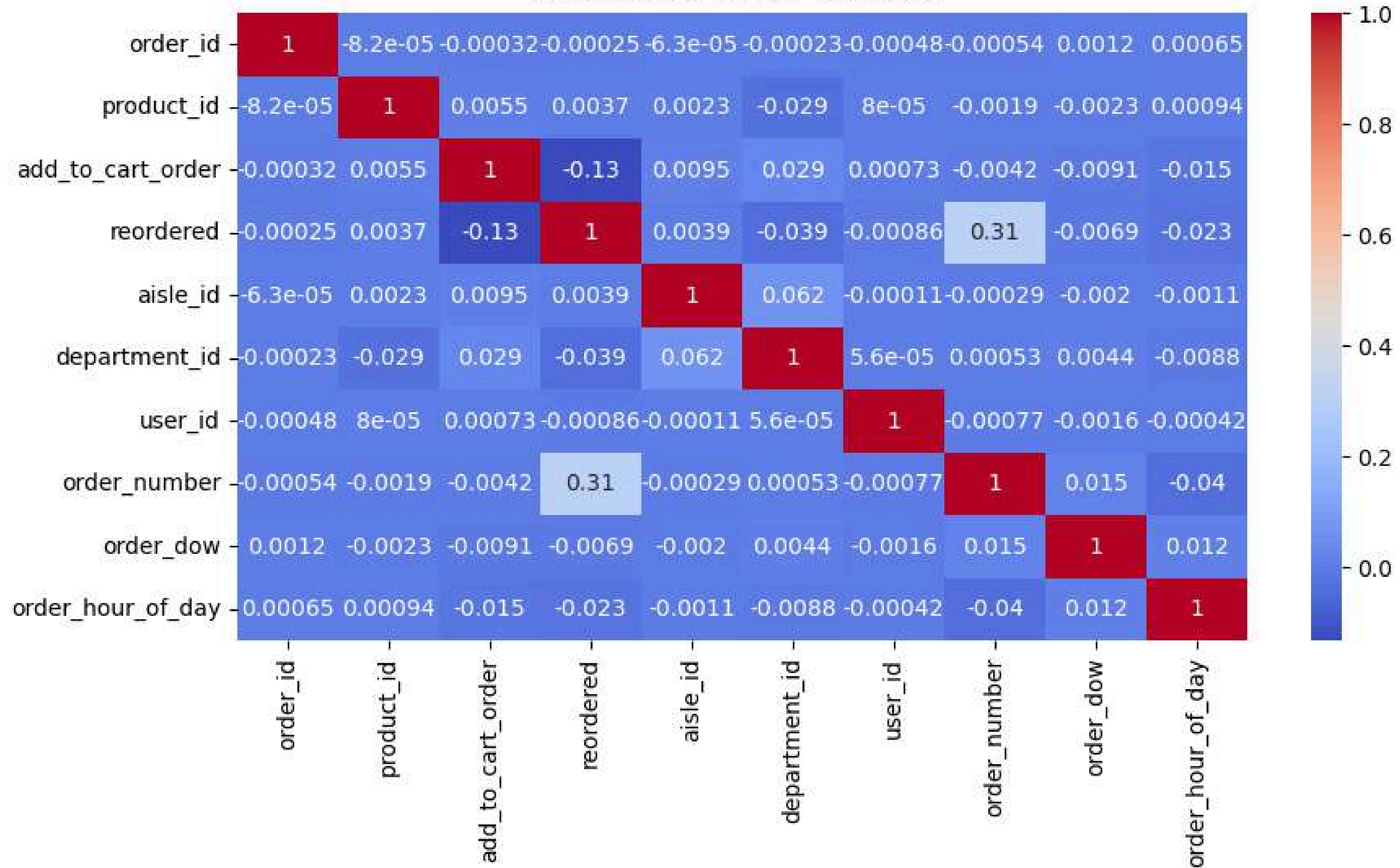
طريقة بالمنوال (Mode Imputation)



