

# Documento Final Instacart

## Pipelines

### pipeline\_instacartdb:

1. Data Loader: Se encarga de sacar los datos de MySQL. Lo hace sacando dfs de cada tabla
2. Data Exporter: Se encarga de exportar los datos hacia el schema RAW. Lo hace enviando los dfs a formato sql

### pipeline\_instacartsnowf:

1. Data Loader: Se encarga de sacar los datos del schema RAW de Snowflake. Lo hace sacando dfs de cada tabla
2. Transformer: Se encarga de realizar las transformaciones necesarias descritas en el plan de acción. Las realiza en los dfs con pandas
3. Data Exporter: Se encarga de exportar los datos hacia el schema CLEAN. Lo hace enviando los dfs a formato sql

## Plan de acción de acuerdo al EDA

### 1. Limpieza y transformación de datos

- **Valores nulos:**

- Tabla INSTACART\_ORDERS: Para DAYS\_SINCE\_PRIOR\_ORDER, los valores nulos corresponden a la primera compra de los usuarios y llenarlos con 0.
- Tabla ORDER\_PRODUCTS: Para ADD\_TO\_CART\_ORDER, los valores nulos pueden ser eliminados
- Tabla PRODUCTS: Para PRODUCT\_NAME, asignar una etiqueta como Desconocido.

- **Duplicados:**

- Remover los 15 registros duplicados en INSTACART\_ORDERS.

- **Conversión de tipos:**

- Convertir DAYS\_SINCE\_PRIOR\_ORDER (Tabla INSTACART\_ORDERS) y ADD\_TO\_CART\_ORDER (Tabla ORDER\_PRODUCTS) a enteros (int)

- **Outliers:**

- Representan comportamientos válidos de los clientes; por lo que, no se modificaran

## **Normalización**

Si necesitaramos un modelo de machine learning

1. **Tabla:** INSTACART\_ORDERS

- ORDER\_HOUR\_OF\_DAY: Escalar los valores entre 0 y 1.

## **Modelo Dimensional (Star Schema)**

### **Tablas de Dimensiones:**

1. **dim\_products**

- PRODUCT\_ID (PK)
- PRODUCT\_NAME
- AISLE\_ID (FK)
- DEPARTMENT\_ID (FK)

