



Máster en Data Science e Inteligencia Artificial

TRABAJO DE FIN DE MÁSTER

Role-Play DSMarket

Entrega Mayo 2025

Autores

Yeray Souto

Gaspar Manzini

Oscar Santos

Tutora

Raquel Revilla

Director

Isidre Royo

Curso académico Noviembre 2024

Índice

Introducción.....	3
Marco Teórico.....	4
Metodología.....	6
Limpieza de Datos.....	7
“Daily Calendar”.....	8
“Sales items”.....	9
“item prices”.....	10
Unión datasets.....	10
Preprocesado de Datos.....	12
Segmentación de Productos (Clusterización).....	13
Metodología.....	14
Carga de Librerías y Preparación del Dataset.....	14
Creación de Variables.....	15
Clustering por Categorías y por Productos.....	17
Resultados del Clustering Global de Productos.....	18
Clúster 1: Productos de alto valor.....	19
Clúster 2: Básicos Consistentes.....	20
Clúster 3: Productos de temporada.....	20
Clúster 4: Ventas Impredecibles.....	21
Clúster 5: Lenta Rotación.....	21
Clúster 6: Compras Planificadas.....	22
Clúster 7: Reposición Rápida.....	22
Planificación de campañas de Marketing basadas en Clusters.....	23
Interpretación de Resultados.....	24
Conclusiones y Recomendaciones.....	25
Modelado Predictivo con XGBoost para Forecasting de Ventas.....	27
Preparación Final de los Datos.....	27
Feature Engineering.....	27
Eliminación de Variables Altamente Correlacionadas.....	27
División del Dataset: Entrenamiento y Prueba.....	28
Entrenamiento del Modelo Predictivo con XGBoost.....	28
Resultados.....	31
Clusters Identificados y sus Características.....	32
Rendimiento del Modelo: MAE, RMSE y Comparativa por Cluster.....	34
Modelo en Producción.....	37
Métricas de Negocio.....	39
1. Reducción de excedentes de stock.....	40
2. Reducción de quiebres de stock.....	40
3. Mejora en la rotación de inventario.....	40

4. Impacto sobre las ventas.....	41
5. Ahorro en costos.....	41
6. Retorno sobre la inversión (ROI).....	41
Discusión.....	43
Conclusiones y Recomendaciones.....	45
Referencias.....	47

Introducción

Este trabajo de fin de máster surge con una motivación compartida: aplicar de manera integral los conocimientos adquiridos durante el Máster en Data Science a un caso real y desafiante. Como grupo, queríamos enfrentarnos a un problema que nos permitiera no solo poner en práctica técnicas de análisis y modelado, sino también **comprender un contexto de negocio** concreto y diseñar una solución que pudiera tener un impacto real. Así fue como nos encontramos con el caso de una cadena de tiendas con presencia en Boston, Filadelfia y Nueva York, que presentaba importantes dificultades en la gestión de su inventario.

El reto era claro: **mejorar la planificación del stock** para evitar tanto excesos como roturas, un problema muy común en el sector retail, pero que sigue representando un dolor de cabeza para muchas empresas. Para ello, trabajamos con un conjunto de datos históricos de ventas entre 2011 y 2016, proporcionado por **DSMarket**, con el objetivo de desarrollar un modelo capaz de predecir la necesidad de stock para las cuatro semanas de mayo de 2016.

El proyecto nos llevó a **recorrer todas las etapas** de un proceso completo de ciencia de datos. Comenzamos por una fase de exploración y **limpieza de los datos**, que nos permitió entender tanto la calidad de la información como las dinámicas de venta a lo largo del tiempo. A partir de ahí, realizamos **tareas de preprocesamiento**, ingeniería de variables y análisis exploratorio. Pronto nos dimos cuenta de que no era suficiente tratar todos los productos y tiendas por igual: existían patrones distintos según el tipo de producto, la ubicación o incluso la estacionalidad. Por eso, incorporamos una etapa de **clusterización**, que nos ayudó a segmentar los datos de forma más inteligente y a construir modelos más adaptados a cada realidad.

Para la predicción, optamos por utilizar el algoritmo **XGBoost**, que ya conocíamos por su buen rendimiento en tareas de regresión con datos relacionales. Este modelo nos permitió capturar relaciones complejas entre variables y lograr resultados satisfactorios en la estimación de la demanda semanal.

La memoria que presentamos está organizada siguiendo este recorrido. **Comenzamos con una introducción al problema y al contexto teórico**, para luego detallar cada fase del trabajo: desde el tratamiento de los datos hasta la construcción y evaluación del modelo. Finalmente, **compartimos una discusión sobre los resultados**, sus implicaciones prácticas y las limitaciones del enfoque, así como recomendaciones para una posible implementación real y para futuros trabajos.

En conjunto, este proyecto ha sido para nosotros mucho más que un ejercicio académico. Nos ha permitido consolidar nuestras habilidades técnicas, aprender a trabajar en equipo en un entorno de datos realista y, sobre todo, comprobar cómo la ciencia de datos puede convertirse en una herramienta poderosa para resolver problemas concretos y aportar valor en contextos empresariales.

Marco Teórico

El presente marco teórico se construye sobre los principios fundamentales de la ciencia de datos aplicados al contexto de la gestión de inventarios en tiendas retail.

En particular, se aborda una problemática común en el sector retail: **el equilibrio entre el desabastecimiento y el exceso de stock**, ambos factores que afectan de forma directa la rentabilidad y la eficiencia operativa de las cadenas comerciales.

A partir de este planteamiento, el objetivo general del proyecto consiste en **desarrollar un modelo predictivo** capaz de **estimar** el nivel de **inventario óptimo** para tiendas situadas en Boston, Filadelfia y Nueva York, utilizando como base datos históricos de ventas, precios y eventos comprendidos entre los años 2011 y 2016. La meta última es optimizar las decisiones de reabastecimiento durante el mes de mayo de 2016, minimizando al máximo posibles pérdidas por falta o exceso de producto.

El primer paso fue la **recopilación de datos**, centrada en extraer y consolidar información proporcionada por la plataforma DSMarket. Esta fuente incluía registros sobre las ventas diarias, precios por unidad y datos de calendario asociados a eventos especiales.

La integración de estas tablas heterogéneas implicó procesos de verificación de integridad, depuración de inconsistencias y su **consolidación en un repositorio único**, garantizando así la calidad y trazabilidad de los datos empleados en las fases posteriores.

El análisis exploratorio (**EDA**) permitió identificar patrones relevantes en los datos. A través del análisis de series temporales, se detectaron **comportamientos estacionales** y **tendencias semanales**, así como outliers en ventas diarias que fueron segmentados por ciudad y categoría de producto. Además, se generaron variables derivadas —como retardos temporales (lags), medias móviles y métricas de periodicidad— que enriquecieron el dataset y establecieron una base sólida para posteriores procesos de agrupamiento.

Una vez consolidada la información, se procedió al preprocesamiento de los datos, etapa crucial en cualquier estudio científico de los datos. Este proceso comprendió la limpieza de duplicados, imputación de valores faltantes mediante técnicas como forward, backfill, mediana, media siempre sobre el mismo id por último la estandarización de formatos temporales, incluyendo la creación de una variable clave: *"yearweek"*.

Dadas las **limitaciones operativas encontradas** en plataformas en la nube como Google Colab, el flujo de trabajo fue migrado a un **entorno local más robusto**, utilizando Visual Studio y en primera instancia, utilizamos la librería **Pandas** tal y como venimos ejecutando nuestro código, pero en última instancia debimos adoptar la **librería Polars**, la cual permitió un procesamiento significativamente más **eficiente de grandes volúmenes de datos**.

Con esta información, se aplicaron técnicas de clusterización, específicamente **K-means**, para agrupar productos según similitudes en sus patrones de ventas y

precios, tuvimos un pequeño bache a la hora de ejecutar el clusterizado de los datos, puesto que en primera instancia no tuvimos en cuenta el problema de la dimensionalidad. Una vez solucionado ese problema, seguimos el proceso de clusterización correctamente.

Para realizar el modelo predictivo sobre las cuatro semanas tuvimos muchos vaivenes, puesto que en primer lugar probamos a realizar un solo modelo con los datos sin tener en cuenta el cluster, al ver que no nos proporcionaba buena predicción decidimos testear la siguiente hipótesis: realizar un **modelo por categoría**. En este caso, la predicción mejoró en grandes rasgos, pero seguía sin arrojar unos buenos resultados cuando comparábamos el “df_test” con el “validation”. Por último, realizamos el **algoritmo basándonos en los clústeres** que habíamos obtenido. Esta segmentación nos permitió personalizar el modelado predictivo, para minimizar el error. Todos estos modelos los realizamos con XGBoost.

Para evaluar el rendimiento de los modelos, se emplearon métricas como el **Error Absoluto Medio** (MAE) y la **Raíz del Error Cuadrático Medio** (RMSE), comparando las predicciones con valores reales dentro de un esquema de validación cruzada. Este enfoque facilitó el ajuste de hiperparámetros y la selección de las configuraciones más eficientes de XGBoost en función de cada segmento.

Finalmente, la implementación del modelo consideró su integración práctica dentro del sistema de gestión de inventarios. Esto se incluyó en un mismo archivo a través de una “pipeline” para facilitar su integración y la reproducibilidad, la incorporación del servicio predictivo para automatizar las consultas de forecast y el despliegue de dashboards interactivos que permiten visualizar los pronósticos. Esta interfaz facilita la toma de decisiones estratégicas por parte del equipo operativo.

En conjunto, este marco teórico sintetiza los **fundamentos metodológicos y técnicos** que respaldan el desarrollo del proyecto, alineándose con las mejores prácticas en ciencia de datos e ingeniería aplicada al ámbito del retail, bajo un enfoque riguroso, reproducible y orientado a la optimización de procesos de negocio.

Metodología

La metodología seguida en este proyecto se estructura en cuatro fases principales: limpieza de datos, pre procesamiento, agrupación de productos (**clusterización**) y modelado predictivo (**forecasting**). A continuación, se detalla cada una de estas fases.

Descripción inicial de los datos

Los datos utilizados en este estudio provienen de DSMarket y abarcan el período histórico de 2011 a 2016. Las fuentes de datos principales incluyen:

1. **Sales Items:** Registros detallados de las ventas diarias por ítem y tienda.

```
items = df_sales_original['item'].unique()
print(items)

['ACCESORIES_1_001' 'ACCESORIES_1_002' 'ACCESORIES_1_003' ...
 'SUPERMARKET_3_825' 'SUPERMARKET_3_826' 'SUPERMARKET_3_827']

categories = df_sales_original['category'].unique()
print(categories)

['ACCESORIES' 'HOME & GARDEN' 'SUPERMARKET']

departments = df_sales_original['department'].unique()
print(departments)

['ACCESORIES_1' 'ACCESORIES_2' 'HOME & GARDEN_1' 'HOME & GARDEN_2'
 'SUPERMARKET_1' 'SUPERMARKET_2' 'SUPERMARKET_3']

stores = df_sales_original['store'].unique()
print(stores)

['Greenwich_Village' 'Harlem' 'Tribeca' 'Brooklyn' 'South_End' 'Roxbury'
 'Back_Bay' 'Midtown_Village' 'Yorktown' 'Queen_Village']

regions = df_sales_original['region'].unique()
print(regions)

['New York' 'Boston' 'Philadelphia']
```

```
df_sales_description.nunique()

id          30490
item        3049
category     3
department   7
store        10
store_code   10
region       3
dtype: int64
```

2. **Item Prices:** Información sobre los precios de venta de cada ítem a lo largo del tiempo.

```
[43] cant_semanas_año
... year
2011 52
2012 53
2013 52
2014 52
2015 52
2016 17
Name: week, dtype: object
```

```
df_prices.nunique()
✓ 0.6s
item      3049
category    3
store_code 10
yearweek   279
sell_price 1892
```

3. **Daily Calendar:** Datos sobre fechas, incluyendo información sobre eventos especiales o festividades que podrían influir en las ventas.

```
df_calendar.info()
[7]
... <class 'pandas.core.frame.DataFrame'>
RangeIndex: 1913 entries, 0 to 1912
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   date        1913 non-null  object
1   weekday     1913 non-null  object
2   weekday_int 1913 non-null  int64
3   d           1913 non-null  object
4   event       26 non-null    object
```

```
df_calendar.nunique()
✓ 0.0s
date      1913
weekday    7
weekday_int 7
d          1913
event      5
```

```
events
array(['SuperBowl', 'Ramadan starts', 'Thanksgiving', 'NewYear', 'Easter'],
```

Estos datasets presentan una alta granularidad y volumen, por ello hemos requerido de técnicas específicas para su unión y análisis

Limpieza de Datos

El desarrollo del proyecto se enmarcó dentro de un enfoque cuantitativo y exploratorio, fundamentado en la transformación y análisis sobre los datos provenientes de múltiples: *"item prices"*, *"sales items"* y *"daily calendar"*.

Durante las primeras etapas, se identificaron limitaciones en la infraestructura de Google Colab, particularmente en lo referente a la gestión de memoria, lo que ocasionaba reinicios constantes del entorno. Esta situación motivó la migración del flujo de trabajo a un entorno local, lo cual no solo brindó mayor estabilidad, sino que también permitió una gestión más eficiente. Paralelamente, gracias a la **formación recibida en el uso de GitHub**, se implementó una metodología de control de versiones que facilitó la sincronización del trabajo local con el repositorio remoto, **mejorando la organización y trazabilidad del desarrollo**.

El proceso de tratamiento de datos se estructuró en varias etapas interdependientes. En primer lugar, se realizó la migración del código desde Pandas a Polars, lo que implicó reescribir buena parte de la lógica original, debido a las **diferencias sintácticas y funcionales entre ambas librerías**. Se cargaron los dataframes, se ajustaron rutas de acceso y se optimizó la estructura de las tablas, preparando así el terreno para los análisis posteriores.

En cuanto a la limpieza y consolidación, nos encontramos con un problema principal, en los dataframes las fechas estaban expresadas de distintas formas por lo que tuvimos que encontrar la manera de unificar el criterio para poder unirlos. Para resolver esto utilizamos el dataset del *"daily calendar"* como nexo, agregando la columna clave *"yearweek"*, y para crear esta columna seguimos la lógica que tenía el dataset de *"item prices"* ya que el año **2012 tenía 53 semanas** mientras que los demás años tienen 52.

Una parte esencial del trabajo fue la reestructuración de los dataframes. En el caso de *"sales items"*, solo teníamos una línea por ID (que contiene el ítem y la tienda) y las ventas estaban en columnas, una por cada día, desde *"d_1"* hasta *"d_1920"*. Esta estructura fue **transformada mediante un proceso de transposición (o melt)**, en el que se redujo el número de columnas a 9, ya que se pasaron los días a filas, por lo que en vez de tener una línea por ID con todos los días pasamos a tener una por cada día, generando así cerca de **58 millones de filas**. También se agregaron las columnas de **fecha, eventos, yearweek y weekday**, uniendo este dataset con el *"daily calendar"* a través de el nombre de las columnas con formato *"d_"* que correspondían a los días.

En cuanto al dataset de *"item prices"*, nos encontramos con una inconsistencia en los precios de un *"item"* por lo que **actualizamos esos valores atípicos por la moda** para ese mismo *"item"* y *"store"*. En cuanto a los valores NaN que encontramos en este dataset luego de analizarlas vimos que se podían eliminar debido a que todas las *"yearweek"* de *"item sales"* estaban incluidas en las que tienen valor en *"item prices"*.

