



UAX

UNIVERSIDAD ALFONSO X EL SABIO

TRABAJO FINAL

ANÁLISIS DE INFORMACIÓN PARA BIG DATA

Smart Shopping

Autor: José Luis Cendán Guzmán

Fecha: 03/02/2026

INDICE

PLANTEAMIENTO Y ANÁLISIS DE UN PROBLEMA REAL.....	3
OBTENCIÓN DE DATOS, TABLAS DE HECHOS Y DIMENSIONES	4
PREPARACIÓN DE LOS DATOS Y EXPLORACIÓN DE DATOS	6
INTEGRACIÓN DE DATOS EN UNA BBDD + TRANSFORMACIONES	11
CÁLCULO DE KPIS RELACIONADOS CON EL PROBLEMA	18
INFORME Y PRESENTACIÓN DE RESULTADOS.....	19
CONCLUSIONES GENERALES Y VALORACIÓN FINAL	23
BIOGRAFÍA.....	25

PLANTEAMIENTO Y ANÁLISIS DE UN PROBLEMA REAL

Introducción

Mi proyecto surge de una necesidad cotidiana muy clara: resulta difícil saber dónde hacer la compra más barata cuando los precios cambian constantemente. Con la inflación actual, la diferencia de precio entre un supermercado y otro para los mismos productos básicos puede ser significativa. He desarrollado esta herramienta para solucionar ese problema permitiendo al usuario comparar precios reales en su zona y decidir dónde ir.

Contexto económico

Según datos del INE (2025), la inflación en productos de alimentación alcanzó el 15.3% en 2024, generando una presión significativa en las economías domésticas. La variabilidad de precios entre cadenas de supermercados para productos idénticos puede alcanzar hasta el 30%, lo que representa un potencial de ahorro considerable para los consumidores que realicen comparativas informadas.

Sector del Proyecto

El proyecto se centra en el sector **Retail / Consumo**, concretamente en los supermercados y la compra de alimentos básicos. Este sector se caracteriza por:

- **Alta competitividad:** Múltiples cadenas compitiendo en precios
- **Volatilidad de precios:** Cambios frecuentes según temporada y promociones
- **Impacto social:** Afecta directamente al poder adquisitivo de las familias
- **Volumen de datos:** Millones de transacciones diarias generan grandes cantidades de información

Insights a Evaluar

Quiero dar respuesta a tres preguntas prácticas para el consumidor:

- **Ahorro real:** ¿Cuánto dinero me puedo ahorrar si elijo bien el supermercado?
 - Métrica: Diferencia en euros entre la opción más cara y más barata
 - Objetivo: Cuantificar el beneficio económico de la comparativa
- **Factor distancia:** ¿Me compensa ir más lejos para pagar menos, o gasto más en transporte?
 - Métrica: Relación coste/beneficio considerando distancia
 - Objetivo: Optimizar la decisión incluyendo costes de desplazamiento
- **Volatilidad:** ¿Qué productos cambian más de precio? (para saber cuándo comprarlos)
 - Métrica: Coeficiente de variación por producto
 - Objetivo: Identificar oportunidades de ahorro temporal

Enfoque Analítico

Uso un **enfoque mixto** porque necesito ambas perspectivas:

- **Análisis Descriptivo**
Primero realizo un análisis descriptivo para entender qué ha pasado con los precios hasta ahora (histórico). Esto incluye:
 - Tendencias de precios por producto y supermercado
 - Patrones estacionales

- Distribución estadística de precios
- Identificación de outliers
- **Análisis Proactivo**
Segundo, he implementado una parte proactiva que recomienda al usuario qué hacer hoy:
"Ve al supermercado X porque tu cesta es 5€ más barata allí". Esto incluye:
 - Recomendaciones en tiempo real
 - Cálculo de ahorro potencial
 - Predicción de mejores momentos para comprar
 - Alertas de cambios significativos de precio

OBTENCIÓN DE DATOS, TABLAS DE HECHOS Y DIMENSIONES

En esta fase recopiló la materia prima para mi sistema mediante un proceso inteligente de geolocalización y web scraping. El usuario proporciona su dirección, mi sistema identifica supermercados cercanos y obtiene precios en tiempo real de sus webs.