2. **dim\_aisles**

- AISLE\_ID (PK)
- AISLE

3. **dim\_departments**

- DEPARTMENT\_ID (PK)
- DEPARTMENT

### **Tabla de Hechos:**

1. **fct\_order\_products**

- ORDER\_ID (FK)
- PRODUCT\_ID (FK)
- ADD\_TO\_CART\_ORDER
- REORDERED

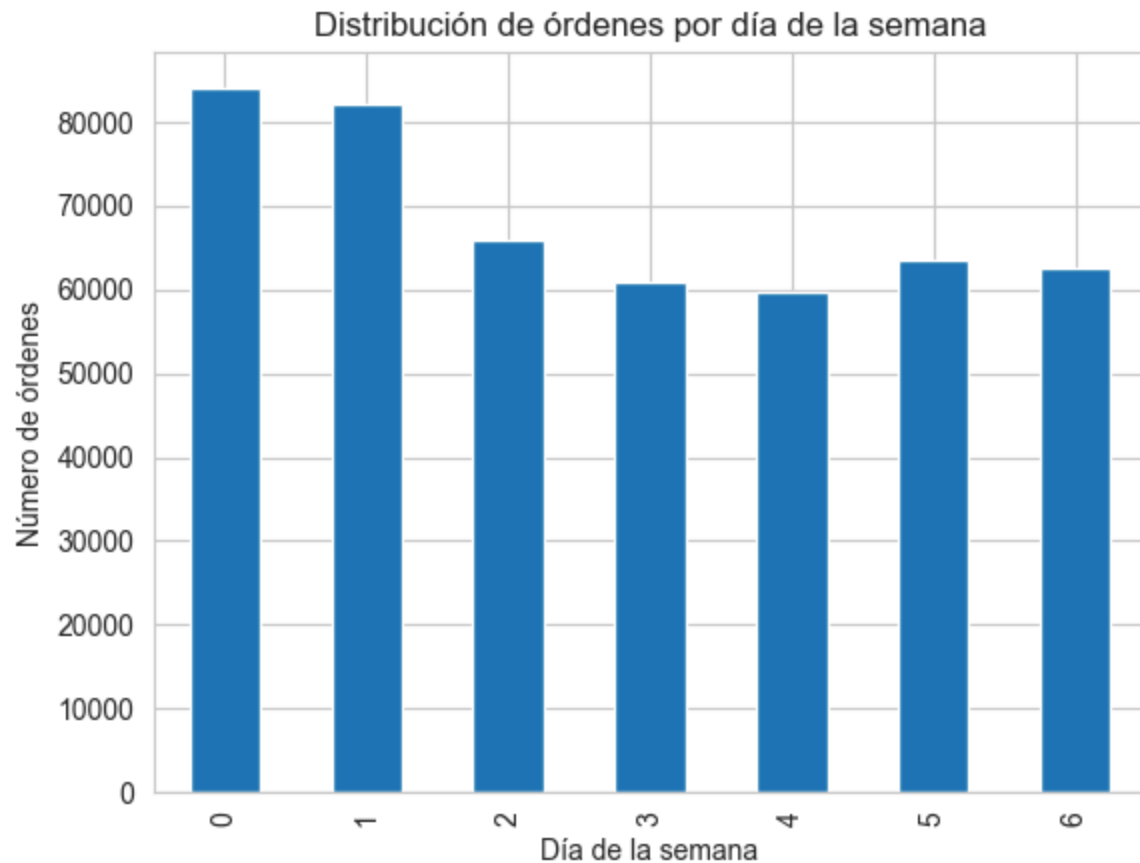
## 2. **fct\_instacart\_orders**

- ORDER\_ID (PK)
- USER\_ID (PK)
- ORDER\_NUMBER
- ORDER\_DOW
- ORDER\_HOUR\_OF\_DAY
- DAYS\_SINCE\_PRIOR\_ORDER

## **Preguntas**

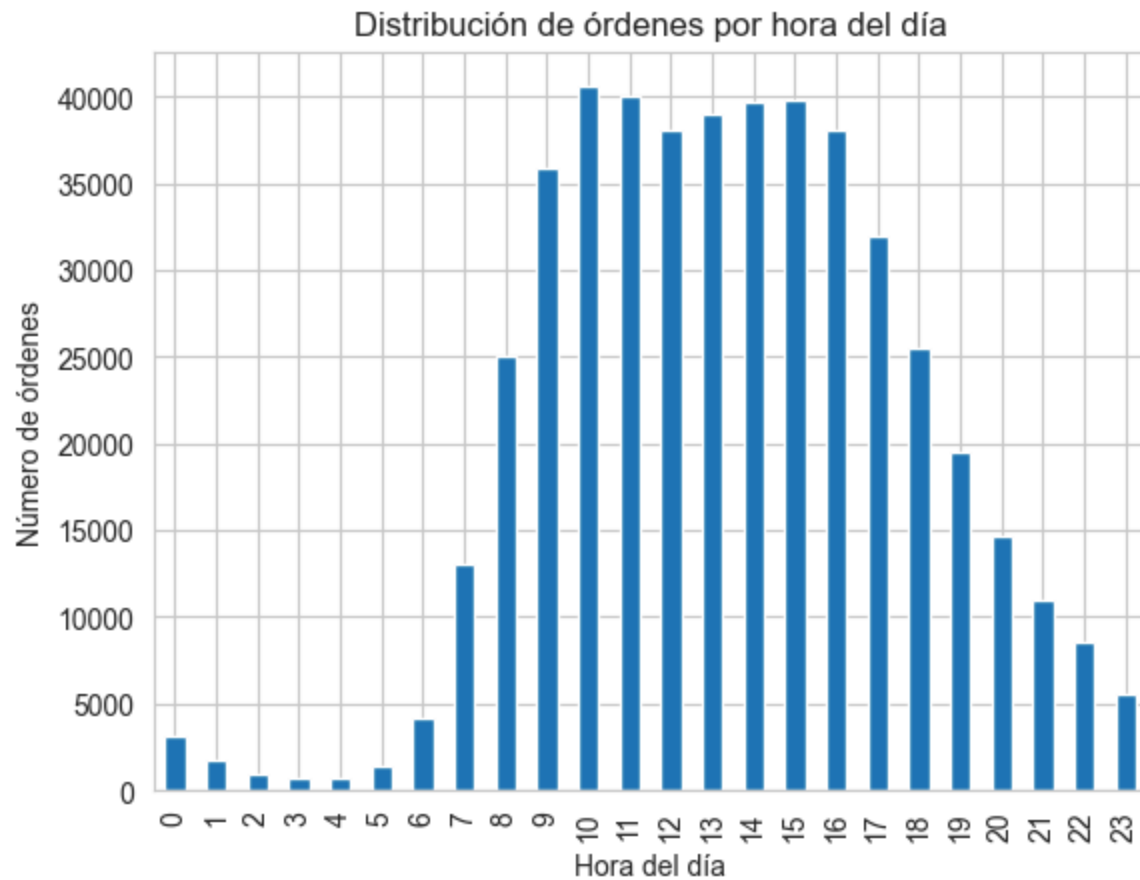
### **1. Comportamiento de compra según día de la semana**

- El gráfico muestra que el día con más órdenes es el **domingo (0)**, seguido del **sábado (6)**. Los días de la semana (lunes a viernes) tienen menos órdenes en comparación con los fines de semana.
- Los **plátanos (Banana)** son los más comprados en casi todos los días, seguidos de productos como **fresas orgánicas (Organic Strawberries)**, **espinacas orgánicas (Organic Baby Spinach)**, y **aguacates orgánicos (Organic Hass Avocado)**.