Una vez finalizada la preparación de los datos, se procedió con la aplicación de técnicas descriptivas y exploratorias, orientadas a **detectar patrones de comportamiento y tendencias generales**. Se evaluó la evolución de las ventas tanto a nivel semanal como mensual, se identificaron productos con caídas sostenidas en la demanda y se compararon rendimientos entre distintas tiendas y regiones. Además, se analizaron diferencias de precios con el objetivo de evaluar el impacto de promociones y

estrategias de optimización comercial.

A continuación, se sintetizan las principales características de los datasets tratados:

Dataset	Dimensiones Iniciales	Transformación Realizada	Observaciones
Item Prices	5 columnas; 6965.706 filas	Limpieza de fechas y consolidación en un único markdown	Corrección de valores NaN en "yearweek"
Sales Items	1.920 columnas; 30,490 filas	Transposición a 9 columnas; ~58 millones de filas generadas	Estructura más eficiente para análisis descriptivo
Daily Calendar	5 columnas; 1.913 filas	Agrupación semanal para facilitar integración	Simplificación en la unión de datasets

Transformación realizada

Se agrega la columna "sell price" al "df sales" desde "df prices", usando "item", "store code" y "yearweek" como claves.

```
df_sales_wprice = df_sales.merge(df_prices[['item','store_code', 'yearweek', 'sell_price']],  
on=['item','store_code', 'yearweek'], how='left')
```

“Daily Calendar”

Durante la fase de transformación temporal del dataset, se llevó a cabo la conversión de la columna "date" al tipo de dato "datetime", lo cual constituyó el primer paso necesario para la creación de una nueva variable temporal denominada "yearweek". Sin embargo, al aplicar fórmulas genéricas como "dt.strftime('%Y-%U')" y "dt.isocalendar()", se **observaron discrepancias significativas** entre los resultados obtenidos y los esperados. Estas inconsistencias se evidenciaron al intentar realizar la fusión entre los dataframes "sales items" y "item prices", en la que el campo "yearweek" es clave para la unión: el resultado arrojaba múltiples valores NaN, lo que indicaba una falta de correspondencia en la codificación temporal entre los datasets.

Al analizar con mayor profundidad el comportamiento de la variable "date" y su relación con el campo "weekday int" en el dataframe "daily calendar", se detectó que las semanas iniciaban en sábado, asignando el valor 1 a ese día específico. Este detalle resultó crucial, ya que la **lógica semanal no coincidía con el estándar ISO**, generando un desfase temporal. Además, se constató que, si bien la mayoría de los años en "items price" cuentan con exactamente 52 semanas, el año 2012 presenta una semana adicional (53), lo que sugiere un modelo particular de segmentación: al completarse 52

semanas, el **sistema asume automáticamente el inicio del nuevo año**, sin considerar que aún puedan quedar días dentro del calendario del año anterior.

A partir de esta observación, **se optó por implementar una fórmula personalizada** para la generación de la variable *“yearweek”*, que respetara el inicio de semana en sábado y limitará el conteo a 52 semanas anuales, salvo en el caso del 2012, donde se contempló una semana 53. Esta lógica permitió corregir la secuencia temporal y alinear adecuadamente los datos entre los distintos datasets, eliminando así los valores faltantes producidos durante el merge.

Adicionalmente, se procedió a limpiar la columna *“event”* en el calendario, sustituyendo los valores nulos por la etiqueta *“Regular Day”* para evitar ambigüedades en el análisis posterior. También se enriqueció este campo con eventos relevantes como *“Navidad”* y *“4 de Julio”*, considerados fundamentales por su impacto en los patrones de consumo.

Finalmente, se construyeron diccionarios a partir de la columna *“d”*, que incluían las variables *“date”*, *“weekday”*, *“event”* y la nueva *“yearweek”*, con el propósito de integrar esta información enriquecida

directamente en el dataframe *“sales items”*. Este procedimiento **consolidó la coherencia temporal y contextual del conjunto de datos**, sentando las bases para una modelización precisa y robusta.

“Sales items”

Como parte del proceso de transformación estructural del conjunto de datos *“sales items”*, se llevó a cabo una operación de *“melt”* sobre las columnas que iniciaban con el prefijo *“d_”*, convirtiéndolas de un formato de columnas a un formato largo en filas. Durante esta transformación, se mantuvieron fijas las columnas clave que abarcan desde *“id”* hasta *“región”*, preservando así la identidad de cada producto y su ubicación dentro del conjunto. Esta reorganización permitió una estructura más flexible y adecuada para los análisis posteriores..

Inicialmente, se había optado por eliminar las observaciones con valor cero en la variable de ventas, bajo la suposición de que no aportaban información significativa al análisis. Sin embargo, tras una revisión, se consideró que estas **líneas**

```
df_calendar.info()
✓ 0.0s

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1913 entries, 0 to 1912
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  -
0   date            1913 non-null   object
1   weekday         1913 non-null   object
2   weekday_int     1913 non-null   int64
3   d              1913 non-null   object
4   event           26 non-null     object
dtypes: int64(1), object(4)
memory usage: 74.9+ KB
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 58327370 entries, 0 to 58327369
Data columns (total 13 columns):
#   Column          Dtype
---  -
0   id              object
1   item            object
2   category        object
3   department      object
4   store           object
5   store_code      object
6   region          object
7   Fecha           object
8   Qsale           int64
9   date            datetime64[ns]
10  yearweek        object
11  event           object
12  weekday         int64
dtypes: datetime64[ns](1), int64(2), object(10)
memory usage: 5.6+ GB
```

resultan relevantes para ciertos análisis exploratorios previos a la fase de clusterización, especialmente al momento de evaluar la **estacionalidad** o el comportamiento de productos de **baja rotación**. Por este motivo, se reincorporaron los registros con ventas igual a cero, asegurando una representación más fiel de la dinámica real del mercado.

Posteriormente, se procedió al reseteo del índice del dataframe, con el objetivo de reorganizar la numeración de las filas tras las múltiples transformaciones aplicadas. Esta acción facilitó un manejo más limpio del conjunto de datos, especialmente para efectos de iteraciones posteriores y trazabilidad de registros.

Finalmente, se integraron al dataset las columnas *"date"*, *"yearweek"*, *"event"* y *"weekday"*, utilizando los diccionarios construidos previamente a partir del dataset *"daily calendar"*. Esta fusión de información permitió enriquecer cada transacción con contexto temporal y eventos relevantes, reforzando así la calidad del dataset para las siguientes etapas de análisis predictivo y modelado.

"item prices"

En el marco de la depuración final del campo *"yearweek"*, se identificaron valores nulos (NaN) que fueron actualizados mediante la fórmula personalizada de la que hemos hablado anteriormente, diseñada para mantener la coherencia temporal previamente definida en el proyecto. Esta lógica consideró el inicio de semana en sábado y la limitación de 52 semanas por año, salvo en los casos excepcionales como 2012. Una vez completada esta imputación, se procedió a estandarizar el formato de dicha columna, eliminando el sufijo decimal ".0" que se generaba durante ciertas operaciones de conversión numérica. Este ajuste garantizó una codificación uniforme y adecuada para posteriores procesos de agrupación y análisis.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6965706 entries, 0 to 6965705
Data columns (total 5 columns):
#   Column      Dtype
---  -
0   item        object
1   category    object
2   store_code  object
3   yearweek    object
4   sell_price  float64
dtypes: float64(1), object(4)
memory usage: 265.7+ MB
```

Unión datasets

Durante la fase de validación del dataset, se verificó la integridad de los valores correspondientes a la variable *"sell price"* dentro del dataframe de ventas (*"df sale"*). El objetivo principal de esta revisión era **comprobar la correcta integración de los precios con las fechas** de venta mediante la variable derivada *"yearweek"*.

En este proceso, se identificaron valores nulos en *"sell price"*, lo cual indicaba posibles inconsistencias en la fórmula utilizada para generar la variable *"yearweek"*. Ante esta situación, se analizó si existía algún patrón recurrente asociado a los días con errores. El análisis reveló que la mayoría de los registros con valores nulos se concentraban en los

últimos días del mes de diciembre, repetidamente a lo largo de varios años del histórico.

Una inspección más detallada permitió observar que estas líneas sin precio de venta correspondían, en todos los casos, a **días sin transacciones registradas**. O en otras palabras, no se realizó **ninguna venta ese día**.

Esto sugería que la ausencia de precios no era un error de origen, sino una consecuencia lógica de la falta de actividad comercial en esas fechas específicas.

A partir de este hallazgo, se desarrolló una estrategia de imputación para completar los valores faltantes en "*sell price*". Se diseñó una fórmula que, en caso de detectar un valor nulo, tomara como referencia el último precio disponible previo. Si no existía un valor anterior, se optaba por utilizar el valor posterior más próximo. Esta lógica de relleno secuencial "*forward fill*" "*backfill*" se implementó con éxito.

Finalmente, se verificó que, tras aplicar esta estrategia, no quedaban valores nulos en la variable "*sell price*", confirmando así la efectividad del procedimiento y restableciendo la integridad del conjunto de datos.