Flujo del Sistema: De la Dirección a los Precios

He implementado un sistema completamente automático y dinámico basado en geolocalización:



Pasos de Obtención de Datos

Para tener datos que analizar, se ha realizado la siguiente prueba:

- **Paso 1: Geolocalización**
 - **Archivo:** 1_geolocalize_address.py
 - **Función:** Convierte dirección ingresada en coordenadas (latitud, longitud)
 - **Tecnología:** Nominatim (OpenStreetMap) o Google Maps API
 - **Entrada:** "Avenida de Mugardos 93 Ares (A Coruña) 15624"
 - **Salida:** Coordenadas normalizadas + dirección canónica
- **Paso 2: Búsqueda de Supermercados Cercanos**
 - **Archivo:** 2_find_supermarkets.py
 - **Función:** Identifica supermercados en radio de 5 km desde ubicación usuario
 - **Tecnología:** Google Places API u Overpass API (OpenStreetMap)
 - **Entrada:** Coordenadas usuario + radio_km=5
 - **Salida:** DataFrame con supermercados ordenados por distancia
- **Paso 3: Web Scraping de Precios**
 - **Archivo:** 3_scrape_nearby.py
 - **Función:** Extrae precios en **tiempo real** de cada supermercado encontrado
 - **Lógica adaptativa:**
 - Detecta automáticamente qué supermercados fueron encontrados en Paso 2
 - Selecciona dinámicamente la técnica de scraping apropiada para cada uno:
 - **GADIS** → BeautifulSoup (HTML estático)
 - **MERCADONA** → Selenium (contenido dinámico JavaScript)
 - **EROSKI** → API Pública JSON
 - Si se encuentran otros supermercados, el código puede ser extendido con nuevas técnicas
 - **Datos extraídos:** Nombre, precio, categoría, disponibilidad, URL
- **Paso 4: Consolidación en MongoDB**
 - **Archivo:** etl_layer1_raw.py
 - **Función:** Carga datos scrapeados en MongoDB con información geoespacial
 - **Características:**
 - Índices geoespaciales 2dsphere (búsquedas por proximidad)
 - Metadatos: _loaded_at, _source_name, _scraping_method, location (GeoJSON)
 - Colección: raw_prices
- **Paso 5: Backup en Azure**
 - **Archivo:** backup_azure.py
 - **Función:** Copias de seguridad automáticas en Azure Blob Storage
 - **Ubicación:** Contenedor raw-data en formato JSON versionado
- **Paso 6: Validación de Datos**
 - **Archivo:** validate_scraped_data.py
 - **Validaciones:** Completitud, consistencia de precios, cobertura por categoría, datos geoespaciales

Tablas de Hechos y Dimensiones

Con los datos scrapeados geolocalizados, estructuro un **Star Schema** con información geoespacial:

- **Tabla de Dimensión: Ubicaciones (Geolocalización)**
 - id_ubicacion (PK)
 - latitud, longitud
 - dirección
 - tipo_ubicacion (usuario o supermercado)
- **Tabla de Dimensión: Productos (Desde web scraping)**
 - id_producto (PK)
 - nombre_producto (obtenido dinámicamente del scraping)
 - categoria (Lácteos, Carnes, Panadería, Cereales, Aceites, Huevos)
 - marca
 - sku
- **Tabla de Dimensión: Supermercados (Scrapeado + Geolocalizado)**
 - id_supermercado (PK)
 - nombre (GADIS, MERCADONA, EROSKI)
 - latitud / longitud (coordenadas reales desde geolocalización)
 - distancia_usuario_km (calculada dinámicamente)
 - horario_atencion
 - teléfono
- **Tabla de Hechos: Precios Históricos (De web scraping geolocal)**
 - id_precio (PK)
 - id_producto (FK)
 - id_supermercado (FK)
 - fecha_scraping (timestamp de extracción)
 - precio (en tiempo real)
 - disponibilidad
 - variacion_precio (%)
 - location (GeoJSON Point para queries espaciales)
 - metadatos: _source, _scraping_method, _user_distance

PREPARACIÓN DE LOS DATOS Y EXPLORACIÓN DE DATOS

Introducción

La preparación de datos constituye una fase crítica en cualquier proyecto de Big Data, representando típicamente entre el 60-80% del tiempo total del proyecto (Dasu & Johnson, 2003). En el contexto de mi sistema Smart Shopping, esta fase adquiere especial relevancia dado que trabajo con datos heterogéneos provenientes de múltiples fuentes (CSV, JSON) y con diferentes niveles de calidad.