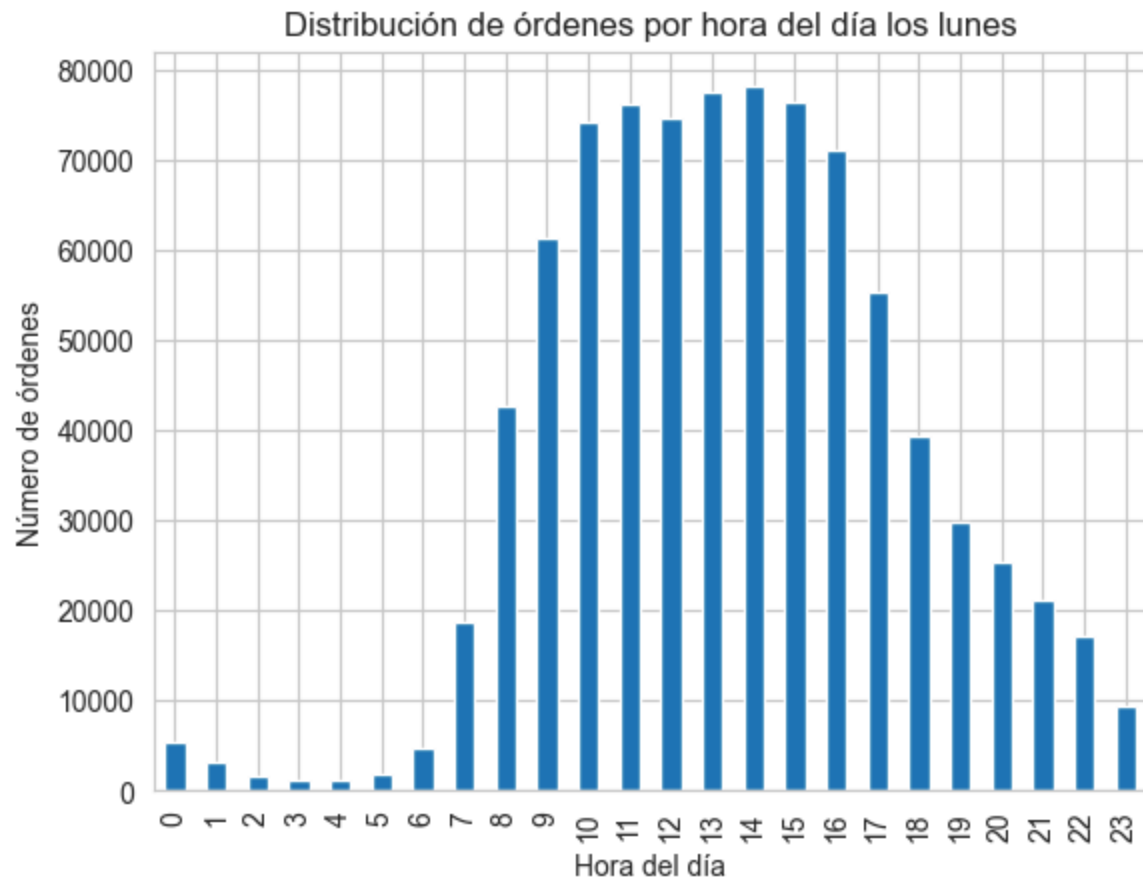
## 2. Comportamiento de compra según hora del día

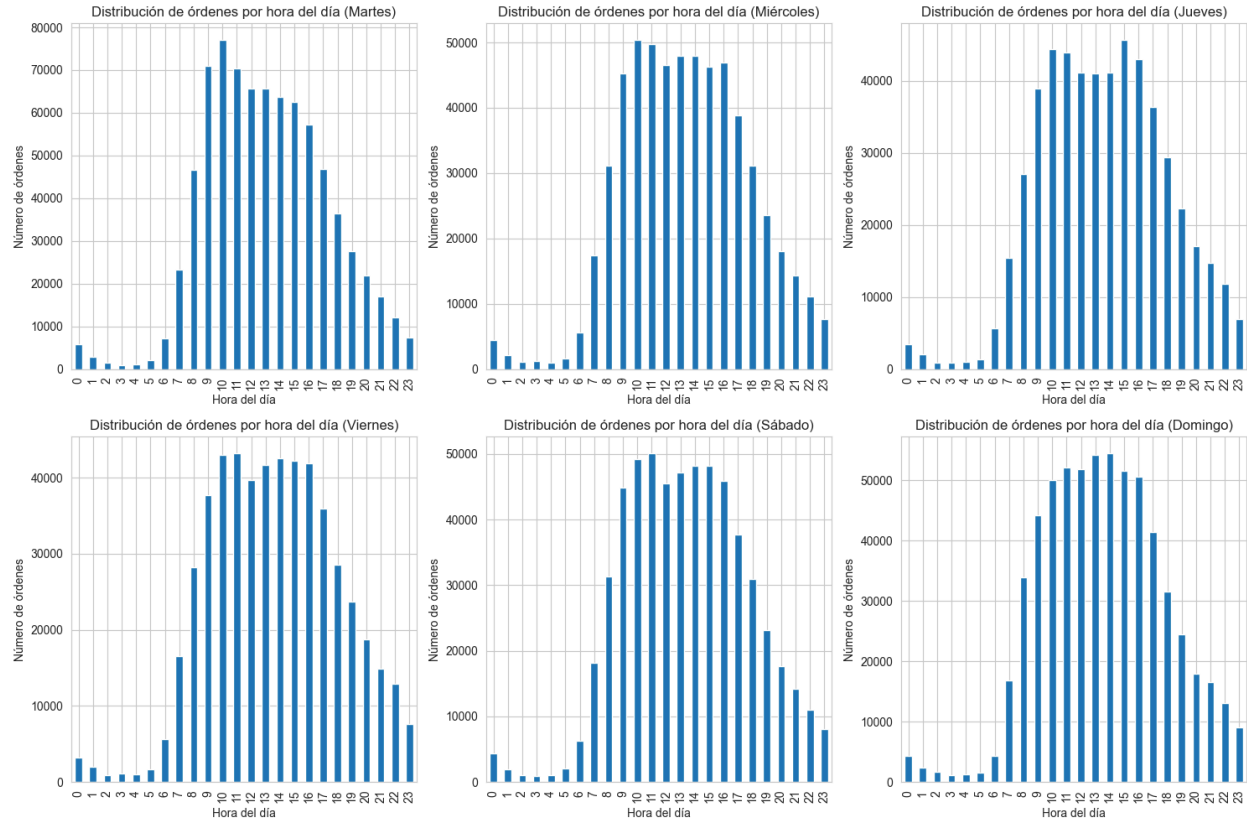
- El gráfico muestra que la mayoría de las órdenes se realizan entre las **8:00 y las 18:00 horas**, con un pico alrededor de las **10:00 a.m.**. Las horas con menos órdenes son las primeras horas de la mañana (antes de las 8:00) y las últimas horas de la noche (después de las 20:00).
- Se analizan los productos más comprados en diferentes franjas horarias (mañana, tarde, y de 9 a 17 horas). Los productos más populares en todas las franjas horarias son los mismos: **plátanos, fresas orgánicas, espinacas orgánicas, y aguacates orgánicos**.