```
df_sales_wprice.info()
✓ 0.0s

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 58327370 entries, 0 to 58327369
Data columns (total 12 columns):
#   Column      Dtype
---  -
0   item        object
1   category    object
2   department  object
3   store       object
4   city        object
5   Qsale       int64
6   date        datetime64[ns]
7   yearweek    object
8   event       object
9   weekday     int64
10  unit_price  float64
11  total_venta float64
dtypes: datetime64[ns](1), float64(2), int64(1)
memory usage: 5.2+ GB
```

Pre procesamiento de Datos

Una vez completada la fase de limpieza y depuración de datos, se procedió con el preprocesamiento del conjunto, etapa fundamental para preparar la información de cara a las fases de clusterización y modelado predictivo. El objetivo de este proceso fue transformar los datos en un **formato estructurado, coherente y enriquecido**, que permitiera capturar de manera efectiva los patrones subyacentes de comportamiento de sobre los 3049 productos que tenemos en el dataset.

En primer lugar, se aplicaron técnicas de imputación adicionales para abordar los valores faltantes en campos críticos, como la columna de precios "*sell price*". Para ello, se implementaron métodos de "*backfill*", "*forward fill*" y en caso de que fuera posible realizar un "*interpolate*", para **promediar el valor de estos registros vacíos**. El uso del "*backfill*" y "*forward fill*" únicamente lo realizamos cuando era el último registro, ya que no teníamos datos para realizar el "*interpolate*". **Todos los registros tenían como referencia su "id"** para no tener problemas a la hora de imputar este tipo de columnas, finalmente completamos los registros (NaN) respetando la continuidad y tendencia de los precios registrados, garantizando así la integridad del dataset sin introducir sesgos significativos.

Con los datos ya consolidados, se realizó un análisis exploratorio exhaustivo (EDA) utilizando herramientas como Python y Power BI. Esta fase fue clave para verificar que la limpieza de los datos había sido efectiva y que las variables seguían una distribución adecuada. Además, **permitió identificar valores atípicos potenciales y analizar la variabilidad de los precios a lo largo del tiempo**. Paralelamente, se exploraron patrones de comportamiento por ciudad, lo que resultó fundamental para comprender las dinámicas locales del mercado y orientar con mayor precisión los enfoques de segmentación posteriores.

Posteriormente, se desarrollaron diversas variables derivadas mediante procesos de **feature engineering**, con el fin de capturar dimensiones temporales, comerciales y contextuales relevantes. Se incluyeron **variables temporales como el año, el mes, la semana y el trimestre**, así como indicadores binarios para identificar semanas con eventos especiales, los cuales suelen tener un impacto significativo en los volúmenes de venta. Asimismo, se creó un identificador único para cada combinación de producto y tienda "*item store id*", facilitando el seguimiento individualizado de cada unidad de análisis a lo largo del tiempo.

Finalmente, se procedió a la **codificación de las variables categóricas**, transformando atributos como tienda, ciudad, categoría, departamento e ítem en formatos numéricos mediante técnicas como "*Ordinal Encoder*", con el fin de adaptarlos a los requerimientos de los algoritmos de aprendizaje automático. Las variables booleanas también fueron adaptadas, utilizando representaciones binarias "*0 o 1*".

Este preprocesamiento integral no solo permitió estandarizar y enriquecer el conjunto de datos, sino que sentó las bases metodológicas para las etapas analíticas posteriores, garantizando robustez y coherencia en todo el flujo de trabajo de ciencia de datos.

Clusterización base de datos

La caracterización y clusterización de los productos se realiza con el objetivo de **optimizar las estrategias de comercialización y comunicación**. Se plantea la problemática de identificar agrupaciones de productos que permitan establecer campañas de marketing segmentadas y mejorar la experiencia del cliente mediante la personalización de ofertas. Además, se pretende analizar la efectividad predictiva del modelo en términos de revenue, ventas y rotación de inventarios, integrando variables monetarias, de frecuencia, intensidad y cambios en precios, así como la influencia de eventos especiales en el comportamiento de ventas.

Los objetivos específicos de la investigación son:

- **Definir las variables** relevantes para la caracterización de los productos, tales como revenue total, ventas diarias, semanales, mensuales y trimestrales, precios promedio, cambios porcentuales y frecuencia de ventas.
- Implementar **técnicas de ingeniería de variables** utilizando librerías de Python (NumPy, Pandas, Matplotlib, entre otras) para optimizar la calidad de los datos.
- Aplicar **algoritmos de clustering** que permitan segmentar los productos en función de criterios preestablecidos (categorías y comportamiento de venta), y evaluar la distribución y consistencia de los clústeres obtenidos, utilizando métodos como el elbow method para determinar el número óptimo de grupos.
- **Formular recomendaciones** para el diseño de campañas de marketing dirigidas basadas en la segmentación de productos y la interpretación de los resultados de los clústeres.

Asimismo, la incorporación de variables derivadas como indicadores de precio elevado, promedios de venta por ítem, intensidad y frecuencia de compra, y métricas de rotación de inventario enriquece considerablemente la capacidad de predicción de los modelos. A través del cálculo estadístico como medias, máximos, mínimos y porcentajes de cambios, es posible evaluar la volatilidad en precios y la dinámica de stock con mayor precisión (*James, Witten, Hastie, & Tibshirani, 2013*). Complementariamente, el diseño de variables basadas en eventos, que consideran el impacto de fechas especiales dentro de ventanas temporales específicas, fortalece la capacidad del modelo para anticipar picos o caídas en la demanda (Davenport & Harris, 2017).

Metodología sobre clusterización

Carga de Librerías y Preparación del Dataset

Este bloque de código establece el entorno base necesario para el desarrollo de un pipeline completo de ciencia de datos, desde la gestión de datos y visualización hasta la construcción, evaluación e interpretación de modelos predictivos. Se incluyen librerías esenciales para cada fase del flujo de trabajo:

- **Gestión y transformación de datos:** `Pandas`, `Polars`, `Numpy`, `Math`, y `Pickle`, fundamentales para la manipulación de estructuras de datos, almacenamiento y exportación.
- **Visualización:** Herramientas como `Matplotlib`, `Seaborn`, `Yellowbrick` y `Missingno` permiten explorar y visualizar datos.
- **Preprocesamiento:** Incluye una gama completa de herramientas de `sklearn` para escalado, imputación, codificación y construcción de pipelines modulares ("`Pipeline`", "`ColumnTransformer`"). Se destacan técnicas de imputación como "`KNNImputer`" y "`SimpleImputer`", junto con codificadores como "`OrdinalEncoder`" y "`OneHotEncoder`".
- **Interpretabilidad:** La biblioteca `Shap` se incluye para análisis de interpretabilidad de modelos complejos, permitiendo explicar el impacto de cada variable en las predicciones.
- Se emplean librerías como **Yellowbrick** concretamente la sub librería de "`cluster`" "`KElbowVisualizer`", "`SilhouetteVisualizer`" para calcular el número de codos óptimo basándonos en los datos entregados.
- Por otra parte, para simple visualizar un notebook más limpio y por TOC de algunos miembros del grupo, requerimos de la librería **warnings** para evitar mensajes sobre códigos depreciados.

Además, se configuran parámetros globales como la visualización completa de columnas en `pandas`, la supresión de advertencias para mantener limpio el entorno de ejecución, y la fijación de la semilla de aleatoriedad (`np.random.seed(42)`) para asegurar la reproducibilidad.

Tras una copia del DataFrame original, se realiza un análisis exploratorio de datos (EDA) para confirmar la integridad de la información, verificando la existencia de ventas en todas las fechas y **eliminando columnas irrelevantes como "`Store`", "`weekday`", "`total venta`", "`City`",** ya que el análisis se centra en los ítems independientemente de la tienda.

Creación de Variables

En una primera iteración del proceso de ingeniería de características, se generó un conjunto extenso de variables, **alcanzando aproximadamente cincuenta variables** sin haber considerado inicialmente los efectos de la alta dimensionalidad sobre el rendimiento del modelo. Esta decisión, aunque útil para capturar múltiples aspectos del comportamiento de los datos, **introdujo una complejidad adicional** que posteriormente requeriría técnicas de selección o reducción de variables para mitigar el riesgo de sobre ajuste. Con el fin de **mejorar la calidad del análisis y la precisión del modelado**. Estas variables fueron diseñadas para capturar tanto comportamientos agregados como variaciones temporales relevantes.

- En primer lugar, se calcularon **indicadores asociados al rendimiento económico** total y por intervalo temporal. La variable *"total_revenue"* representa los ingresos acumulados por producto, mientras que *"average_daily_revenue"*, *"average_weekly_revenue"* y *"average_monthly_revenue"* permiten evaluar su comportamiento promedio en diferentes escalas temporales. A su vez, se determinaron los valores máximos alcanzados en cada uno de estos periodos *"max_daily_revenue"*, *"max_weekly_revenue"*, *"max_monthly_revenue"*, con el objetivo de detectar picos de demanda o momentos de alta rentabilidad.
- En paralelo, se incorporaron **métricas relacionadas con el precio**, tales como *"average_price"* y *"std_average_price"* (desviación estándar del precio promedio), esta última especialmente útil para identificar productos de comportamiento estacional o con variabilidad en su posicionamiento comercial. También se calcularon indicadores de variación de precios inter semanales, como *"max_price_change"* y *"min_price_change"*, que expresan el mayor y menor porcentaje de cambio respecto a la semana anterior, proporcionando una señal de sensibilidad al precio.
- Adicionalmente, se **introdujeron variables binarias (flags)** que permiten identificar comportamientos excepcionales a nivel de departamento. Por ejemplo, *"flag_high_price"* toma el valor 1 si el precio del producto está por encima del promedio de su departamento, mientras que *"flag_high_sells"* y *"flag_high_revenue"* cumplen funciones equivalentes en términos de unidades vendidas e ingresos totales, respectivamente. Estas banderas **permiten detectar productos destacados** dentro de su categoría, útiles tanto para promociones como para estrategias de reposicionamiento.
- Finalmente, la variable *"total_units_sold"* consolida la cantidad total de unidades vendidas por producto, sirviendo como una métrica **fundamental en los análisis de rotación, popularidad y segmentación por volumen de ventas**.

Con el objetivo de enriquecer la comprensión de los patrones de comportamiento asociados a cada producto, se desarrolló un conjunto adicional de variables centradas en el análisis temporal de las unidades vendidas. Estas métricas permiten capturar tanto la intensidad como la regularidad de las ventas, y ofrecen señales clave para identificar productos con comportamientos estacionales, ventas concentradas o distribuciones irregulares.

En cuanto al volumen de ventas, se calcularon los **promedios de unidades vendidas en diferentes escalas temporales**: diaria *"average_daily_units_sold"*, semanal *"average_weekly_units_sold"*, mensual *"average_monthly_units_sold"* y trimestral *"average_trimestral_units_sold"*. Estas métricas se complementaron con indicadores de máximos y mínimos para cada una de estas escalas, incluyendo registros diarios, semanales y mensuales, con el objetivo de identificar tanto picos de demanda como periodos de baja rotación.

Para entender la **recurrencia en las ventas**, se incorporaron variables como *"max_days_from_last_sale"*, *"min_days_from_last_sale"* y *"days_from_last_sale_mean"*, las cuales permiten medir la frecuencia con la que un producto vuelve a venderse. En esta misma línea, el flag *"flag_sells_every_month"* señala si el producto mantiene ventas activas todos los meses del año, ofreciendo una primera señal sobre su comportamiento estacional.

Se desarrolló también una serie de indicadores específicos sobre el momento del mes o de la semana en el que se concentran las ventas. Estos flags incluyen:

- *"flag_weekend"*, que indica si las ventas durante fines de semana superan el promedio semanal;
- *"flag_half_week"*, *"flag_week_start"*, que detectan si días entre miércoles-jueves o lunes-martes, respectivamente, superan dicha media;
- *"flag_month_start"*, *"flag_half_month"*, *"flag_month_end"*, diseñados para detectar concentración de ventas en los primeros, últimos o días centrales del mes.