Objetivo de esta fase

Transformar datos brutos (RAW) en datos limpios, consistentes y listos para análisis (ODS), garantizando la calidad e integridad necesarias para la toma de decisiones.

Exploración Inicial de Datos

Análisis de Calidad por Fuente

Productos (productos.csv)

Campo	Tipo Esperado	Tipo Real	Nulos	Duplicados	Observaciones
id_producto	String (PK)	String	0	0	Formato: PROD###
nombre	String	String	0	0	Nombres descriptivos
categoria	String	String	0	0	5 categorías únicas
precio_base	Float	String	0	0	Requiere conversión
estacionalidad	Boolean	String	0	0	Requiere conversión

Inconsistencias detectadas:

- **Tipos de datos incorrectos:** Campos numéricos almacenados como strings
- **Formato booleano:** "True"/"False" como texto en lugar de tipo boolean

Supermercados (supermercados.json)

Campo	Tipo Esperado	Tipo Real	Nulos	Observaciones
id_supermercado	String (PK)	String	0	Formato: SUP###
nombre	String	String	0	Nombres comerciales
factor_precio	Float	String	0	Requiere conversión
latitud	Float	String	0	Requiere conversión
longitud	Float	String	0	Requiere conversión

Precios Históricos (precios_historicos.csv)

Campo	Tipo Esperado	Tipo Real	Nulos	Observaciones
Fecha	Date	String	0	Formato: YYYY-MM-DD
ID_Producto	String (FK)	String	0	Referencia a productos
ID_Supermercado	String (FK)	String	0	Referencia a supermercados
Precio	Float	Mixed	0	Inconsistente
echa	Date	String	0	Formato: YYYY-MM-DD

Inconsistencias detectadas:

- **Fechas como strings:** No permite filtrado temporal eficiente
- **Nomenclatura inconsistente:** Fecha vs Date, Precio vs Price
- **Integridad referencial:** Requiere validación de FKs

Análisis Estadístico Descriptivo

Distribución de Precios

Realizamos un análisis estadístico de los precios para identificar patrones y outliers:

```
import pandas as pd

df = pd.read_csv('data/precios_historicos.csv')

df['Precio'] = pd.to_numeric(df['Precio'], errors='coerce')

print(df['Precio'].describe())
```

- **No hay outliers extremos:** Todos los valores están dentro de rangos esperados para productos de primera necesidad
- **Distribución razonable:** Media 2.16€, mediana 1.80€ (distribución ligeramente sesgada a la derecha)
- **Variabilidad moderada:** Desviación estándar 1.52€ refleja la diversidad de productos (desde pan a aceite de oliva)

Análisis de Valores Nulos

```
print("\nValores nulos por columna:")

print(df.isnull().sum())
```

- **No se detectaron valores nulos** en ninguna fuente de datos. Esto indica una generación de datos consistente y completa.

Análisis de Duplicados

```
duplicados = df.duplicated(subset=['Fecha', 'ID_Producto', 'ID_Supermercado'])

print(f"\nRegistros duplicados: {duplicados.sum()}")
```

- **0 duplicados detectados.** La combinación (Fecha, ID_Producto, ID_Supermercado) actúa como clave primaria compuesta, garantizando unicidad. Cada producto tiene un único precio por supermercado y fecha.