### 3. Comportamiento según hora del día y día de la semana

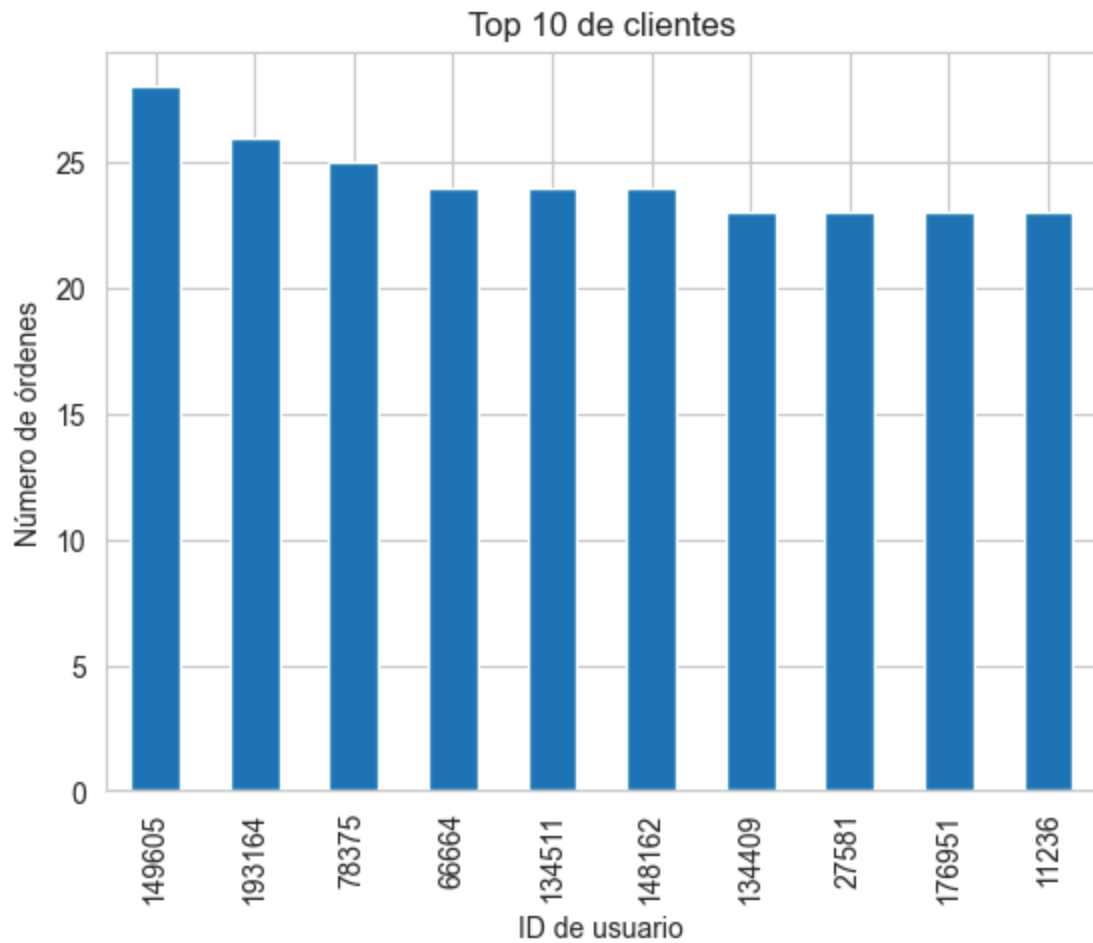
- Los patrones de compra por hora son similares en todos los días de la semana, con un pico en las horas de la mañana (alrededor de las 10:00 a.m.). Sin embargo, los fines de semana (sábado y domingo) tienen un volumen de órdenes ligeramente más alto en comparación con los días de la semana.