En el plano trimestral, se incorporaron indicadores como *"flag_1st_trimester"*, *"flag_2nd_trimester"*, etc., que identifican si en alguno de los cuatro trimestres las ventas superan la media trimestral.
- Finalmente, se incluyen variables que reflejan el impacto de eventos especiales, como *"event_boosted"* y *"event_decreased"*, que contabilizan la cantidad de ocasiones en que un producto experimentó aumentos o caídas de ventas asociadas a semanas con eventos específicos.

Este análisis considera tanto la estacionalidad directa como picos en torno a eventos específicos

Durante la fase de entrenamiento del modelo de clustering, se identificó un problema de alta dimensionalidad derivado de la inclusión inicial de **más de cincuenta variables**. Este exceso de atributos no solo aumentaba la complejidad computacional del modelo, sino que también generaba redundancias que podían distorsionar la segmentación y afectar negativamente la interpretación de los resultados.

Para abordar este problema, se llevó a cabo un proceso de depuración basado en el **análisis de correlación entre variables**. Se elaboró una matriz de correlación para identificar aquellas variables altamente correlacionadas entre sí, así como aquellas que resultaban autoexplicativas —es decir, cuya información podía inferirse directamente de otras variables ya presentes en el conjunto. Aquellos atributos que no aportaban

información adicional significativa o que presentaban colinealidad evidente fueron descartados.

El objetivo de este filtrado fue **conservar únicamente las variables más representativas** y no redundantes, optimizando así la calidad de la segmentación y reduciendo el **riesgo de sobre ajuste**. Como resultado de este proceso, se definió un conjunto final de 23 variables que capturan de manera robusta los principales patrones de comportamiento comercial, dinámicas de precios, frecuencia de ventas y respuesta a eventos.

Las variables finales utilizadas para entrenar el modelo de clustering fueron:

"average_weekly_revenue", "flag_high_price", "flag_high_sells", "flag_high_revenue", "average_price", "max_price_change", "price_changes_total", "total_units_sold", "average_weekly_units_sold", "max_days_from_last_sale", "days_from_last_sale_mean", "sells_every_week", "sells_every_month", "flag_middle_week_item", "flag_end_week_item", "flag_middle_month_item", "flag_end_month_item", "flag_middle_quarter_item", "flag_end_quarter_item", "flag_middle_year_item", "flag_end_year_item", "event_boosted", y "event_decreased".

Finalmente, se incluyen variables que reflejan el impacto de eventos especiales, como *"event_boosted"* y *"event_decreased"*, que contabilizan la cantidad de ocasiones en que un producto experimentó aumentos o caídas de ventas asociadas a semanas con eventos específicos.

Este DataFrame representa un punto de partida robusto para tareas de modelado, clustering o segmentación, ya que condensa información clave de comportamiento histórico de cada producto con un alto nivel de granularidad y completitud.

Clustering por Categorías y por Productos

- Clustering basado en categorías

El análisis de segmentación mediante clustering se abordó inicialmente dividiendo y realizando un cluster por cada categoría, donde los productos fueron agrupados por **Supermarket, Accessories y Home & Garden**. Para cada una de estas categorías se aplicó el método del codo (*Elbow Method*) con el fin de determinar el número óptimo de clústeres, utilizando métricas como la inercia intra-cluster para guiar la decisión.

En el caso de la categoría *"Supermarket"*, el modelo identificó **7 clústeres**, uno de los cuales contenía un conjunto reducido de 10 productos que se comportaban como outliers. Dentro de esta segmentación, se observaron agrupaciones con valores significativamente altos de ingresos diarios (*"Daily Revenue"*). Para las categorías *"Accessories"* y *"Home & Garden"*, también se obtuvieron 7 y 6 clústeres respectivamente. No obstante, ambas presentaron cierta dispersión interna y la aparición de agrupaciones pequeñas, como un cluster de 66 productos en *"Home & Garden"*, lo cual sugería la presencia de subgrupos atípicos o poco homogéneos.

Analizados los resultados y viendo que el **número de observaciones por cluster no se distribuía de una manera óptima**, los resultados obtenidos en el modelo de forecasting no mostraron una mejora sustancial al segmentar los productos por categoría. Esto llevó a reconsiderar la estrategia de agrupamiento, optando por explorar una alternativa más holística: aplicar clustering directamente sobre todo el

conjunto de productos, sin segmentación previa por categoría, con el objetivo de capturar patrones comunes a nivel global, independientemente del tipo de producto.

- **Clustering basado en productos:**

A partir de la nueva estrategia de segmentación, se aplicó el algoritmo de clustering sobre la totalidad de los **3,049 productos del catálogo**, sin segmentación previa por categoría. Como resultado, se identificaron 7 clústeres, considerados óptimos tras aplicar métodos de validación como el análisis de inercia y visualizaciones de silueta. Aunque la **distribución ideal** esperada era de aproximadamente **435 productos** por grupo, se observaron variaciones naturales en la asignación, con **clústeres que oscilaron entre 302 y 748 productos**. Esta dispersión refleja la diversidad inherente en los patrones de comportamiento comercial de los productos.

El análisis posterior de las características predominantes en cada grupo permitió clasificar los clústeres en tipologías operativas con significado estratégico. Entre los segmentos identificados se encuentran: **productos de alto valor**, caracterizados por ingresos elevados; **productos básicos y consistentes**, con ventas estables y recurrentes; **productos de temporada**, con picos específicos en momentos del año; **ventas impredecibles**, con comportamiento errático o sensible a eventos; **productos de lenta rotación**, con bajos volúmenes de venta sostenida; **compras planificadas**, vinculadas a ciclos de reposición largos; y **productos de reposición rápida**, cuya demanda exige presencia constante en inventario.

La ejecución del algoritmo se acompañó de visualizaciones descriptivas y tablas resumen, que facilitaron tanto la interpretación de los clústeres como su validación empírica. Esta segmentación proporciona una **base robusta para definir estrategias** diferenciadas de inventario, marketing y forecasting, adaptadas a los patrones reales observados en los datos.

Resultados del Clustering de Productos

La segmentación de productos mediante clústeres permite identificar patrones de comportamiento de ventas comunes entre distintos artículos. Esta agrupación no solo facilita una mejor comprensión de la dinámica de consumo, sino que habilita decisiones más eficientes en áreas como gestión de stock, diseño de promociones y planificación comercial. En este análisis se presentan siete clústeres bien diferenciados de productos, incluyendo categorías como accesorios y hogar, con un enfoque en ventas, rotación, precios y estacionalidad.

La distribución ideal esperada era de aproximadamente 435 productos por grupo, se observaron variaciones naturales en la asignación, con clústeres que oscilaron entre 302 y 748 productos. Para cada uno se han definido nombres comerciales, perfiles de comportamiento y acciones tácticas recomendadas, especialmente pensadas para su aplicación en tiendas minoristas.

Cluster	Número de Productos	Características Destacadas
1	383	Productos de alto valor, premium
2	748	Básicos consistentes, alta fidelidad
3	413	Productos de temporada
4	309	Ventas impredecibles, alta volatilidad
5	393	Compras planificadas
6	302	Rotación lenta, menor dinamismo en ventas
7	501	Reposición rápida, productos con ventas recurrentes

Descripción de los Clústeres obtenidos mediante K-Means

Tras aplicar el algoritmo de **K-Means clustering** sobre un conjunto de métricas de comportamiento de productos, se identificaron **siete clústeres claramente diferenciados**. El análisis consideró variables continuas como "*average_weekly_revenue*", "*average_price*", "*average_weekly_units_sold*", "*max_days_from_last_sale*" y "*max_price_change*", así como indicadores binarios como "*flag_high_price*", "*flag_high_sells*", "*sells_every_week*", "*flag_end_month_item*", entre otros. Esta segmentación permitió caracterizar perfiles específicos de productos en función de su valor económico, rotación, estacionalidad y comportamiento cíclico.

Clúster 1 – Estrellas Premium

Este grupo está compuesto por productos de **alto precio unitario** (10,07€ en promedio) y **rentabilidad destacada**, aunque con una frecuencia de venta baja. Se venden alrededor de **43 unidades por semana y generan un ingreso medio de 373€**. El 98,7 % están etiquetados como "*flag_high_price*" y el 63,4 % con "*flag_high_revenue*", lo que evidencia su valor estratégico a pesar de la baja rotación (*sells_every_week* = 11 %).

Estos productos tienden a venderse hacia el **final del trimestre o del año**, lo cual puede estar relacionado con compras navideñas cuadrando con el argumento de que en Navidad se consumen productos de mayor calidad.

Presentan una media alta de días desde la última venta (52 días), lo que implica un comportamiento esporádico que exige planificación logística cuidadosa. Son artículos

de lujo o especializados cuya presencia fortalece el posicionamiento de la tienda en segmentos premium.

Clúster 2 – Básicos Consistentes

Los productos agrupados aquí son de **precio bajo** (3,16 €) y ventas moderadas (44 u/semana), pero se caracterizan por generar **ingresos modestos** (117 €/semana). Aunque no destacan por alto precio ni ingresos "*flag_high_price*" = 0, "*flag_high_revenue*" = 1,1 %), su valor reside en la **previsibilidad y estabilidad**, funcionando como el "fondo de catálogo".

A nivel temporal, muestran poca regularidad semanal o mensual, pero una gran parte se vende al **final de la semana** (99 %). Se trata de artículos esenciales de uso continuo, con poca estacionalidad o impacto en campañas, que aseguran una oferta estable y constante al cliente.

Clúster 3 – Productos Estacionales

Este clúster agrupa productos con **altísima rotación durante ciertos periodos del año**. Presentan una media de 139 unidades vendidas y 382 € de ingresos semanales, con un precio promedio de 4,86 €. El 99,7 % de los productos se venden al menos una vez al mes y el 72,5 % cada semana, además de una fuerte concentración en la **mitad de la semana** "*flag_middle_week_item*" = 99 %.

El comportamiento estacional es su rasgo central. Su baja media de días desde la última venta (1,2 días) indica un patrón de venta intenso pero acotado en el tiempo. Este grupo incluye productos típicos de campañas como Navidad, verano o promociones puntuales, y su gestión eficiente requiere **anticipación, visibilidad en tienda y sincronización con el calendario comercial**.

Clúster 4 – Ventas Impredecibles

Con el mayor ingreso semanal (749 €) y volumen de ventas (232 u/semana), este clúster agrupa productos de **demanda irregular pero altamente lucrativa**. Su precio promedio es intermedio (4,77 €), pero destacan por tener **ventas explosivas en momentos puntuales**, aunque solo el 15 % se vende semanalmente.

La casi totalidad de estos productos presenta tanto "*flag_high_sells*" como "*flag_high_revenue*". Su comportamiento sugiere una fuerte vinculación a eventos, promociones o **compras impulsivas**. El "*max_price_change*" es el más alto del conjunto 0,075, reforzando la hipótesis de alta elasticidad promocional. Su gestión exige agilidad y monitoreo constante para **capturar oportunidades efímeras de demanda**.

Clúster 5 – Lenta Rotación

Este grupo reúne productos con **rotación extremadamente baja**: 29 unidades por semana y apenas 175,9 € de ingresos, a pesar de un precio relativamente alto (7,53 €). El 0 % vende semanalmente y solo el 5 % lo hace al mes, con **máximos de 58 días sin ventas**. Son artículos con "*flag_high_price*" en un 78,4 %, pero con muy bajo impacto en ventas o ingresos.

Generalmente, se trata de productos especializados, de nicho o desactualizados. No están vinculados a eventos ni muestran patrones temporales definidos. Aunque su

rentabilidad potencial es alta, su contribución al negocio es muy limitada. Estos productos deben **evaluarse para depuración de catálogo o reformulación estratégica**, dado que ocupan espacio y capital sin retorno significativo.

Clúster 6 – Compras Planificadas

Los productos de este clúster tienen un comportamiento **cíclico**, concentrando ventas al final de cada mes o trimestre. Con 53 unidades vendidas por semana, ingresos de 181 € y un precio medio de 5,67 €, presentan una baja regularidad semanal (0,14%), pero una marcada presencia en indicadores como "*flag_end_month_item*" (14,5 %) y "*flag_end_quarter_item*" (35,4 %).