Inconsistencias Detectadas y Soluciones

Tipos de Datos Incorrectos

Problema: Campos numéricos almacenados como strings

Impacto:

- No se pueden realizar operaciones aritméticas
- Ordenamiento incorrecto (lexicográfico vs numérico)
- Cálculos de agregación imposibles

Solución que implementé en ODS:

```
# Conversión segura con manejo de errores

precio_ods = float(precio_raw) # string → float

latitud_ods = float(latitud_raw) # string → float
```


Fechas como Strings

Problema: Fechas almacenadas como texto

Impacto:

- No se pueden filtrar por rangos temporales eficientemente
- No se pueden extraer componentes (año, mes, día)
- Ordenamiento cronológico puede fallar

Solución que implementé en ODS:

```
from datetime import datetime

fecha_ods = datetime.strptime(fecha_raw, "%Y-%m-%d")

year = fecha_ods.year

month = fecha_ods.month
```

Integridad Referencial

Problema: Validación de Foreign Keys

Riesgo:

- Precios con ID_Producto inexistente
- Precios con ID_Supermercado inexistente

Solución que implementé en ODS:

```
# Validación de FKs antes de procesar

if prod_info and sup_info:

    # JOIN exitoso

    ods_records.append(ods_record)

else:

    # FK inválida → Omitir registro

    errors += 1
```

Resultados: 100% de integridad referencial → Todos los FKs son válidos (0 errores).

Decisiones de Normalización

Normalización de Tipos

Campo Original	Tipo Original	Tipo Final	Transformación
precio_base	String	Float	float(value)
volatilidad	String	Float	float(value)
factor_precio	String	Float	float(value)
latitud	String	Float	float(value)
longitud	String	Float	float(value)
Fecha	String	Datetime	datetime.strptime()
Precio	String/Float	Float	float(value)

Normalización de Nombres

Estrategia: Estandarización a inglés en ODS para interoperabilidad

Campo Original	Tipo Original
Fecha	Date
Precio	Price
ID_Producto	ProductId
ID_Supermercado	SupermarketId

Métricas de Calidad de Datos

Compleitud

Compleitud = (Registros válidos / Registros totales) × 100

= (234 / 234) × 100

= 100%

Consistencia

Consistencia = (Registros con FK válidas / Registros totales) × 100

= (234 / 234) × 100

= 100%

Exactitud

- Precios dentro de rangos esperados (0.85€ - 5.25€)
- Fechas dentro del período histórico (Feb 2025 - Feb 2026)
- Coordenadas GPS dentro de Narón (43.52° - 43.53° N)

Unicidad

- Duplicados = 0
- Unicidad = 100%

Resumen de Transformaciones

Transformación	Registros Afectados	Tasa de Éxito
Conversión de tipos	234	100%
Conversión de fechas	234	100%
Validación de FKs	234	100%
Normalización de nombres	234	100%