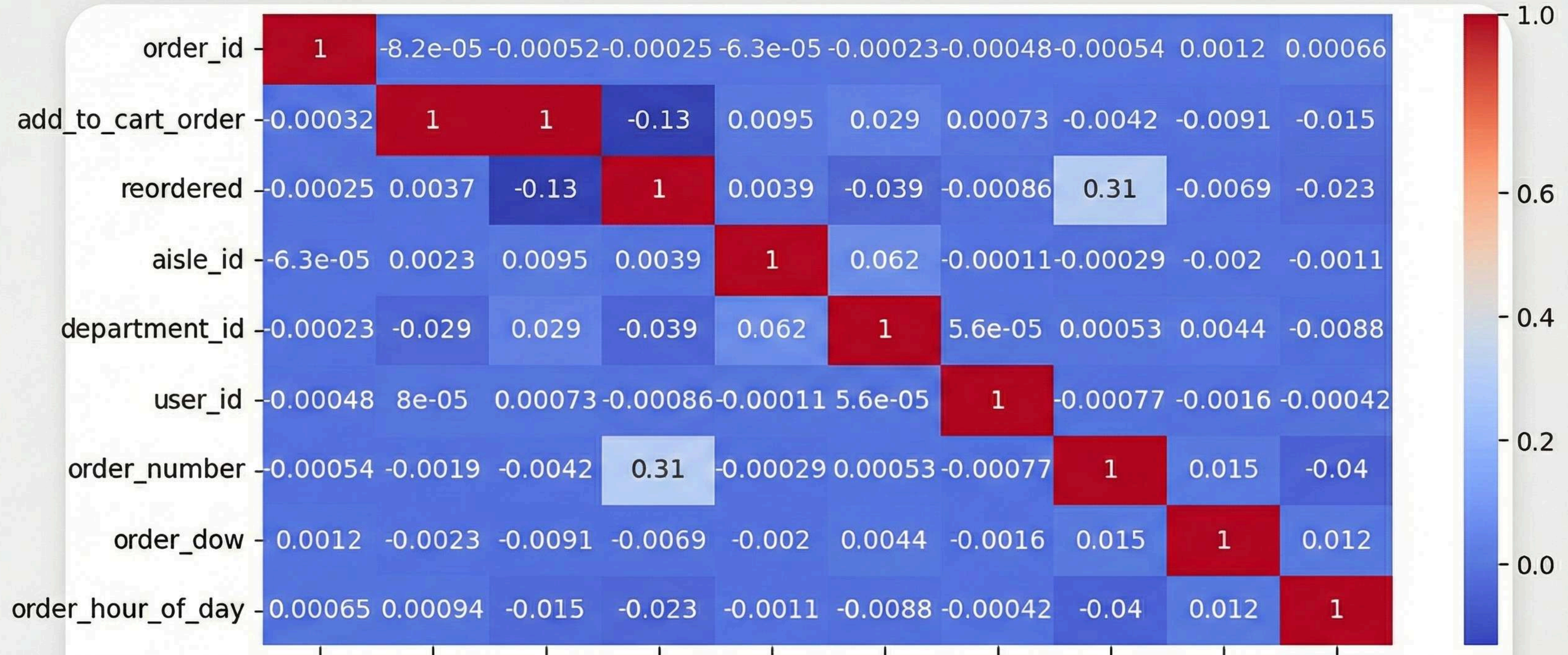
ماذا تخبرنا البيانات؟

استخدام القيمة الأكثر تكراراً (المنوال) يُظهر قمة حادة جداً وغير طبيعية عند 7 أيام، مما قد يؤدي إلى تحيز في التحليل والمبالغة في تقدير تكرار هذه الفترة الزمنية المحددة.