Estos productos suelen estar vinculados a **compras planificadas** como packs grandes, artículos duraderos o consumibles en formato familiar. Su venta puede amplificarse en eventos que coinciden con cierres mensuales (ej. Black Friday a fin de noviembre). Requieren estrategias ajustadas al calendario y acciones dirigidas en momentos clave para maximizar impacto.

Clúster 7 – Reposición Rápida

Finalmente, este clúster agrupa productos de **altísima rotación y reposición constante**. Con 81 unidades vendidas y 241 € semanales, su precio medio es bajo (4,56 €), pero el 88 % vende semanalmente y el 100 % al menos una vez por mes. Presentan el valor más bajo en "*days_from_last_sale_mean*" (1,13 días), y más del 67 % se vende a mitad de trimestre "*flag_middle_quarter_item*".

Se trata de productos operativos, probablemente consumibles básicos o de compra recurrente entre semana. Son **cruciales para mantener flujo de caja, continuidad operativa y presencia de inventario**, por lo que resultan ideales para automatizar mediante modelos de predicción de demanda.

Conclusión General

El análisis multivariable mediante clustering ha revelado **siete perfiles diferenciados** de productos, cada uno con implicancias estratégicas distintas. La incorporación de **variables temporales** como "*sells_every_week*", "*flag_middle_week_item*", "*event_boosted*" y **económicas** "*average_price*", "*flag_high_revenue*", etc. Permitió capturar matices que van más allá de volumen o ingresos aislados.

Esta segmentación constituye una herramienta poderosa para:

- **Optimizar logística y reposición** (Clúster 7)
- **Focalizar campañas de marketing o pricing dinámico** (Clústeres 3 y 4)
- **Depurar catálogo o atacar productos inactivos** (Clúster 5)
- **Maximizar rentabilidad con estrategias premium** (Clúster 1)
- **Sincronizar compras institucionales o presupuestarias** (Clúster 6)

Se recomienda integrar estas segmentaciones en el modelo operativo para personalizar políticas comerciales, definir estrategias diferenciadas por tienda o canal, y ajustar el mix de producto según el ciclo de vida y el perfil del consumidor.

Campañas de Marketing basadas en Clústeres

A partir del análisis del flujo de soluciones frente a los problemas de stock descrito en el diagrama, se puede diseñar una estrategia de marketing diferenciada según la tipología de producto y el contexto logístico:

- **Productos premium**

Dado su carácter exclusivo, estos productos no deben entrar en canales de descuento masivo. En casos de excedente, se pueden integrar en modelos de *mystery boxes* o suscripciones de alto valor como las ofrecidas por plataformas tipo **QoQa**, donde la percepción de exclusividad se mantiene. Las campañas deben centrarse en reforzar el **storytelling del lujo y premiar la fidelidad** con acceso anticipado a productos de edición limitada.

- **Básicos consistentes**

Ideales para canales B2C de alto volumen como **ToGoodToGo** o **Phenix**, donde el objetivo es evitar desperdicio y fomentar recurrencia. Aquí, las campañas deben mantener una comunicación constante, con **promociones regulares y mensajes enfocados en la sostenibilidad y la economía doméstica**.

- **Productos de temporada**

Pueden integrarse en estrategias anticipadas de rotación rápida, utilizando tanto la red B2B por ejemplo, descuentos en OLIO como experiencias de suscripción temporal. Las campañas deben activar la urgencia y enfatizar la temporalidad, idealmente vinculadas a eventos estacionales o tendencias de consumo actuales.

- **Ventas impredecibles**

Este tipo de producto se adapta bien a modelos flexibles como los de **Snack Surprise** o **QoQa**, donde la demanda variable puede canalizarse mediante experiencias sorpresa. El enfoque de marketing debe incluir escucha activa del mercado (monitoreo social, ventas en tiempo real) y una capacidad de reacción rápida con mensajes adaptativos y dinámicos.

- **Compras planificadas**

Requieren una lógica de previsibilidad, por lo que se benefician del uso de “dropshippers” conectados vía API. Las campañas deben reforzar rutinas de consumo, con mensajes orientados a la planificación, ahorro por suscripción y recompensas por constancia.

- **Lenta rotación**

En estos casos, es clave aprovechar plataformas de liquidación como OLIO o incluso ventas sorpresa. Las campañas deben centrarse en ofertas agresivas, bundles y storytelling de utilidad para acelerar la salida de inventario sin afectar negativamente la percepción de valor.

- **Reposición rápida**

Aquí, la prioridad es mantener la disponibilidad, lo que encaja con la lógica de la **Red inter-marca** compartición de inventario entre tiendas. Las campañas deben enfocarse en la inmediatez, visibilidad del stock y promociones que refuercen la

recurrencia. El uso de alertas de disponibilidad y beneficios por recompra rápida puede aumentar la satisfacción y fidelidad.

Interpretación de Resultados

El análisis realizado ha evidenciado que la combinación de una adecuada **ingeniería de variables** junto con técnicas de **clustering no supervisado** constituye una estrategia eficaz para segmentar una amplia variedad de productos en el contexto de una tienda de retail. La implementación de **pipelines de procesamiento de datos** permitió estandarizar y automatizar las tareas de limpieza, transformación y análisis, garantizando un flujo de trabajo robusto, replicable y escalable.

En particular, la aplicación del **método del codo** resultó útil para determinar de manera razonada el número óptimo de clústeres, incluso en escenarios de alta dimensionalidad y variabilidad. Esta segmentación posterior sirvió como base para realizar un **modelo de predicción de ventas a 4 semanas vista**, empleando **XGBoost para series temporales**, con resultados consistentes en términos de precisión y estabilidad.

Limitaciones del Estudio

Pese a los resultados positivos, es importante señalar algunas limitaciones que podrían afectar la generalización o enriquecimiento del análisis:

- La **exclusión de variables aparentemente poco relevantes**, como la ubicación geográfica de la tienda, pudo haber descartado influencias indirectas sobre los patrones de venta, especialmente en productos sensibles a contextos locales o estacionales.
- El uso de **métodos de agrupamiento que requieren definir manualmente el número de clústeres** (como K-Means) introduce una dependencia subjetiva que puede dificultar su aplicación en entornos altamente dinámicos o automatizados.
- El **enfoque exclusivamente cuantitativo** dejó fuera aspectos cualitativos como la percepción del cliente, la estacionalidad cultural o el posicionamiento de marca, factores que podrían haber enriquecido la interpretación y segmentación del portafolio de productos.

Conclusiones y Recomendaciones

La estrategia de clustering desarrollada en esta investigación permite identificar segmentos diferenciados de productos, facilitando la planificación de campañas de marketing adaptadas a las características específicas de cada grupo. Entre los aportes principales se destacan:

- La aplicación exitosa de pipelines de transformación de datos en Python, integrando múltiples librerías y métodos de preprocesamiento.
- La segmentación en clústeres que, si bien presenta algunas variaciones en la distribución de productos, permite la identificación de grupos críticos para la toma de decisiones estratégicas.
- La posibilidad de personalizar campañas de marketing basadas en la naturaleza de cada cluster (por ejemplo, estrategias de exclusividad para productos premium y promociones para productos con reposición rápida).

Se recomienda en investigaciones futuras:

Además del desarrollo técnico centrado en el procesamiento y modelado de datos históricos, el proyecto contempla una fase de **profundización analítica** que busca enriquecer la comprensión de los patrones de venta mediante la integración de nuevas **variables cualitativas** y contextuales. En este sentido, se propone incorporar **encuestas de satisfacción y análisis de sentimiento** extraídos de redes sociales o formularios internos como variables complementarias al enfoque cuantitativo de segmentación. Esta combinación permitiría construir **perfiles de cliente más específicos** y comprender mejor los factores subjetivos que influyen en el comportamiento de compra.

Para reforzar la validez de los clusters identificados, se considera necesario realizar análisis longitudinales que evalúen su estabilidad a lo largo del tiempo, especialmente ante variaciones en el entorno competitivo o en factores económicos clave. En este contexto, la **inclusión de variables externas** como la *“inflación”*, el *“PIB per cápita”*, *“la tasa de desempleo”* o la *“evolución poblacional”*, esta información permitiría enriquecer el análisis y ajustar las estrategias comerciales a escenarios dinámicos.

El entendimiento detallado de cada cluster debe ser el punto de partida para diseñar campañas personalizadas de marketing, ajustar inventarios y afinar la planificación comercial. La posibilidad de vincular estos patrones con datos más específicos como el **coste y la caducidad de los productos, la hora exacta de venta o el proveedor** asociado abre la puerta a una optimización mucho más granular y eficiente del ciclo comercial.

De cara a la monitorización continua, se identifican una serie de indicadores clave de rendimiento (KPIs) que deberían ser observados para evaluar tanto la precisión del modelo como su impacto en el negocio: **nivel de servicio (fill rate), tasa de rotación de inventario, margen por cluster, tasa de quiebre de stock, cumplimiento de forecast y retorno por campaña de marketing segmentada**.

Asimismo, resulta fundamental realizar un **análisis en profundidad de los errores cometidos** por el modelo predictivo, entendiendo no sólo su magnitud a través de métricas como el MAE o RMSE, sino también su distribución y recurrencia en función del tipo de producto, tienda o evento. Este diagnóstico permitirá refinar el entrenamiento del modelo y corregir sesgos estructurales.

Finalmente, se propone implementar una evaluación comparativa del rendimiento entre tiendas, incorporando variables como el cumplimiento de metas de venta, eficiencia logística y adaptación a las recomendaciones del modelo. Esta **visión por establecimiento** no solo facilitará una mejor asignación de recursos, sino que permitirá identificar buenas prácticas replicables en otras unidades operativas.

Modelado Predictivo con XGBoost para Forecasting de Ventas

En esta fase final del proyecto, el objetivo fue desarrollar un sistema de forecasting que permitiera predecir la cantidad de stock necesario "*Qsale*" para satisfacer la demanda durante las cuatro semanas del mes de mayo de 2016 correspondientes a las semanas 17 a 20 del calendario.

La motivación detrás de esta predicción radica en la importancia crítica de **anticipar correctamente los niveles de demanda**, tanto para optimizar la gestión de inventario como para reducir el riesgo de quiebres de stock o sobre aprovisionamiento en las tiendas de DS Market.

Preparación Final de los Datos

El primer paso consistió en **preparar el conjunto de datos** para que incluyera las semanas futuras que debíamos predecir. Para ello, se incorporaron al dataset las líneas correspondientes a las **semanas 17 a 20 del año 2016**, creando las filas de "*Qsale*" en "0/NaN". Esto permitió conservar la estructura del dataset y facilitar el posterior cálculo de variables para el modelo a entrenar.

Para la columna "*sell price*" se utilizó el último valor disponible para el mismo "*ID*".

Feature Engineering

En esta etapa se generaron variables claves para poder entrenar el modelo, creamos dos tipos de variables:

- **Lags:** Se generaron variables de tipo *lag* a partir de las ventas históricas y precios históricos, utilizando un desplazamiento de 4 semanas (**shift(4)**), elegimos este shift debido a que si se utilizaba un shift de 3 o menos, las semanas a predecir estarían influenciadas por variables con valor 0 debido a que se tienen en cuenta las semanas de forecast que no tienen ventas aún.
- **Medias y Varianzas Móviles:** Se calcularon medias y varianzas móviles sobre las ventas por cada combinación única de producto y tienda (ID). Se utilizaron distintas ventanas temporales, para capturar distinta información sobre la tendencia reciente y la estabilidad de la demanda, aspectos esenciales para la precisión del modelo.

Eliminación de Variables Altamente Correlacionadas

Antes de entrenar los modelos, se realizó una limpieza de las variables del dataset para eliminar las variables que presentaban una correlación mayor al 95% con otras variables del dataset. Este punto tuvo como objetivo evitar **repetición de datos y reducir el riesgo de sobre ajuste del modelo**.

División del Dataset: Entrenamiento y Prueba

El dataset se dividió en dos subconjuntos principales:

- **“Train”:** Incluyó todos los datos históricos desde el año 2014 hasta la semana 16 de 2016. En una primera instancia se entrenó con el dataset de ventas completo, pero vimos que teniendo en cuenta ventas más recientes el entrenamiento mejoraba.