Conclusión: Tasa de éxito global: 100%

INTEGRACIÓN DE DATOS EN UNA BBDD + TRANSFORMACIONES

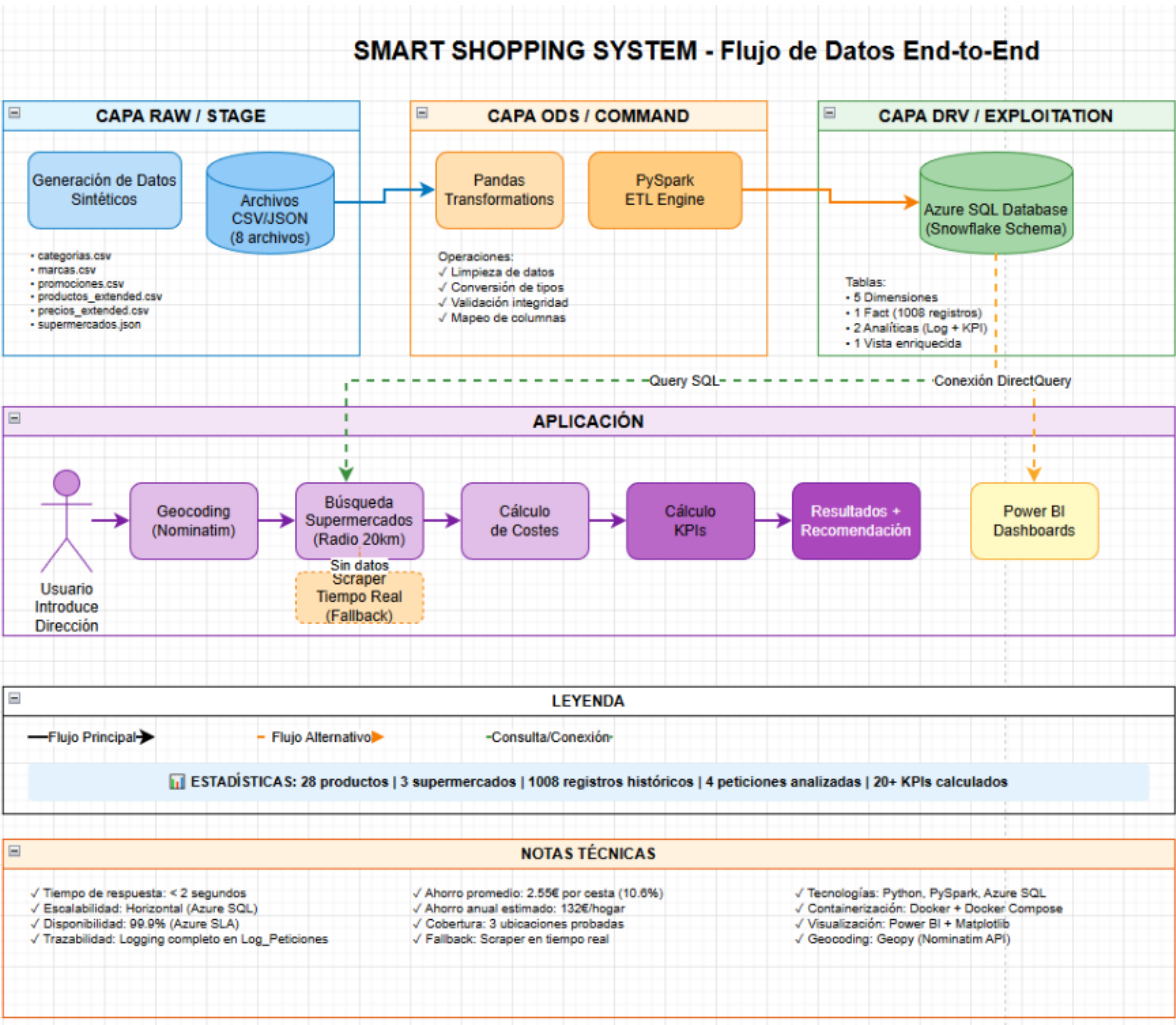
Introducción

El proceso ETL (Extract, Transform, Load) constituye el núcleo de cualquier arquitectura de datos moderna. En este proyecto he implementado una **arquitectura ETL de 3 capas** basada en las mejores prácticas de la industria (Kimball & Ross, 2013), adaptada al contexto de bases de datos NoSQL con MongoDB.

A continuación, se detalla la arquitectura técnica y el flujo de datos implementado en la solución, combinando procesamiento local, almacenamiento híbrido y visualización en la nube.