## 4. Distribución de las órdenes hechas por los clientes

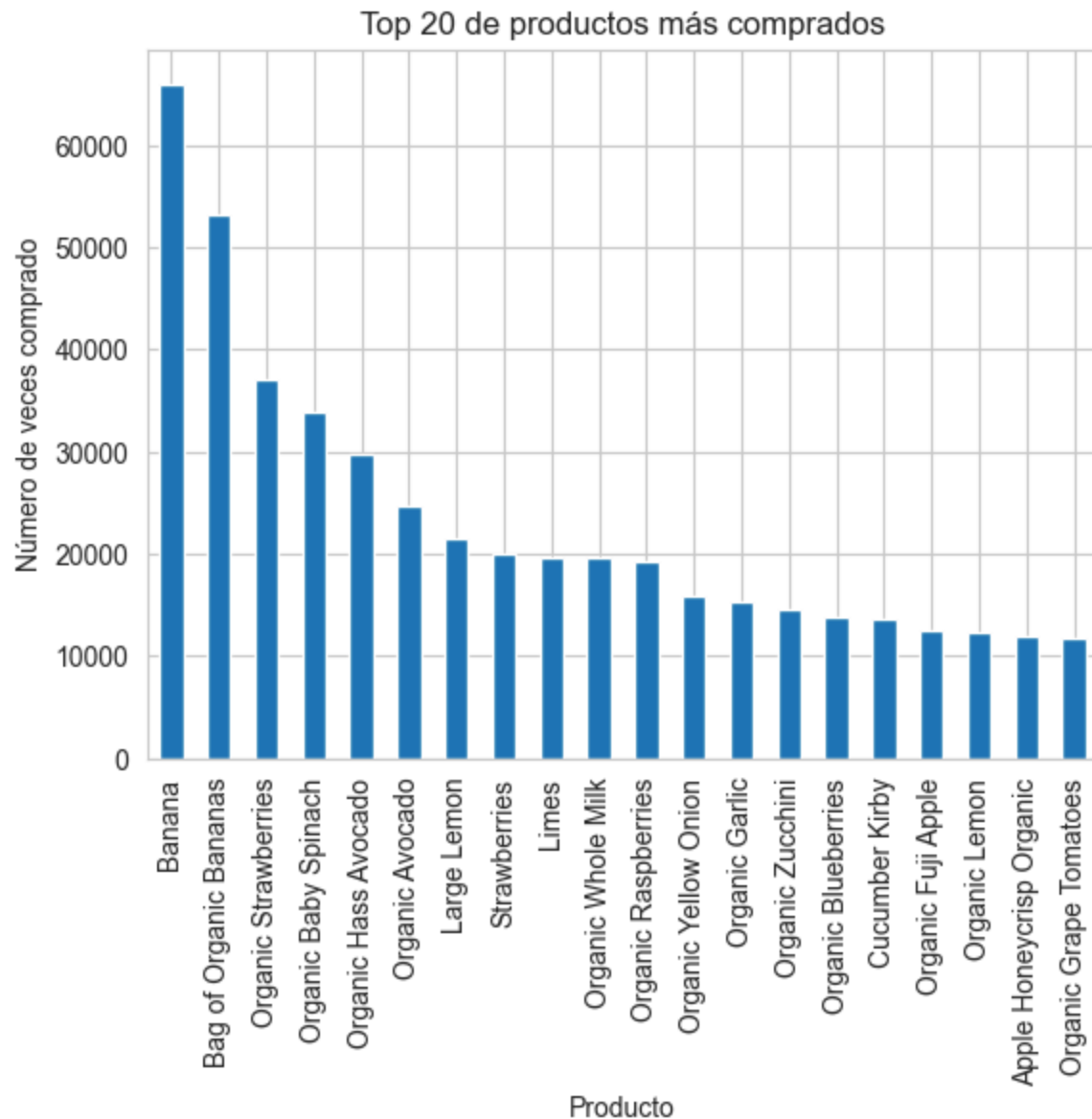
- En promedio, un cliente realiza **3.04 órdenes**. Además, del **top 10 de clientes** con más órdenes y otro del **top 10 de clientes con menos órdenes**. Esto sugiere que hay una variabilidad significativa en el número de órdenes por cliente, con algunos clientes realizando muchas más órdenes que otros.



## 5. Top 20 productos más frecuentes

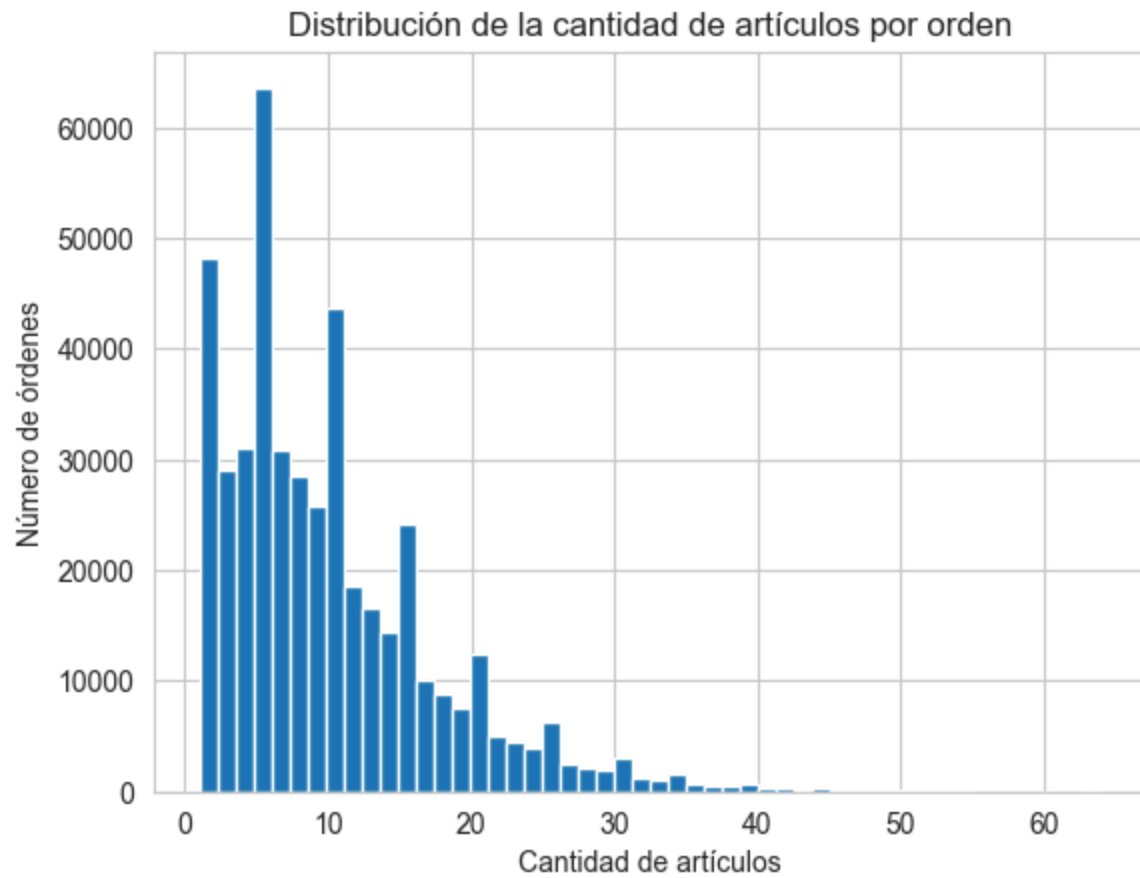
Los productos más populares son:





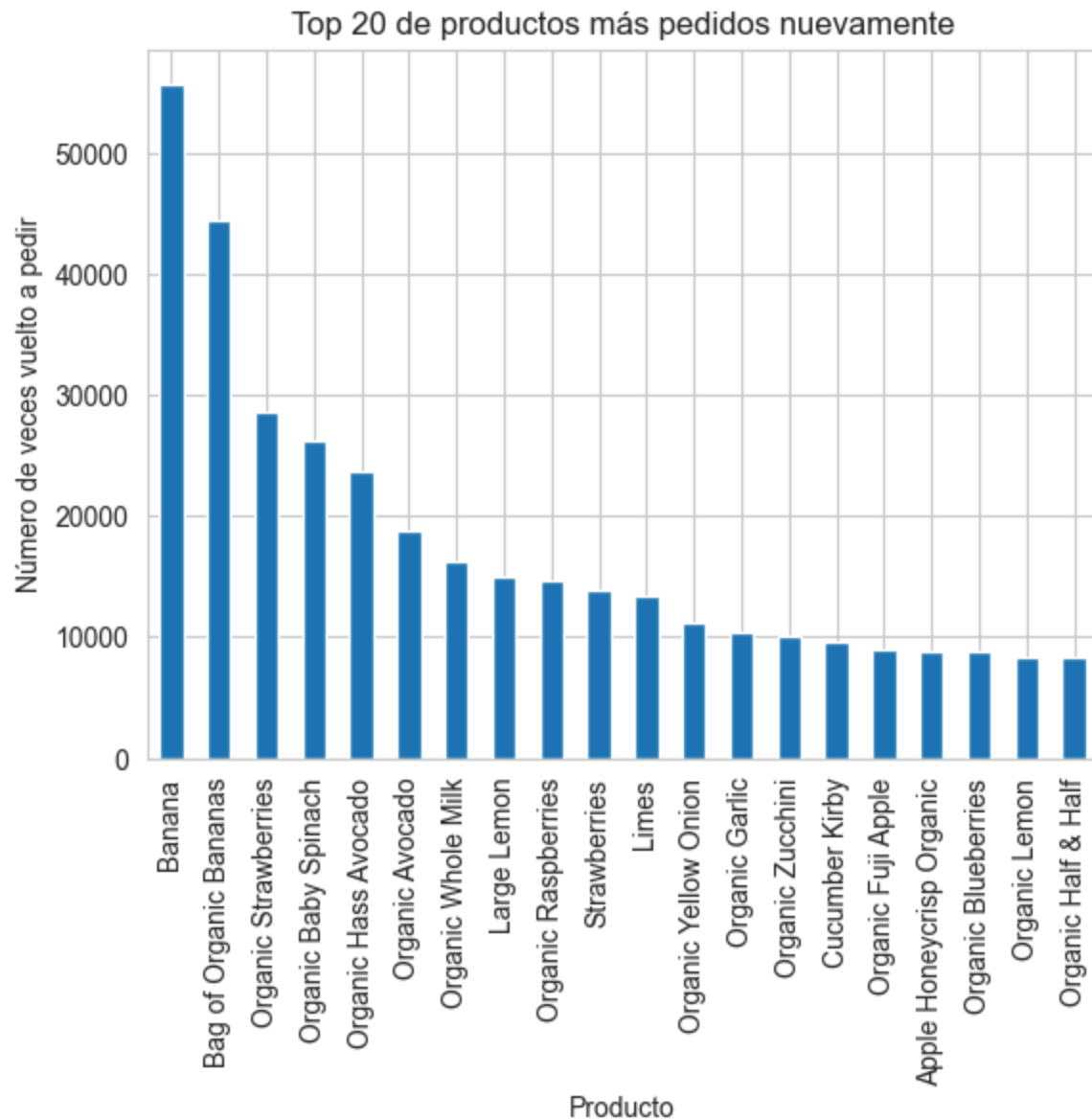
## 6. ¿Cuántos artículos se compran generalmente en un pedido?

- En promedio, se compran **10.1 artículos por orden**. La distribución de la cantidad de artículos por orden se muestra en un histograma, donde la mayoría de las órdenes tienen entre **5 y 15 artículos**.