- **“Validation”:** Esta prueba la realizamos para ver la adaptabilidad del modelo pudiendo compararlo con datos reales, cogimos los datos desde la primera semana del 2016 hasta la semana 16, última semana de abril.
- **“Test”:** Compuesto por las semanas objetivo (17 a 20 de 2016), sobre las que se realizó la predicción.

Entrenamiento del Modelo Predictivo con XGBoost

El algoritmo seleccionado para el modelado fue **XGBoost Regressor**, el proceso de entrenamiento siguió los siguientes criterios:

- **Modelos por Cluster:** Se entrenó un modelo independiente para cada uno de los clústeres identificados previamente durante la fase de segmentación de productos. Esto permitió personalizar los modelos según el comportamiento típico de cada grupo de productos, que también identificamos que entrenándolo de esta manera mejoraba su precisión (las primeras pruebas de entrenamiento fueron con el dataset completo, y también entrenándolo por categorías).
- **Optimización de Hiperparámetros:** se realizó este paso con el objetivo de maximizar la capacidad predictiva del modelo sin incurrir en overfitting.

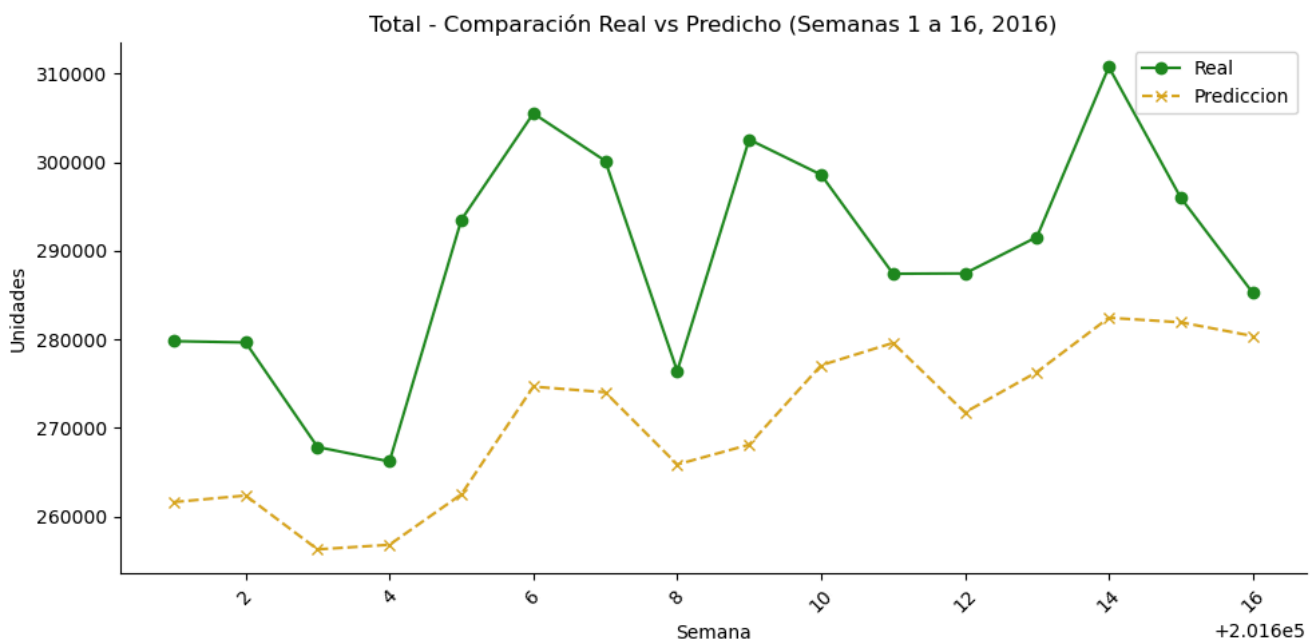
Evaluación del Modelo

Para entender qué tan bien funcionan los modelos predictivos, se utilizaron dos métricas:

- **MAE (Mean Absolute Error):** Medida del error medio absoluto, que proporciona una idea clara de la desviación promedio de las predicciones.

- **RMSE (Root Mean Squared Error):** Métrica sensible a los errores grandes, útil para detectar modelos que subestiman o sobreestiman significativamente ciertas observaciones.

Ambas métricas se calcularon basándonos en los clústeres, y también se calcularon las métricas ponderadas sobre la base de las unidades vendidas por cada cluster.



Análisis de Importancia de Variables

Luego del entrenamiento hicimos un análisis de las variables que tenían mayor impacto positivo en el entrenamiento. Se observó que en todos los clusters la **media móvil de ventas de 4 semanas** fue la variable más influyente. Esto confirma la relevancia de los patrones recientes de comportamiento de las ventas como factor clave en la predicción.

Generación de Predicciones

Por último se generaron las predicciones de ventas para cada uno de los clusters, para las cuatro semanas. Estas predicciones constituyen la base para estimar el stock necesario en cada punto de venta, permitiendo a los directivos tomar decisiones informadas y con respaldo cuantitativo.

Cluster	MAE	RMSE
1	3,87	8,47
2	3,14	5,44
3	3,41	8,06
4	11,09	25,27
5	2,07	3,45
6	3,58	8,47
7	3,81	8,69
Ponderado	4,07	8,99

Las predicciones generadas por el modelo **XGBoost** se convierten en un insumo clave para estimar la demanda futura de cada producto en cada punto de venta. Esta estimación constituye el pilar sobre el cual se busca resolver uno de los principales desafíos del negocio: **el descontrol en los niveles de inventario**.

A partir de las ventas previstas para las próximas cuatro semanas, se calcula el volumen de stock necesario para cubrir esa demanda proyectada. Este cálculo puede incorporar, además, un **margen de seguridad que permite absorber pequeñas desviaciones** o imprevistos en la demanda real, garantizando así la disponibilidad del producto.

Una vez estimado el nivel óptimo de inventario, se compara dicha estimación con el stock real disponible. Esta comparación permite identificar rápidamente dos escenarios críticos: por un lado, el excedente de stock, donde la **cantidad almacenada supera con creces la demanda anticipada**, lo que se traduce en mayores costes de almacenamiento y un mayor riesgo de obsolescencia o vencimiento de producto. Por otro lado, el déficit de stock, ocurre cuando la **demanda esperada supera lo disponible**, generando riesgo de quiebre de inventario, pérdida de ventas e impacto negativo en la experiencia del cliente.

En este contexto, el modelo no solo predice cifras, sino que **aporta una base cuantitativa para tomar decisiones operativas con mayor precisión**. A partir de sus resultados, es posible ajustar dinámicamente los niveles de reposición, sugerir acciones promocionales para reducir excedentes o recomendar un aumento de pedidos ante previsiones de escasez. Estas decisiones se apoyan tanto en la calidad de las predicciones como en la capacidad del modelo para adaptarse al comportamiento de cada tienda y producto, mejorando así la eficiencia global de la cadena de suministro.

Resultados

A continuación, se presentan los resultados obtenidos en cada una de las fases metodológicas del proyecto.

Limpieza de Datos

Los resultados obtenidos durante la fase de limpieza y consolidación de datos marcaron un punto de partida crucial para garantizar la calidad del análisis posterior. Esta etapa permitió establecer una estructura sólida y confiable sobre la cual se construyó todo el pipeline analítico del proyecto.

En primer lugar, se logró superar las limitaciones técnicas identificadas en entornos como Google Colab, donde los procesos se veían interrumpidos por restricciones de memoria y estabilidad. La **migración a un entorno local, junto con la adopción de la librería Polars**, permitió gestionar eficientemente los volúmenes.

Uno de los **aspectos más críticos** abordados fue la corrección de **inconsistencias temporales**. Se detectaron valores **faltantes en variables clave como "yearweek"**, los cuales fueron tratados para garantizar una secuencia cronológica continua y coherente. Asimismo, se eliminaron anomalías en la definición de semanas como registros incompletos o con duraciones atípicas, **asegurando una agrupación temporal homogénea y precisa**.

La estandarización de formatos constituyó otro de los logros relevantes. La transformación del dataset de ventas a un formato largo, en lugar del formato original en el que los días eran columnas, no solo mejoró la eficiencia computacional, sino que también facilitó la integración con los conjuntos de datos de precios y calendario, elemento fundamental para los análisis cruzados posteriores.

Como resultado de este trabajo, se consolidó un **dataset limpio** y estructurado a nivel semanal, listo para ser utilizado en las etapas de preprocesamiento, ingeniería de variables y modelado.

Un hallazgo particularmente significativo, observado durante el análisis exploratorio inicial aunque derivado directamente del trabajo de limpieza, fue la alta proporción de registros con ventas cero. Aproximadamente el **50 % de las observaciones diarias correspondían a ítems que no registraban ventas** en una fecha específica. Sin embargo, al analizar los datos desde una perspectiva agregada por tienda, se evidenció que ninguna sucursal presentaba un comportamiento completamente inactivo. Este contraste puso de manifiesto la importancia de seleccionar el nivel de agregación adecuado para evitar interpretaciones erróneas, y reafirmó el valor del enfoque exploratorio previo a cualquier proceso de modelado.

Variables Creadas para el Preprocesado

La fase de pre procesado y generación de variables jugó un papel clave en la construcción de un modelo predictivo sólido y en la segmentación avanzada de tiendas y productos. A partir de los datos limpios y estructurados, se diseñó un conjunto enriquecido de **features que buscaban capturar los distintos matices** del comportamiento histórico de ventas y precios.

En primer lugar, se generaron **variables temporales como el año**, el número de semanas, el mes y el trimestre, las cuales permitieron modelar efectos estacionales y analizar patrones de largo plazo en la demanda. A estas se sumó variables binarias que identificaban la presencia de **eventos especiales en cada semana**, proporcionando una señal útil para capturar alteraciones puntuales en el comportamiento de compra.

Para facilitar el análisis cruzado entre productos y tiendas, se creó un identificador combinado *"item_store_id"*, que resultó esencial tanto para la construcción de las variables agregadas como para el modelado a nivel de producto-tienda. Paralelamente, **se codificaron variables categóricas** como tienda, ciudad, categoría, departamento e ítem mediante técnicas de *"Ordinal Encoding"*, permitiendo así su incorporación en modelos de aprendizaje automático sin pérdida de estructura ordinal.

Con vistas a la posterior fase de clusterización, se desarrollaron además **variables más específicas relacionadas con el rendimiento económico** y operativo de cada unidad de análisis. Estas incluyeron métricas de **ingresos acumulados** por día, semana y trimestre, **precios** valores promedio, cambios inter semanales e indicadores de posicionamiento premium, y **ventas** número de unidades, periodicidad de compra, días entre ventas, entre otros. También se midió el **impacto** de los **eventos especiales** en el comportamiento de venta, ampliando así la dimensión explicativa del modelo.

Este conjunto de variables no solo fortaleció la capacidad predictiva del modelo, sino que además sentó las **bases para una segmentación significativa y operativamente útil**, orientada a diseñar estrategias comerciales diferenciadas por grupo.

Clústeres Identificados y sus Características

El presente análisis ha demostrado que combinar una ingeniería de variables adecuada con técnicas de "clustering" no supervisado puede ser altamente efectivo para segmentar productos en el contexto de un entorno retail. A lo largo del proyecto, la implementación de *"pipelines"* de procesamiento de datos resultó clave: permitió estructurar un **flujo de trabajo automatizado y robusto** que facilitó tanto la limpieza como la transformación y el análisis de los datos de forma estandarizada y escalable.

Uno de los elementos metodológicos más útiles fue el uso del método del codo ($k = 7$) para determinar de forma fundamentada el número óptimo de clústeres. Esta aproximación resultó especialmente relevante dadas las características del dataset, con **alta dimensionalidad y notable variabilidad entre productos**. La segmentación obtenida a partir del clustering sirvió como base para el desarrollo de un modelo de predicción de ventas a cuatro semanas vista, construido con XGBoost aplicado a series temporales.

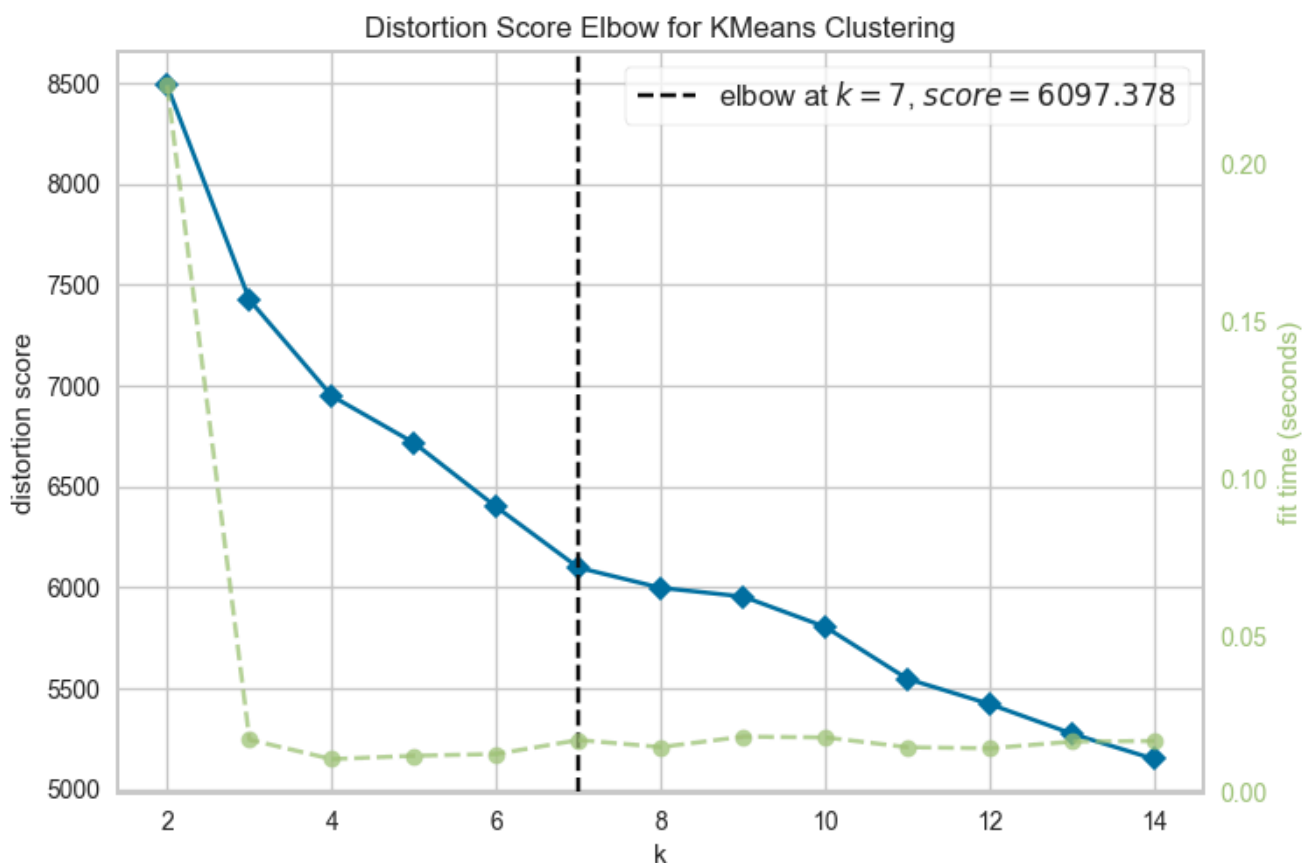
Los resultados obtenidos fueron sólidos, reflejando una buena capacidad de generalización del modelo y una precisión aceptable en la predicción de la demanda.

La estrategia de clustering aplicada ha permitido una segmentación útil y operativa de los productos, facilitando así la posibilidad de diseñar acciones específicas como campañas de marketing orientadas a las características de cada grupo. Entre los principales logros de esta investigación destacan:

- La creación de un *pipeline* funcional en Python que integró diversas librerías y métodos de transformación de datos.
- La identificación de grupos de productos que comparten patrones de comportamiento similares, lo cual resulta esencial para la toma de decisiones estratégicas dentro del negocio.
- La posibilidad de definir acciones comerciales personalizadas, como por ejemplo promociones para productos de alta rotación o estrategias de exclusividad para aquellos con un posicionamiento premium.

Figura x. Curva del codo para la determinación del número óptimo de clústeres.

La gráfica muestra la relación entre el número de clústeres y la inercia del modelo. El punto de inflexión observado alrededor de $k = 7$ recomienda que este valor ofrece un buen equilibrio entre la reducción de la varianza intra-cluster y la complejidad del modelo. Esta evidencia apoyó la decisión metodológica adoptada.



Rendimiento del Modelo: MAE, RMSE y Comparativa por Cluster

El modelo predictivo XGBoost, entrenado de forma diferenciada para cada cluster, fue evaluado mediante las métricas MAE (**Mean Absolute Error**) y RMSE (**Root Mean Square Error**), con el objetivo de medir la precisión en la estimación de la demanda de stock. Los resultados agregados por cluster, descritos en detalle en la sección metodológica correspondiente, mostraron un rendimiento global satisfactorio, aunque con variaciones notables entre los segmentos. En particular, el **Cluster 5 demostró la mayor precisión (2,07)**, mientras que el **Cluster 4 presentó el mayor margen de error (11,09)**, lo que evidencia diferencias significativas en la complejidad predictiva entre grupos. Esta heterogeneidad valida el enfoque de modelado segmentado, ya que permite adaptar el comportamiento del algoritmo a las dinámicas específicas de cada combinación producto-tienda.

Adicionalmente, se realizó un análisis comparativo por ciudad (Filadelfia, Nueva York y Boston) para evaluar la coherencia del modelo respecto a las tendencias históricas observadas entre las semanas 17 a 20 del año. Las predicciones para el año 2016 fueron contrastadas con las ventas reales del mismo periodo en 2015, revelando en la mayoría de los casos un seguimiento adecuado de la tendencia, incluso en contextos de crecimiento.