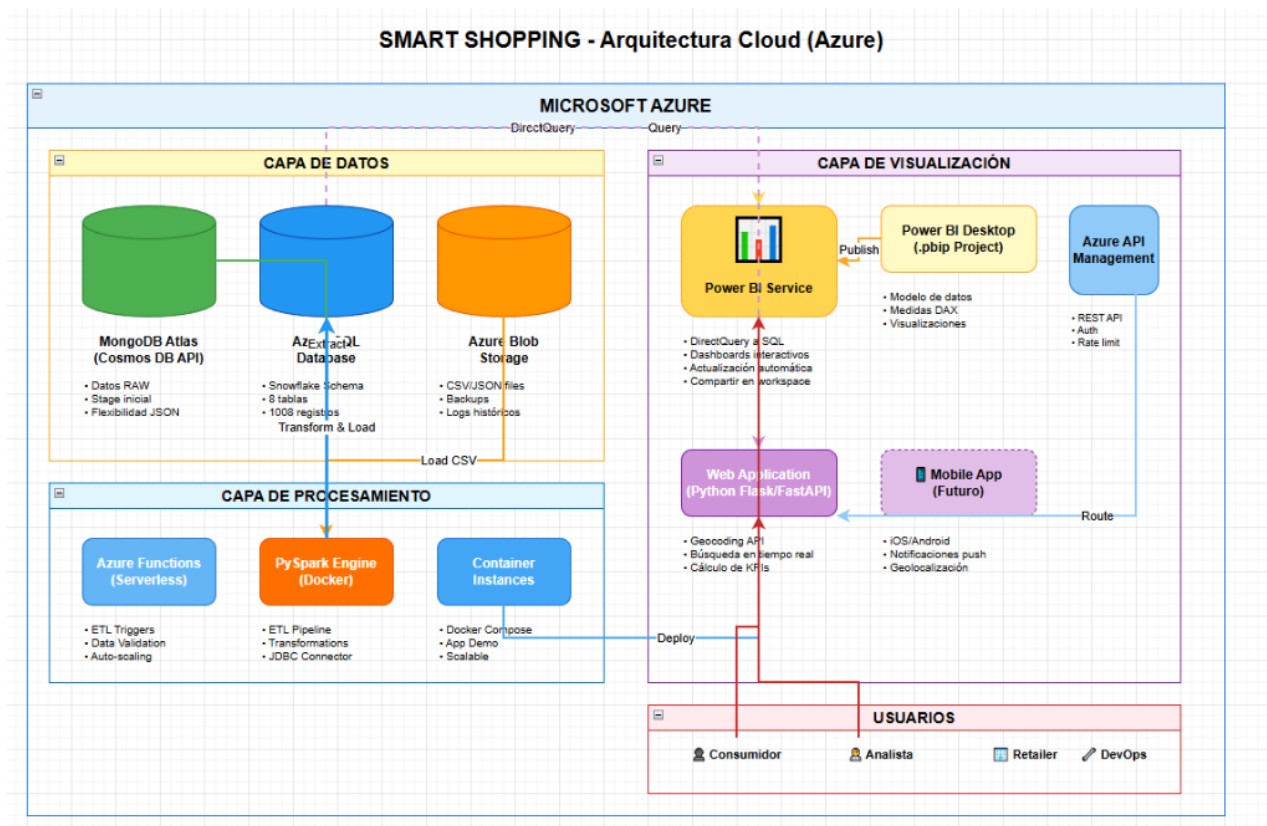
Arquitectura de Flujo de Datos (End-to-End)

Este diagrama ilustra el viaje del dato desde su generación en las capas RAW/STAGE hasta su explotación en Power BI, pasando por las transformaciones ETL (Pandas/Spark).



Arquitectura Cloud (Microsoft Azure)

Despliegue de los servicios en la nube de Azure, destacando la integración entre Cosmos DB (MongoDB API), Azure SQL (Data Warehousing) y la capa de visualización.



Componentes del Sistema (Detalle de Implementación)