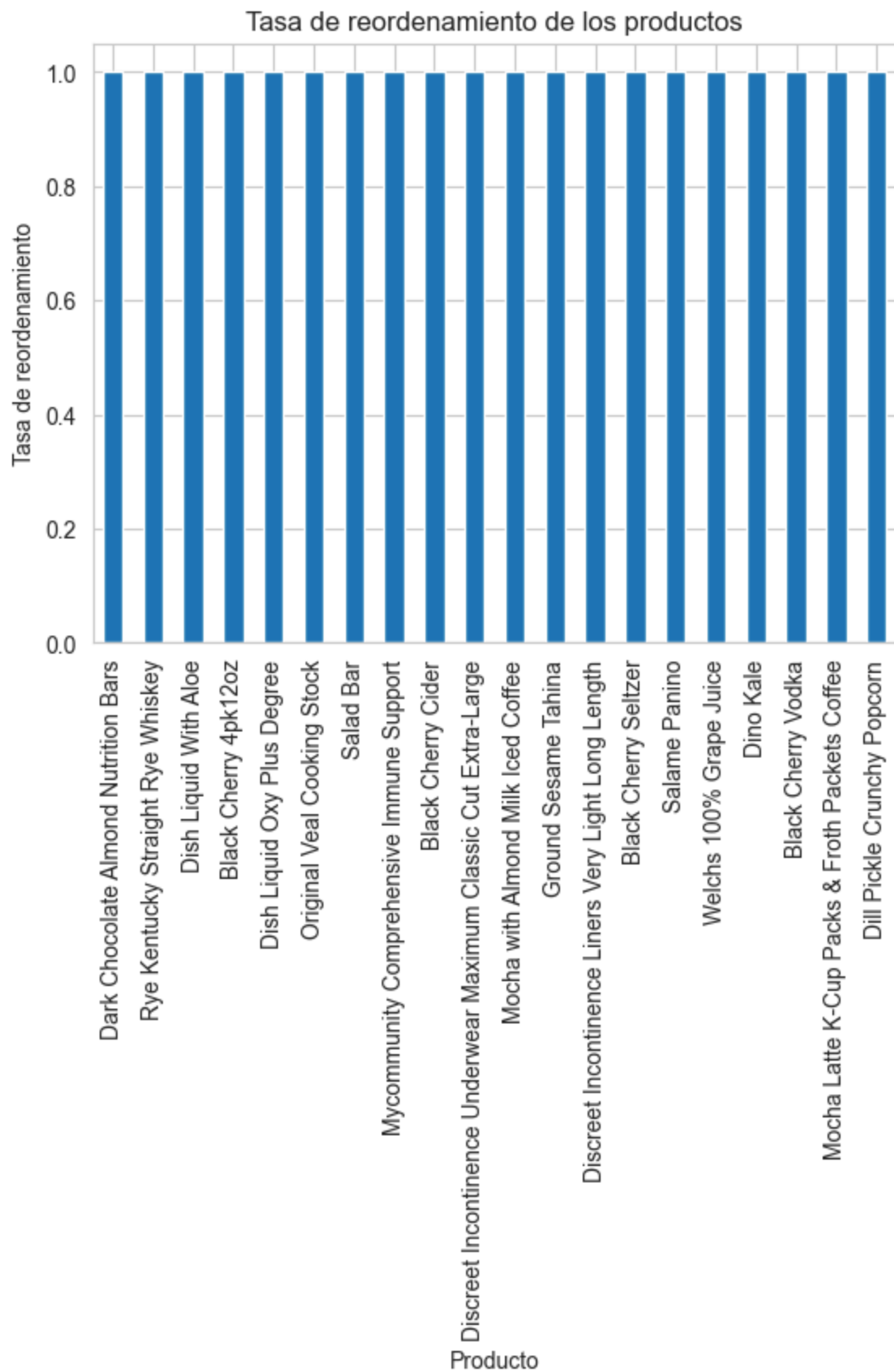


## 7. Top 20 artículos que se vuelven a pedir con más frecuencia

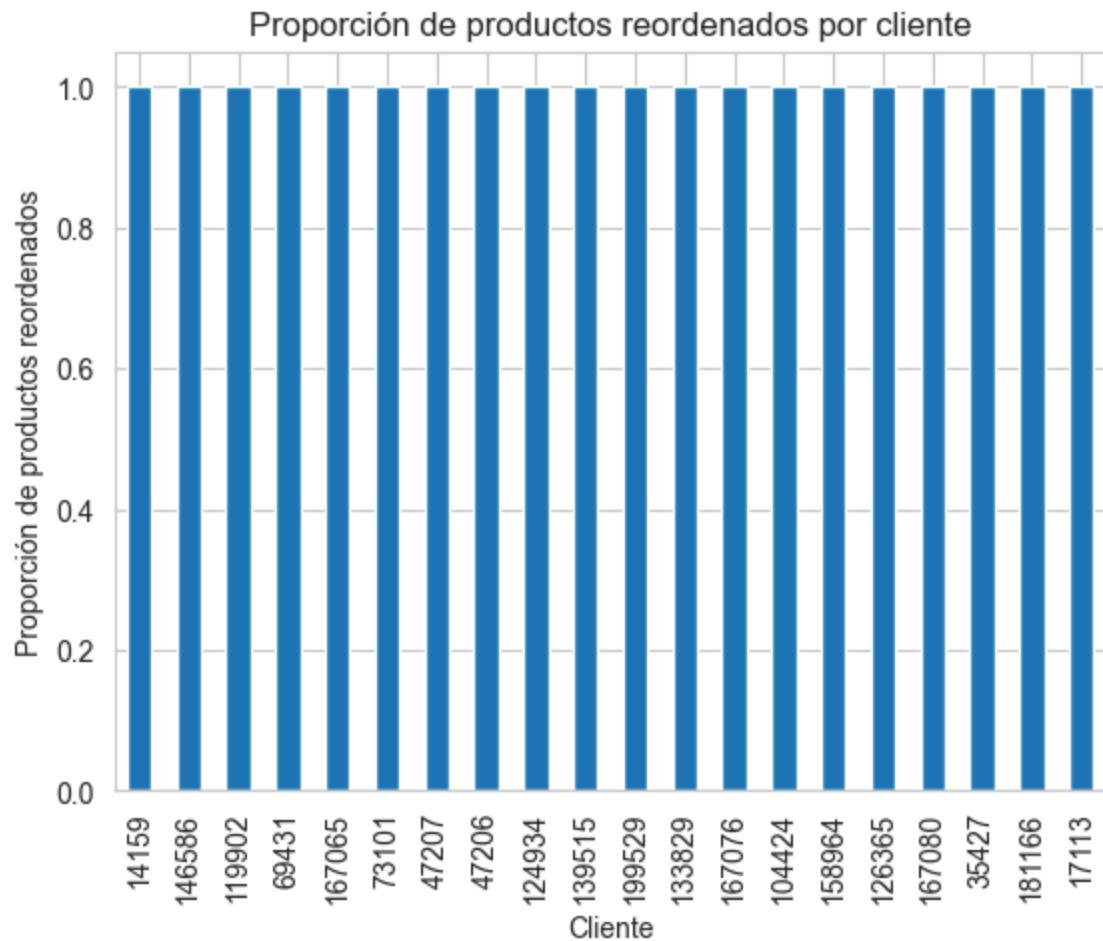
- **Productos más reordenados:**



## 8. Proporción de pedidos que se vuelven a pedir para cada producto

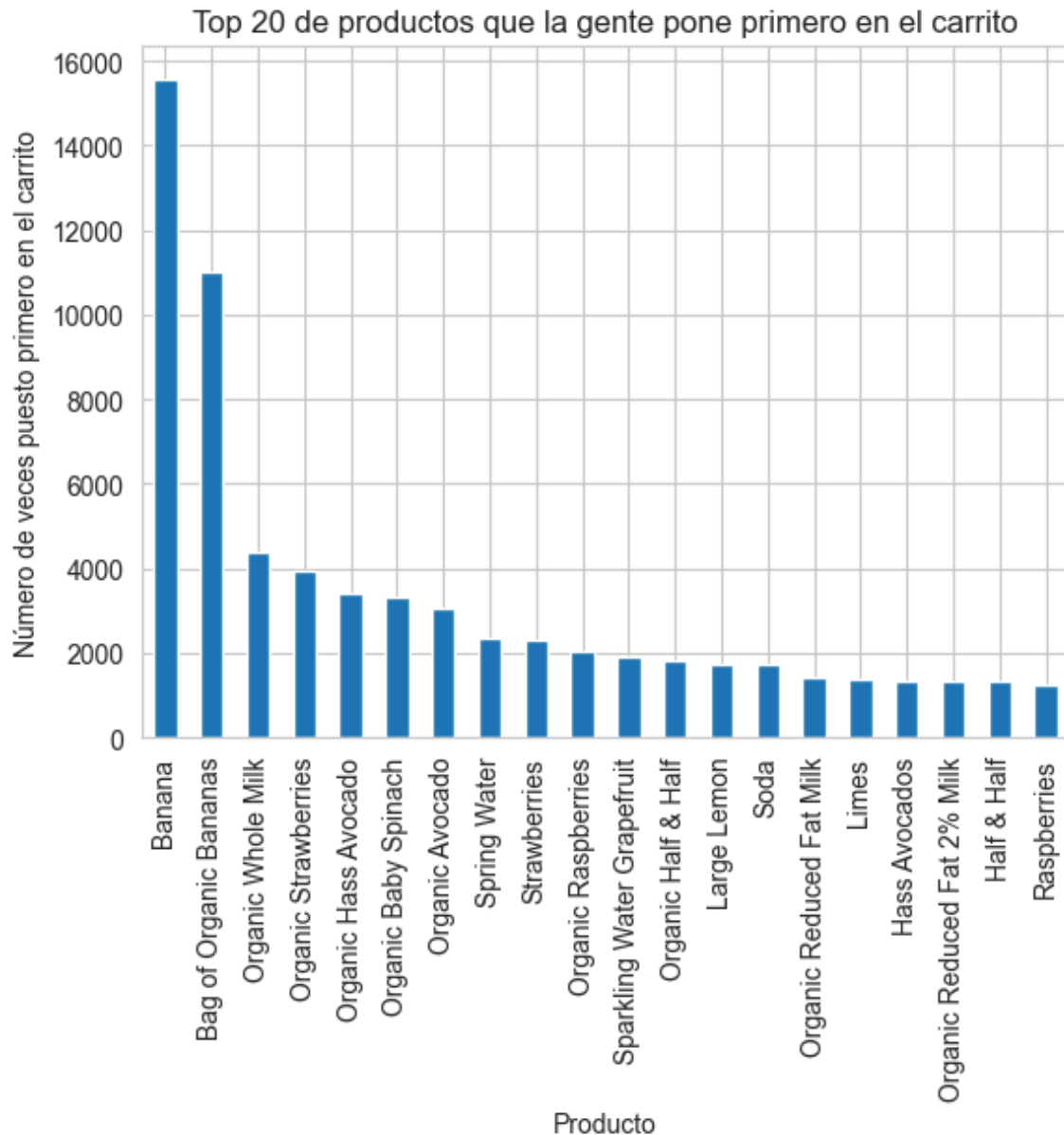


## 9. Proporción de productos pedidos que se vuelven a pedir para cada cliente



## 10. Top 20 artículos que la gente pone primero en el carrito

- Productos más comunes en primera posición:



## Conclusiones

Los fines de semana son los días con mayor actividad de compras, y los productos frescos y orgánicos son los más populares.

La mayoría de las compras se realizan en horas de la mañana, con un pico alrededor de las 10:00 a.m., y los clientes tienden a comprar en promedio 10 artículos por orden.

Los productos más reordenados son consistentes con los productos más

comprados, lo que sugiere una alta fidelidad hacia ciertos artículos. Los clientes muestran un comportamiento heterogéneo en cuanto al número de órdenes realizadas, con un promedio de 3.04 órdenes por cliente, pero con algunos clientes realizando significativamente más compras.

Repositorio Github

<https://github.com/AndresH1234/DataMining/tree/main/INSTACART>