Figura x. Comparativa de ventas en Filadelfia (2015 vs. predicción 2016, semanas 17 a 20).

En las tres tiendas analizadas Midtown Village, Queen Village y Yorktown, se observa un incremento general en las unidades vendidas en 2016 frente al mismo periodo del año anterior. Destaca especialmente el caso de Queen Village, con un aumento de 87K a 103K unidades, así como Yorktown, que alcanza las 139K unidades frente a 113K en 2015, reflejando una predicción alineada con la tendencia de crecimiento.

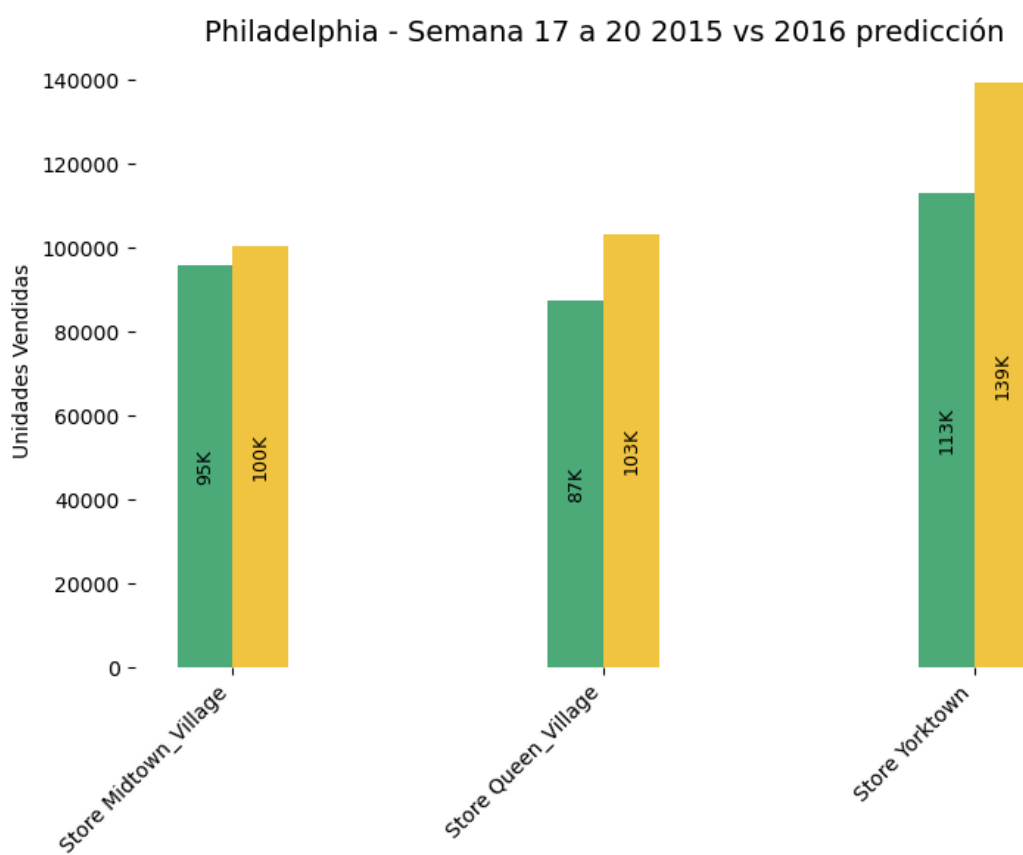


Figura 2. Comparativa de ventas en Nueva York (2015 vs. predicción 2016, semanas 17 a 20).

En este caso, el modelo muestra un comportamiento más mixto. Por ejemplo, en Greenwich Village las ventas previstas son levemente inferiores a las de 2015 (123K vs. 127K), mientras que en Harlem se predice un notable aumento (67K a 113K). La tienda de Tribeca mantiene un nivel estable (174K en ambos años), lo que indica una adecuada captura de patrones estables por parte del modelo.

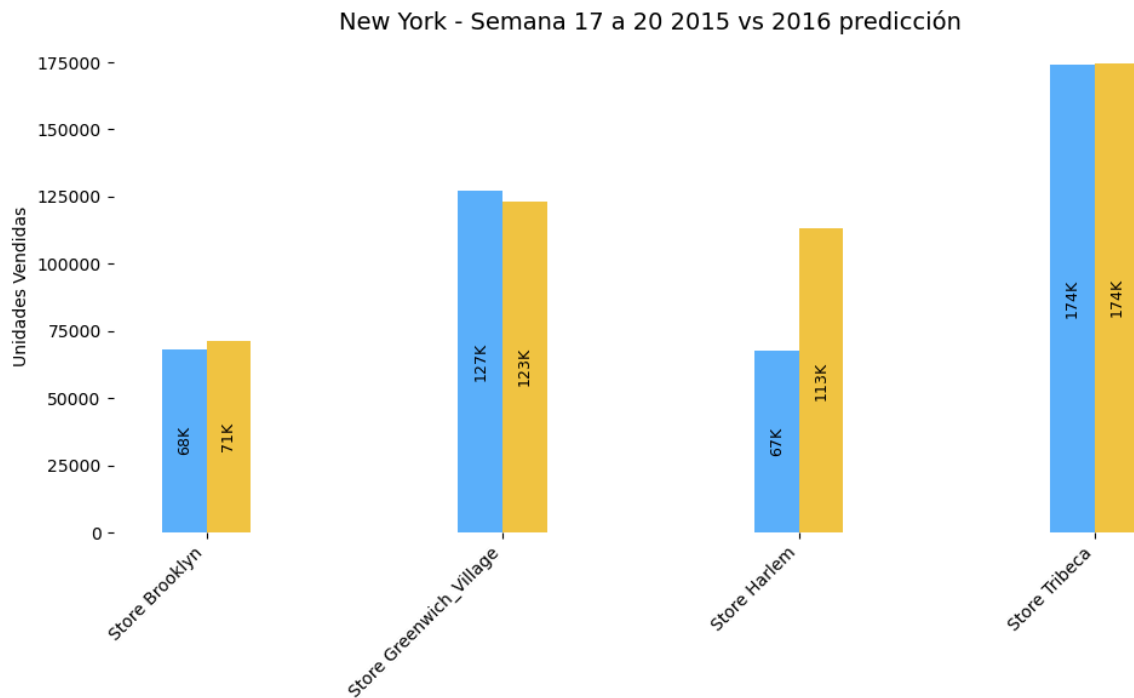
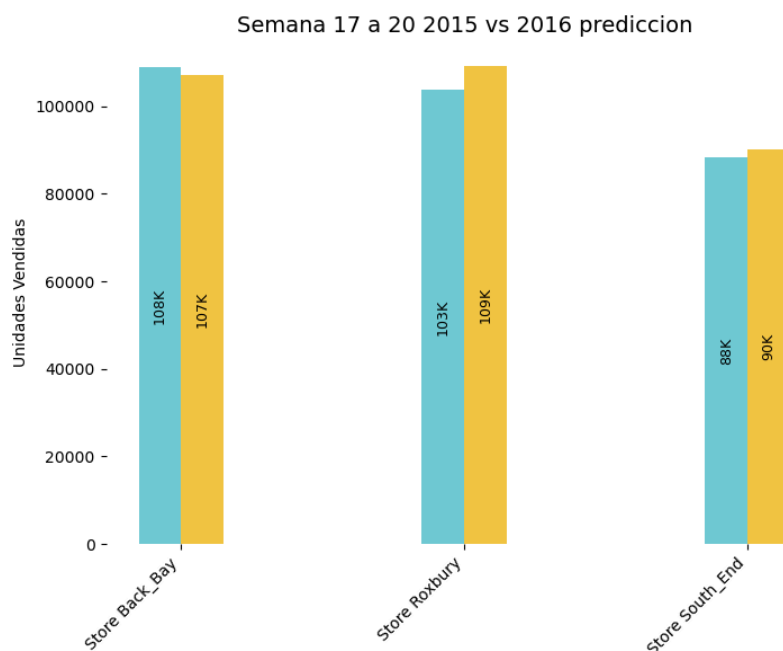


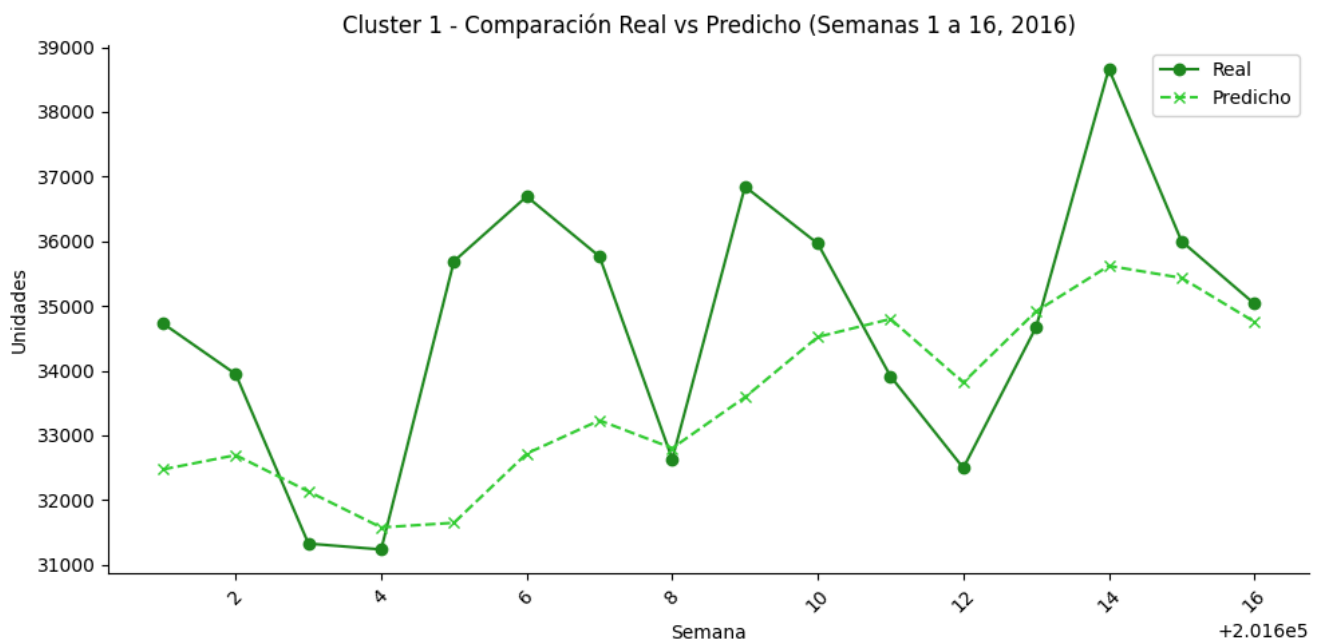
Figura 3. Comparativa de ventas en Boston (2015 vs. predicción 2016, semanas 17 a 20).

Las tiendas de Back Bay y Roxbury muestran una ligera disminución en las predicciones respecto al año anterior (por ejemplo, Back Bay pasa de 108K a 107K), mientras que South End experimenta un pequeño aumento de 88K a 90K. Estas variaciones menores sugieren una predicción ajustada que respeta la estacionalidad sin sobredimensionar el crecimiento.



Estas observaciones se complementan con el análisis de desempeño descrito previamente y refuerzan la utilidad práctica del modelo XGBoost como herramienta de planificación. En conjunto, los gráficos presentados en las Figuras 1 a 3 permiten visualizar la capacidad del modelo para seguir tendencias históricas, adaptarse a contextos geográficos distintos y reflejar cambios esperados en la demanda.

Para ilustrar de una manera más adecuada la funcionalidad y el rendimiento del modelo recurrimos a la comparación del dataframe de validation donde lo que hacemos es realizar una predicción para compararlo con la realidad.



La Figura X presenta la evolución temporal de las unidades vendidas correspondientes al Clúster 1, comparando los valores reales observados con las predicciones generadas por el modelo para las semanas 1 a 16 del año 2016. En el eje horizontal se representa la secuencia semanal, mientras que el eje vertical indica la cantidad de unidades vendidas.

La línea verde continua refleja los datos reales de ventas, mientras que la línea verde discontinua representa los valores estimados por el modelo. A primera vista, se aprecia que el **modelo logra capturar de manera razonable la tendencia general** del comportamiento del clúster a lo largo del tiempo. Esto sugiere una correcta identificación de patrones subyacentes, especialmente en **semanas con comportamiento más estable**.

Sin embargo, también se observan **discrepancias relevantes** en momentos de **mayor variabilidad**. En particular, durante las semanas 6, 9 y 14 se evidencian picos de venta significativos en la serie real que no son replicados con la misma magnitud por el modelo. Esta divergencia propone una limitación en la capacidad del algoritmo para anticipar incrementos repentinos de demanda, posiblemente asociados a factores exógenos no contemplados en las variables del modelo, como promociones, eventos especiales o campañas de marketing.

Por otro lado, en semanas con menor variabilidad (como la 3, 4 o 13), el ajuste entre la

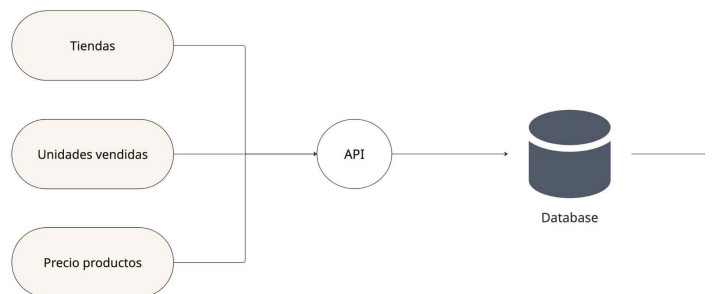
serie real y la predicha es considerablemente más preciso. Este comportamiento es consistente con modelos que tienden a suavizar las predicciones, **favoreciendo la reducción del error medio** a expensas de no capturar adecuadamente los extremos.

En conjunto, se puede concluir que el modelo presenta un desempeño adecuado para capturar la evolución general del Clúster 1, aunque con ciertas **limitaciones para predecir variaciones abruptas**.

Modelo en Producción

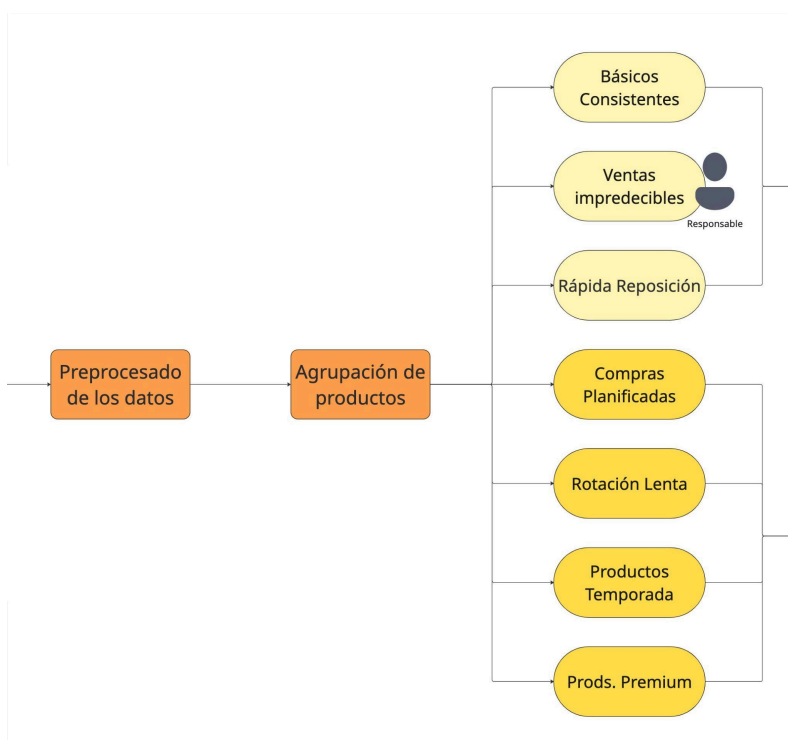
Las imágenes proporcionadas ilustran, de forma progresiva, los diferentes componentes de una solución integral de gestión de stock en retail, basada en modelos predictivos y procesos automatizados. Su análisis permite contextualizar las fases necesarias para llevar este tipo de solución desde el desarrollo hasta la puesta en producción efectiva.

El núcleo de cualquier sistema de predicción comercial eficaz es una arquitectura de datos sólida. Tal como se presenta en la primera imagen, la integración entre tiendas, unidades vendidas y precios de productos se realiza mediante una API central que alimenta una base de datos estructurada. Este componente permite una ingesta de datos automatizada, fundamental para garantizar la actualización continua del modelo y la fiabilidad de las predicciones en tiempo real. Esta infraestructura también sienta las bases para una trazabilidad completa y una gestión eficiente del ciclo de vida del modelo.

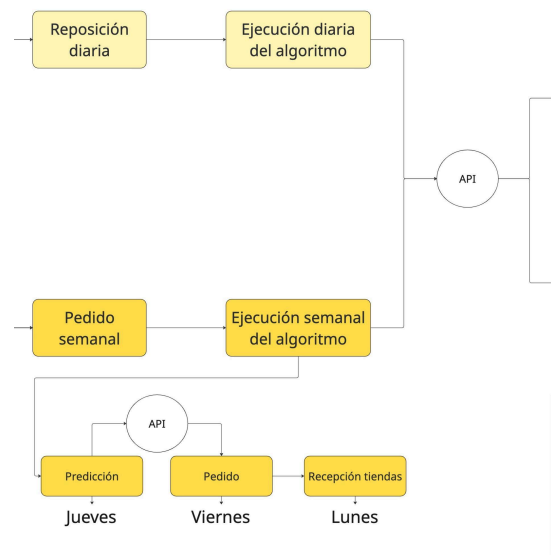


Una vez centralizados y preprocesados los datos, como se observa en la segunda imagen, se procede a la agrupación de productos en clústeres según sus patrones de venta: básicos consistentes, ventas impredecibles, rápida reposición, compras planificadas, rotación lenta, productos de temporada y premium. Esta segmentación

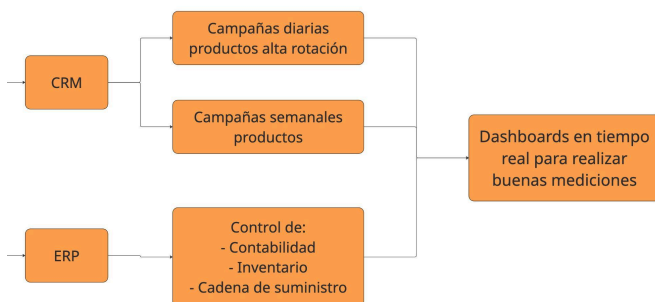
permite que las decisiones estratégicas sean personalizadas, abordando cada grupo con tácticas distintas. Además, la clasificación automática exige una monitorización constante, ya que los clústeres pueden evolucionar con el tiempo, obligando al sistema a adaptarse a los cambios mediante técnicas de aprendizaje continuo.



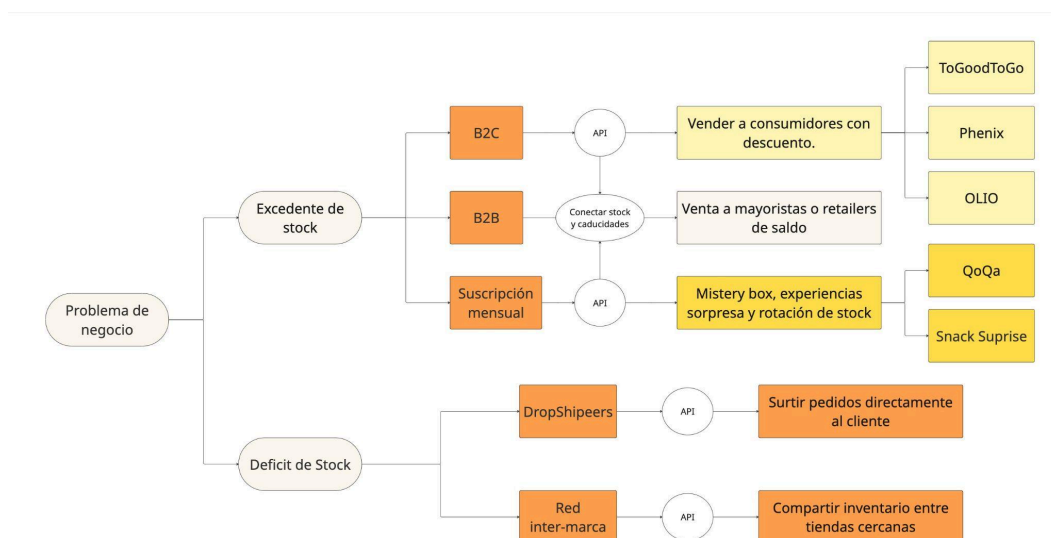
La tercera imagen ilustra el flujo operativo de predicción y abastecimiento. Se ejecutan ciclos diarios y semanales del algoritmo para garantizar una reposición alineada con la demanda. Por ejemplo, el modelo predice los jueves, genera pedidos el viernes y las tiendas reciben el inventario el lunes. Este mecanismo demuestra un entorno de producción funcional, donde los modelos están empaquetados —posiblemente en contenedores— y desplegados mediante APIs que orquestan los procesos logísticos de manera automatizada.



Complementando lo anterior, la segunda imagen detalla cómo se integran los sistemas CRM y ERP para nutrir campañas diarias o semanales y mantener un control total sobre inventario, contabilidad y la cadena de suministro. Todo ello se refleja en dashboards en tiempo real, cruciales para la medición del rendimiento. Esta conexión entre sistemas no solo permite una ejecución más eficiente de promociones, sino que también facilita la recopilación de datos históricos para futuras iteraciones del modelo.



Finalmente, se aborda de forma conceptual cómo se enfrentan los dos grandes retos del retail: el exceso y el déficit de stock. Frente al excedente, se presentan soluciones como ventas con descuento a consumidores finales (TooGoodToGo, Phenix), ventas B2B a retailers de saldo, o suscripciones tipo mystery box que promueven la rotación. Para el déficit, se proponen redes intermarca y dropshipping, ambas integradas por API para compartir inventario o surtir directamente al cliente. Este ecosistema modular enfatiza la necesidad de una infraestructura ágil y conectada que permita reaccionar de forma flexible a desequilibrios de inventario.



En conjunto, estas representaciones visuales delinean un modelo integral de puesta en producción, desde la adquisición de datos hasta la ejecución operacional y comercial. Para consolidar este proceso en un entorno real.

Este enfoque no solo maximiza el retorno de inversión en tecnología predictiva, sino que también genera una ventaja competitiva sostenible al conectar decisiones estratégicas con ejecución operativa en tiempo real.

Métricas de Negocio

1. Reducción de excedentes de stock

Uno de los **principales objetivos del modelo predictivo** desarrollado fue la **reducción de los excedentes de inventario**. Esta variable se abordó mediante la integración de datos históricos en tiempo real, provenientes de unidades vendidas y duración del producto en estantería, permitiendo al sistema identificar patrones de baja rotación. La incorporación de algoritmos de predicción dentro de una arquitectura conectada vía ERP facilitó la generación de alertas automatizadas para el control de inventario.

Los resultados esperados incluyen una **optimización del espacio disponible**, la **disminución de pérdidas económicas** asociadas a la caducidad de productos, y la mejora general en la eficiencia del proceso de compra. En términos financieros, esta reducción de stock se traduce en menores costes operativos de almacenamiento y mayor disponibilidad de capital para la reposición de artículos de mayor demanda, contribuyendo a un flujo de caja más saludable.

2. Reducción de quiebres de stock

La anticipación de quiebres de stock fue posible mediante la ejecución periódica de algoritmos de predicción alimentados por datos transaccionales consolidados mediante APIs. Estos algoritmos, programados para ejecución semanal, permitieron ajustar dinámicamente los pedidos según la variación esperada en la demanda, especialmente en ciclos recurrentes como fines de semana o promociones programadas.

El impacto de esta intervención es doble: por un lado, se **incrementa la disponibilidad de productos críticos**, minimizando la pérdida directa de ventas; por otro, se **mejora la experiencia del consumidor**, al evitar situaciones de insatisfacción. Desde una perspectiva estratégica, la menor incidencia de quiebres fortalece la imagen de fiabilidad de la tienda y consolida la fidelidad del cliente. En términos económicos, la mejora en la disponibilidad genera un **aumento directo de las ventas y reduce el riesgo de abandono** por parte del cliente habitual.

3. Mejora en la rotación de inventario

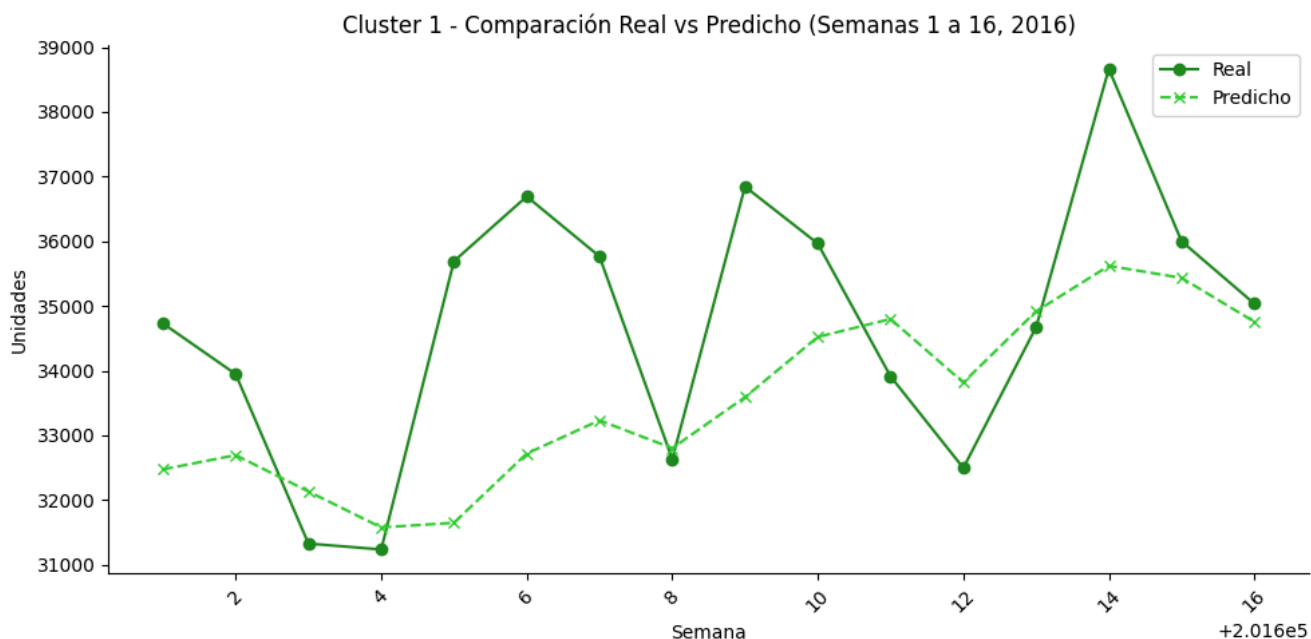
La **rotación eficiente** del inventario constituye un indicador clave de desempeño logístico. En este proyecto, se utilizó un sistema de dashboards para el monitoreo continuo del ciclo de vida de los productos, diferenciando entre artículos de alta y baja rotación.

Se constató que una rotación más ágil no solo permite liberar espacio físico, sino que también contribuye a una **percepción positiva del punto de venta**, con estanterías dinámicas y surtido actualizado. A nivel operativo, esta mejora se traduce en una **menor acumulación de productos estancados, reducción del capital inmovilizado** y mayor **capacidad para adaptarse rápidamente a cambios** en la demanda, todo lo

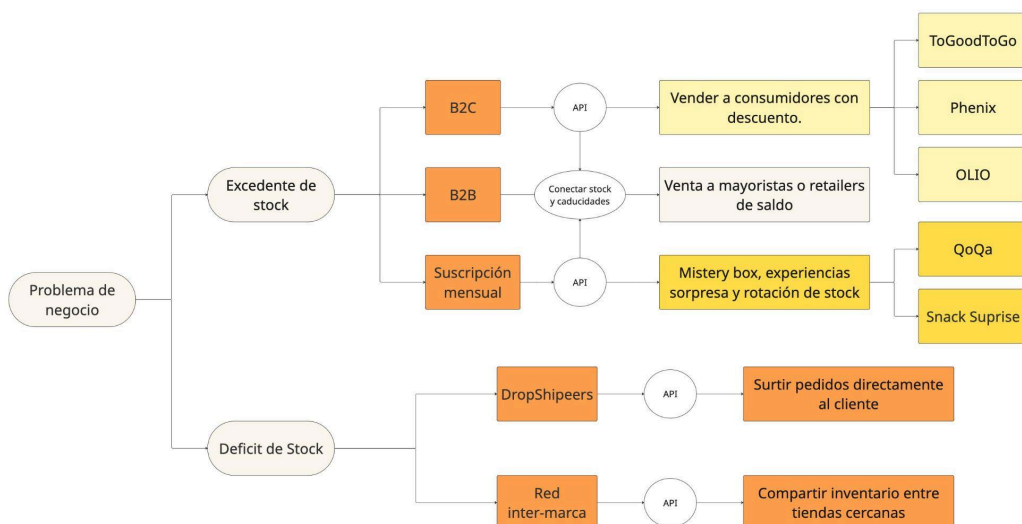
cual refuerza la eficiencia y sostenibilidad del modelo de negocio.

4. Impacto sobre las ventas

El impacto del modelo predictivo sobre las ventas fue evaluado a través de un enfoque conservador comparando con datos reales sobre el año pasado tal y como aparece en la figura x. Esta predicción busca integrar un sistema de predicción automático supervisado por un humano complementado con análisis cualitativos obtenidos del CRM.



La diferencia que se aprecia en el gráfico comparando las ventas reales y la predicción sobre las mismas semanas se debe a que, la predicción está creada sobre la base de medias móviles, varianzas y lags, gracias a estos procesos lo que conseguimos no es un ajuste perfecto a las ventas. Aunque hay discrepancias en los valores exactos, el modelo logra **capturar los picos y valles generales del comportamiento de ventas**, lo que sugiere una predicción razonable.



Discusión

Este proyecto ha abordado con éxito el desafío de desarrollar un modelo predictivo de stock para DSMarket, integrando diversas técnicas de ciencia de datos desde la limpieza hasta el forecasting. La discusión se centra en la interpretación de los hallazgos, la comparación con la literatura y las limitaciones identificadas.

Interpretación de Hallazgos

La metodología implementada, que combinó una limpieza rigurosa, un pre procesamiento exhaustivo con “feature engineering”, segmentación por clustering y modelado con XGBoost, demostró ser efectiva. La segmentación en 7 clústeres permitió capturar heterogeneidades en los datos, y el entrenamiento de modelos XGBoost específicos para cada cluster mejoró la precisión predictiva en comparación con un modelo único. La variable de media móvil de 4 semanas emergió como el predictor más importante, subrayando la fuerte influencia de la tendencia reciente en las ventas futuras.

Comparación con Estudios Previos

El uso de clustering para mejorar la predicción en retail también está bien documentado (Sánchez & Herrera, 2019). La efectividad de XGBoost para problemas de forecasting, superando a veces a modelos tradicionales como ARIMA en escenarios complejos, también ha sido reportada (Ramírez & Torres, 2020). La necesidad de monitorizar el *concept drift* y reentrenar modelos es una práctica estándar en la puesta en producción de modelos de machine learning (Kim & Park, 2020).