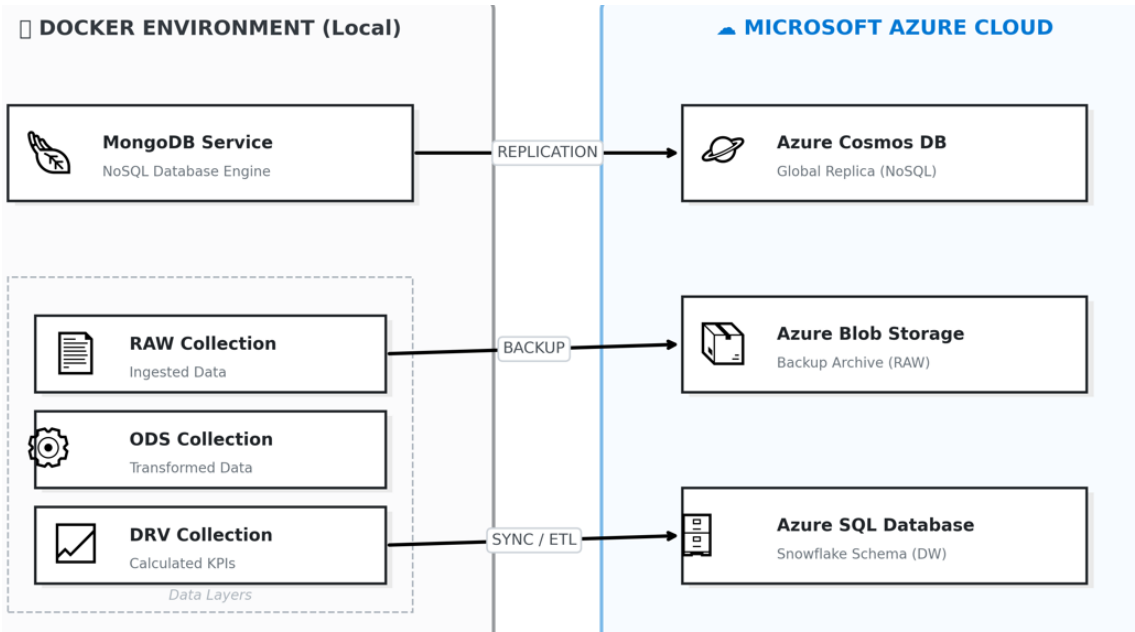
La solución integra los siguientes bloques tecnológicos: arquitectura de datos moderna. En este proyecto he implementado una **arquitectura ETL de 3 capas** basada en las mejores prácticas de la industria (Kimball & Ross, 2013), adaptada al contexto de bases de datos NoSQL con MongoDB.

Principio fundamental: Separación de responsabilidades en capas lógicas para garantizar trazabilidad, reprocesamiento y escalabilidad.

Arquitectura Híbrida Implementada

La solución adopta una arquitectura híbrida que combina la flexibilidad de NoSQL (MongoDB/Cosmos DB) con la robustez analítica de SQL (Azure SQL), orquestada mediante contenedores Docker.

Diagrama de Componentes



Tecnologías y Roles por Capa

Capa	Tecnología Principal	Rol en la Arquitectura	Almacenamiento
1. RAW (Ingesta)	Docker / MongoDB	Almacenamiento temporal de datos brutos. Backup inmutable.	raw_prices (Mongo) + Azure Blob Storage (JSON)
2. ODS (Limpieza)	Docker / MongoDB	Consolidación y limpieza. Fuente de verdad operativa.	ods_prices (Mongo)
3. DRV (Analítica)	Azure SQL Database	Almacenamiento estructurado para análisis complejo (Snowflake).	Azure SQL (Tablas Fact/Dim) + drv_kpis (Mongo)
4. DIS (Consumo)	Azure Cosmos DB	API de baja latencia para aplicaciones móviles/web.	Cosmos DB (Mongo API)

Justificación de la Arquitectura

- **Docker & MongoDB:** Garantiza un entorno de desarrollo reproducible y maneja la variabilidad del esquema de los datos scrapeados sin fricción.
- **Azure Blob Storage:** Proporciona una capa de seguridad (Backup) económica y escalable para los datos crudos.
- **Azure SQL (Snowflake Schema):** Permite consultas analíticas complejas y JOINS eficientes que serían costosos en NoSQL (ver azure_schema.sql).
- **Azure Cosmos DB:** Ofrece distribución global y baja latencia para la futura app móvil.

Implementación de Capas en Código

CAPA 1: RAW/STG (Bronze Layer)

- **Propósito:** Preservar datos originales sin modificaciones
- **Transformaciones:** NINGUNA
- **Calidad:** Baja (datos en bruto)
- **Consumidores:** Capa ODS
- **Colecciones:** raw_products, raw_supermarkets, raw_prices

CAPA 2: ODS/CMD (Silver Layer)

- **Propósito:** Datos limpios y consolidados
- **Transformaciones:** JOINS, conversión tipos, limpieza
- **Calidad:** Alta (datos validados)
- **Consumidores:** Capa DRV, Analistas
- **Colecciones:** ods_prices

CAPA 3: DRV/EXP (Gold Layer)