Numerical Feature Correlation



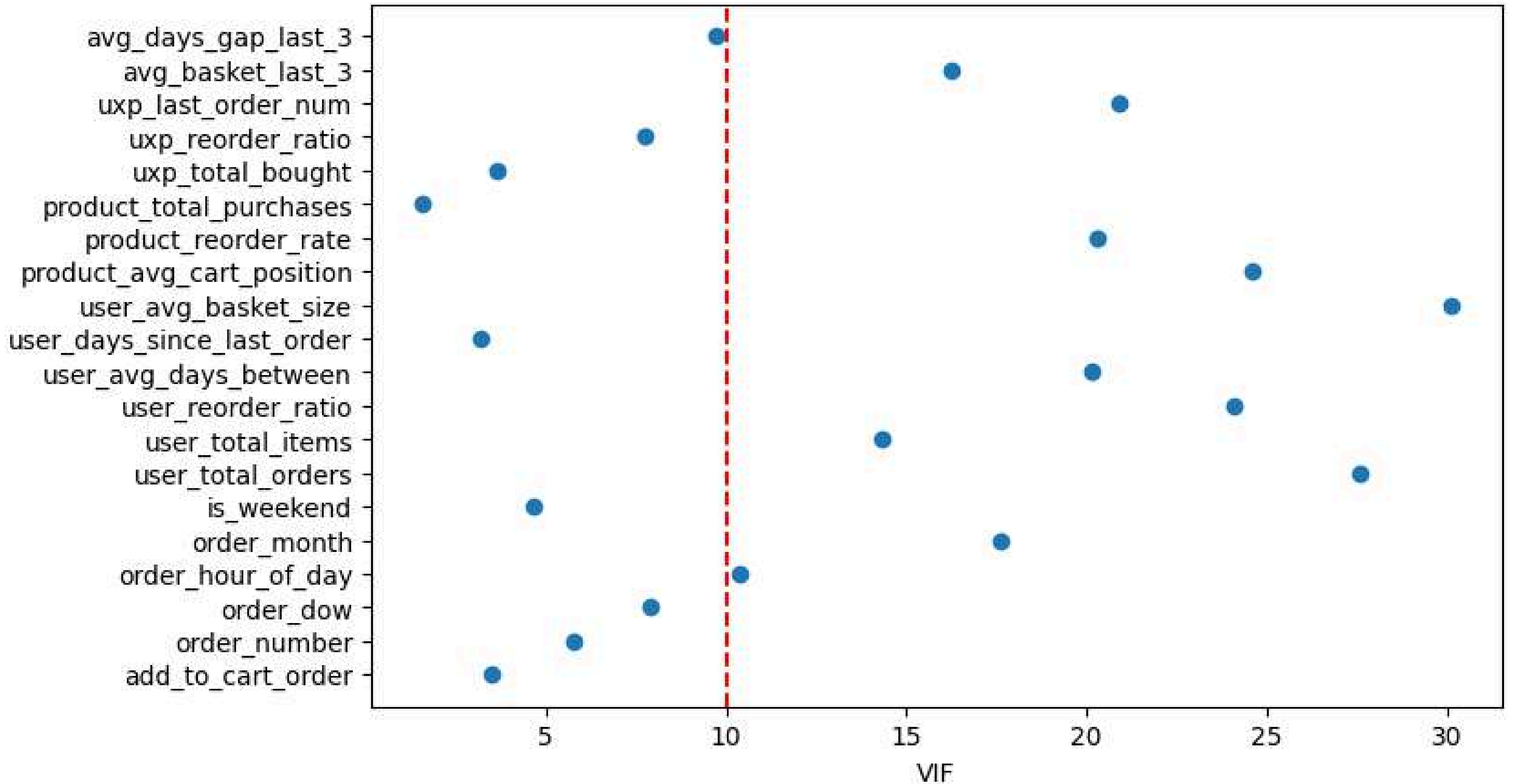
تحليل نسبة إعادة الطلب وعدم إعادته



ماذا تخبرنا البيانات؟

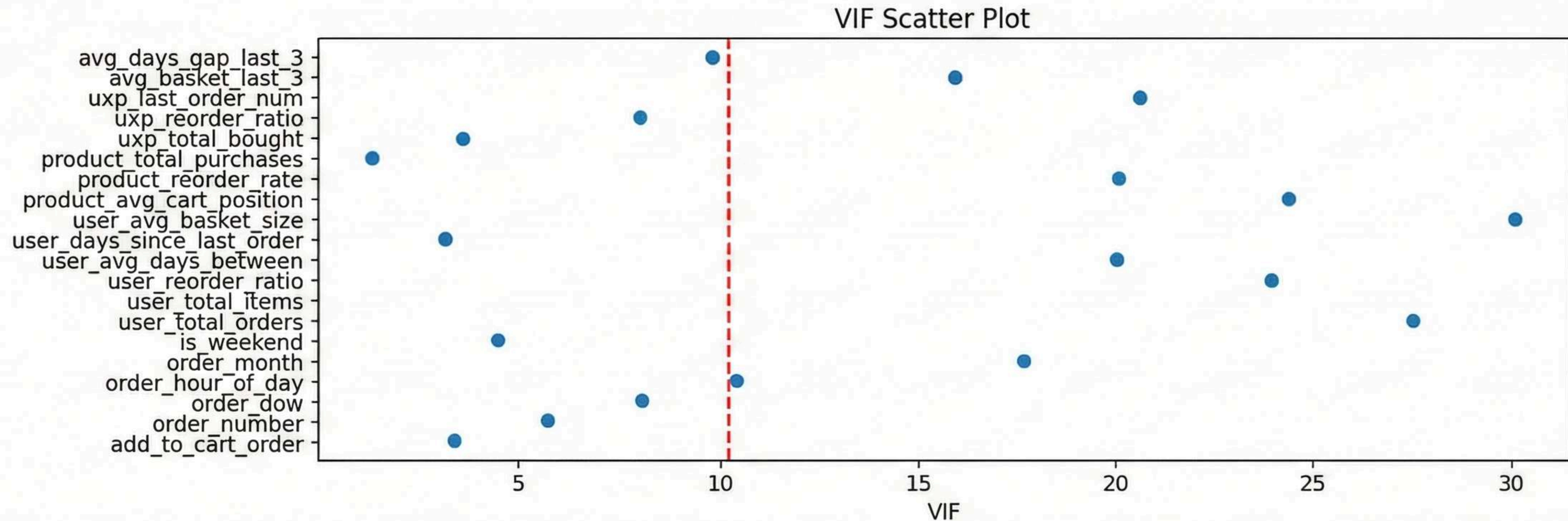
يوضح هذا الرسم البياني العلاقة القوية بين ميزة 'reordered' (إعادة الطلب) وميزة 'order_number' (رقم الطلب)، تشير القيمة العالية (1) باللون الأحمر الداكن إلى وجود ارتباط تام في حالات إعادة الطلب، بينما تشير القيمة المنخفضة (0.31) باللون الأزرق الفاتح إلى نسبة عدم إعادة الطلب. باقي المتغيرات تظهر ارتباطات ضعيفة جدًا.

VIF Scatter Plot



تحليل التعددية الخطية (VIF)

تقييم ارتباط المتغيرات المستقلة للكشف عن التعددية الخطية التي قد تؤثر على دقة النماذج



ماذا تخبرنا البيانات؟

يوضح الرسم البياني قيم عامل تضخم التباين (VIF) لكل متغير مستقل. النقاط التي تقع على يمين الخط الأحمر المتقطع ($VIF > 10$) تشير إلى وجود تعددية خطية عالية قد تؤثر سلباً على نموذج الانحدار، مثل `'user_avg_basket_size'` و `'product_reorder_r'`. بينما المتغيرات على اليسار ($VIF < 10$) تعتبر ذات ارتباط مقبول ولا تمثل مشكلة كبيرة، مثل `'order_dow'` و `'uxp'` و `'product_total_purchases'` و `'is_weekend'`.