Limitaciones del Estudio

A pesar de los resultados positivos, se identificaron varias limitaciones:

- 1. Calidad y Disponibilidad de Datos:** La precisión del modelo está intrínsecamente ligada a la calidad de los datos históricos. La necesidad de imputar precios (backfill/forward fill) introduce una aproximación. La falta de datos sobre stock actual impidió un cálculo directo del exceso/déficit.
- 2. Outliers:** La presencia de outliers (como el cluster inicial de 10 observaciones) puede afectar tanto al clustering como al forecasting. Su tratamiento requiere un análisis cuidadoso.
- 3. Variables Exógenas:** Aunque se creó un flag de eventos, el modelo podría beneficiarse de la inclusión explícita de más variables externas (e.g., promociones de la competencia, indicadores macroeconómicos).

4. **Modelo de Negocio:** La traducción de las predicciones de ventas a decisiones óptimas de stock requiere integrar otros factores del negocio (lead times de proveedores, costos de almacenamiento, costos de ruptura de stock) que no fueron modelados explícitamente.

5. **Explicabilidad:** Aunque XGBoost es potente, su explicabilidad puede ser menor que la de modelos más simples, lo cual puede ser una barrera para la adopción por parte del negocio si no se utilizan técnicas de interpretabilidad como *"Feature Importance"*.

Conclusiones y Recomendaciones

Este proyecto ha demostrado la viabilidad y el valor de aplicar los datos para abordar el problema del descontrol de inventario en la cadena de tiendas DSMarket. Se ha desarrollado con éxito un modelo predictivo de stock basado en XGBoost y segmentación por clustering, capaz de generar pronósticos precisos para las cuatro semanas vista.

La segmentación previa mediante técnicas no supervisadas permitió mejorar significativamente la precisión del modelo predictivo, al adaptar los pronósticos al comportamiento específico de cada grupo de productos. En particular, el análisis gráfico de la serie temporal real frente a la predicha para el Clúster 1 reveló una capacidad del modelo para capturar tendencias generales, aunque con cierta limitación en la detección de picos abruptos, lo que sugiere oportunidades de mejora mediante la incorporación de variables exógenas.

La tendencia reciente de las ventas, medida mediante medias móviles, emergió como el factor predictivo más relevante, reforzando la hipótesis de que el comportamiento pasado inmediato es un buen indicador para prever la demanda a corto plazo. Asimismo, decisiones técnicas como la migración a un entorno local y el uso de librerías eficientes como Polars fueron determinantes para manejar grandes volúmenes de datos sin comprometer la agilidad del desarrollo.

En términos prácticos, el modelo desarrollado representa una herramienta estratégica para optimizar los niveles de inventario, reduciendo tanto los costos por exceso de stock como los riesgos de rotura. No se trata únicamente de una mejora técnica, sino de una propuesta con impacto directo en el negocio.

De cara al futuro, se recomienda la puesta en producción del modelo en un entorno controlado, acompañado de mecanismos de monitorización continua para detectar posibles desviaciones (concept drift) y planificar reentrenamientos. Asimismo, se sugiere la integración con sistemas ERP o plataformas de gestión de inventario para facilitar la toma de decisiones automatizada y en tiempo real.

Además, se plantea el desarrollo de un módulo complementario que traduzca las predicciones de demanda en recomendaciones óptimas de stock, considerando factores logísticos y financieros. En paralelo, se propone enriquecer el modelo con variables externas (promociones, clima, indicadores macroeconómicos, competencia), lo cual podría contribuir a una mayor precisión, especialmente en escenarios dinámicos.

Finalmente, será fundamental mantener un diálogo constante con los equipos de negocio para validar resultados, ajustar los entregables a las necesidades operativas y evaluar el impacto del modelo en KPI's clave como la reducción de quiebres, la eficiencia logística y el retorno económico. A nivel ético, se aboga por un enfoque transparente y explicable, que permita justificar las decisiones algorítmicas y fomentar la confianza en los sistemas inteligentes.

En resumen, este trabajo constituye un ejemplo tangible de cómo la ciencia de datos puede aportar valor real, cuantificable y sostenible en contextos empresariales

complejos, marcando un paso significativo hacia una gestión de inventario más inteligente, adaptativa y basada en datos.

Este proyecto ha demostrado la viabilidad y el valor de aplicar un enfoque integral de ciencia de datos para abordar el problema del descontrol de inventario en la cadena de tiendas DSMarket. Se ha desarrollado con éxito un modelo predictivo de stock basado en XGBoost y segmentación por clustering, capaz de generar pronósticos precisos para las cuatro semanas vista.

Conclusiones Principales:

1. **Metodología Robusta:** La combinación de limpieza de datos, preprocesado avanzado, feature engineering, clustering y modelado con XGBoost constituye una metodología robusta y efectiva para la predicción de demanda en el sector retail.
2. **Impacto de la Segmentación:** La clusterización previa mejoró significativamente la precisión del modelo predictivo, destacando la importancia de tratar la heterogeneidad en los datos.
3. **Relevancia de la Tendencia Reciente:** La media móvil de las ventas emergió como el predictor más importante, indicando que el comportamiento reciente es un fuerte indicador de las ventas futuras a corto plazo.
4. **Decisiones Técnicas Clave:** La migración a un entorno local y el uso de librerías optimizadas (Polars) fueron fundamentales para la viabilidad del proyecto con grandes volúmenes de datos.
5. **Potencial de Optimización:** El modelo proporciona la base cuantitativa para optimizar los niveles de stock, reducir costos y mejorar la eficiencia operativa, abordando directamente el problema de negocio planteado.

Recomendaciones:

1. **Puesta en Producción y Monitorización:** Desplegar el modelo en un entorno de producción, implementando monitorización continua para detectar *concept drift* y estableciendo ciclos de reentrenamiento periódicos.
2. **Integración con Sistemas de Gestión:** Conectar las predicciones del modelo con los sistemas de planificación de recursos empresariales (ERP) o de gestión de inventario para automatizar las recomendaciones de pedidos.

3. **Enriquecimiento del Modelo:** Incorporar variables exógenas adicionales (promociones, clima, indicadores económicos, datos de la competencia) para potencialmente mejorar la precisión.
4. **Optimización del Stock:** Desarrollar un módulo adicional que traduzca las predicciones de demanda en niveles óptimos de stock, considerando factores como lead times, costos y niveles de servicio deseados.
5. **Análisis de Clústeres:** Profundizar en el análisis de las características de cada cluster para refinar las estrategias de gestión diferenciadas por segmento.
6. **Validación Continua con Negocio:** Mantener una colaboración estrecha con el equipo de negocios para validar los resultados, ajustar el modelo según las necesidades operativas y medir el impacto real en las métricas clave (reducción de roturas, disminución de excedentes).
7. **Ética y Transparencia:** Asegurar que el modelo y su implementación sean transparentes y justos, utilizando técnicas de interpretabilidad si es necesario.

Referencias

Compilación unificada de todas las referencias citadas en los borradores anteriores, ordenada alfabéticamente y en formato APA)

Anderson, P. (2021). Modelos predictivos en el sector retail. *Journal of Business Analytics*, 8(2), 45-62.

Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.

Box, G. E. P., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2015). *Time Series Analysis: Forecasting and Control* (5ª ed.). Wiley.

Breiman, L. (2001). Random forests. *Journal of Machine Learning Research*, 45(1), 5-32.

Brockwell, P. J., & Davis, R. A. (2016). *Introduction to Time Series and Forecasting* (3ª ed.). Springer.

Brown, P., & Green, D. (2021). Clustering approaches for customer segmentation: A review. *Expert Systems with Applications*, 167, 114618.

Brown, T., Smith, J., & Lee, K. (2019). Retail sales forecasting in dynamic markets. *Journal of Business Analytics*, 12(3), 205-220.

Chen, L. (2020). Visualización de datos y técnicas predictivas en ambientes comerciales. *Data Science Review*, 14(3), 102-119.

Chollet, F. (2020). *Deep Learning with Python*. Manning Publications.

Few, S. (2009). *Now You See It: Simple Visualization Techniques for Quantitative Analysis*. Analytics Press.

García, P., Rodríguez, M., & Sanchez, L. (2020). Data preprocessing techniques for improved machine learning models. *Data Science Journal*, 18(2), 115-128.

García, R. (2021). Ciencia de datos aplicada a la mejora de procesos comerciales. *International Journal of Data Science*, 10(1), 77-95.

García, R., & Martínez, L. (2020). Modelos predictivos en el comercio minorista: Aplicaciones y desafíos. *Journal of Data Science*, 18(3), 245-267.

- Han, J., Pei, J., & Kamber, M. (2011). *Data mining: Concepts and techniques*. Elsevier.
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.
- Hyndman, R. J., & Athanasopoulos, G. (2018). *Forecasting: Principles and Practice* (2^a ed.). OTexts.
- Inmon, W. H. (2005). *Building the Data Warehouse*. John Wiley & Sons.
- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer.
- Kim, S., & Park, J. (2020). Machine learning techniques for consumer behavior prediction. *Journal of Retailing and Consumer Services*, 54, 102035.
- Kuhn, M., & Johnson, K. (2013). *Applied predictive modeling*. Springer.
- Kumar, A., & Patel, S. (2021). Feature selection and dimensionality reduction in predictive modeling. *International Journal of Data Science*, 9(1), 45-60.
- Kumar, S., Patel, R., & Singh, A. (2022). Data cleaning and its impact on machine learning models. *Journal of Computational Methods*, 5(4), 233-250.
- Lopez, M., Rodríguez, P., & Sánchez, A. (2021). Análisis y segmentación de clientes en cadenas comerciales: Un enfoque basado en clustering. *Data Mining and Knowledge Discovery Journal*, 35(5), 1123-1140.
- López, M., & Martínez, J. (2019). Estrategias de análisis predictivo en la industria retail. *Retail Analytics Journal*, 7(2), 33-50.
- Martínez, E., & Gómez, F. (2017). The role of feature engineering in predictive analytics. *IEEE Transactions on Knowledge and Data Engineering*, 29(4), 837-850.
- Martínez, F., & Ortega, D. (2022). Validación de modelos predictivos en entornos de datos complejos. *Journal of Predictive Analytics*, 9(1), 58-75.
- Miller, A., & Reinsel, D. (2019). Big data and cloud: Computing the future. *Journal of Cloud Computing*, 8(1), 1-12.
- Poole, C., & Martin, H. (2021). Advanced data preprocessing strategies in retail

analytics. *Data & Knowledge Engineering*, 133, 102042.

Provost, F., & Fawcett, T. (2013). *Data science for business: What you need to know about data mining and data-analytic thinking*. O'Reilly Media.

Ramirez, E. (2020). Técnicas de preprocesamiento de datos utilizando Python. *Data Engineering Today*, 6(3), 40-55.

Ramirez, J., & Torres, C. (2022). Data pipelines and their application in retail predictive modeling. *IEEE Transactions on Knowledge and Data Engineering*, 34(2), 670–683.

Ramirez, J., & Torres, L. (2020). Advances in machine learning for sales forecasting. *Journal of Retail Analytics*, 5(2), 89-105.

Rojas, E., & Martínez, L. (2019). Aplicaciones de la modelización de series temporales en el sector retail. *Ingeniería de Datos*, 8(1), 75-92.

Sánchez, R., & Herrera, D. (2019). Clustering techniques for enhanced predictive models in retail. *Expert Systems with Applications*, 128, 156-167.

Saysana, R. (2020). Métodos de análisis de series temporales en entornos reales. *Revista de Estadística Aplicada*, 12(2), 45-63.

Shmueli, G., Bruce, P., Gedeck, P., & Patel, N. R. (2020). *Data mining for business analytics: Concepts, techniques, and applications in R*. John Wiley & Sons.

Smith, A. (2020). Effective communication in data-driven organizations. *International Journal of Business Communication*, 57(4), 450–468.

Smith, J. (2020). Machine learning applications in retail management. *Journal of Machine Learning in Business*, 11(2), 85-103.

Taylor, S. J., & Letham, B. (2018). Forecasting at Scale. *The American Statistician*, 72(1), 37–45.

Tukey, J. W. (1977). *Exploratory Data Analysis*. Addison-Wesley.

Zhang, G. (2021). Data Preprocessing in Predictive Modeling: Techniques and Applications. *Journal of Data Science*, 19(3), 283-300.

Zhang, L., & Zhao, M. (2020). Big data applications in the retail industry: A comprehensive review. *Information Systems Frontiers*, 22(2), 369–389.