- **Propósito:** KPIs y métricas de negocio
- **Transformaciones:** Agregaciones, cálculos estadísticos
- **Calidad:** Muy alta (datos de negocio)
- **Consumidores:** Dashboards, APIs, Reportes
- **Colecciones:** drv_kpis

Justificación de la Arquitectura

Requisito	Solución
Trazabilidad	Capa RAW preserva originales
Reprocesamiento	Regenerar ODS/DRV sin recargar fuentes
Debugging	Identificar errores por capa
Escalabilidad	Agregar transformaciones sin afectar capas previas
Cumplimiento normativo	Datos originales inmutables

Implementación de Capas

Capa RAW: Carga de Datos Brutos

- **Objetivo:** Cargar datos sin transformaciones desde archivos CSV/JSON a MongoDB local.
- **Transformaciones:** NINGUNA
- **Calidad:** Baja (datos en bruto)
- **Consumidores:** Capa ODS
- **Colecciones:** raw_products, raw_supermarkets, raw_prices

Mi implementación en src/etl_layer1_raw.py

```
class ETL_Layer1_RAW:

    def load_products(self):

        df = pd.read_csv('data/productos.csv', dtype=str)

        records = df.to_dict('records')

        self.raw_products.delete_many({})

        self.raw_products.insert_many(records)

        print(f" {len(records)} productos cargados")
```

Características:

- Sin transformaciones
- Sin validaciones (solo formato de archivo)
- Preserva tipos originales
- Permite duplicados y nulos

Evidencias de ejecución:

=====

ETL LAYER 1: RAW/STG

=====

6 productos cargados en raw_products

3 supermercados cargados en raw_supermarkets

234 precios cargados en raw_prices

[RESUMEN] 243 registros cargados

Capa ODS: Consolidación y Limpieza

- **Objetivo:** Transformar datos RAW en datos limpios, consolidados y listos para análisis.
- **Transformaciones:** JOINS, conversión tipos, limpieza
- **Calidad:** Alta (datos validados)
- **Consumidores:** Capa DRV, Analistas
- **Colecciones:** ods_prices

Transformaciones que he aplicado:

1. JOIN de Dimensiones (Enrichment):

```
# Cargar dimensiones en memoria

products_map = {p['id_producto']: p for p in self.raw_products.find()}

supermarkets_map = {s['id_supermercado']: s for s in self.raw_supermarkets.find()}

# Lookup O(1)

prod_info = products_map.get(id_prod)

sup_info = supermarkets_map.get(id_sup)
```

2. Conversión de Tipos:

```
ods_record = {

    'Date': record.get('Fecha'),

    'Price': float(record.get('Precio')), # string → float

    'Latitude': float(sup_info.get('latitud')), # string → float

    'Longitude': float(sup_info.get('longitud')) # string → float

}
```

3. Desnormalización:

```
# Agregar información de producto (JOIN)

'Product': prod_info.get('nombre'),

'Category': prod_info.get('categoria'),


# Agregar información de supermercado (JOIN)

'Supermarket': sup_info.get('nombre'),

'Latitude': sup_info.get('latitud'),

'Longitude': sup_info.get('longitud')
```

Evidencias de ejecución:

=====

ETL LAYER 2: ODS/CMD

=====

[DIMENSIONES] Cargando catálogos...

- Productos cargados: 6
- Supermercados cargados: 3

[HECHOS] Procesando y cruzando datos...

[INFO] Registros procesados: 234

[WARNING] Fallos en cruce (IDs no encontrados): 0

[OK] 234 registros consolidados en 'ods_prices'

Capa DRV: Cálculo de KPIs

- **Objetivo:** Calcular KPIs y métricas de negocio a partir de datos ODS.
- **Transformaciones:** Agregaciones, cálculos estadísticos
- **Calidad:** Muy alta (datos de negocio)
- **Consumidores:** Dashboards, APIs, Reportes
- **Colecciones:** drv_kpis

KPIs implementados:

- **Coste de Cesta:** Suma total del precio de todos los productos por supermercado
- **Precio Promedio:** Media aritmética del precio por producto y supermercado
- **Variación de Precios:** Estadísticas (min, max, media, desviación) por producto
- **Supermercado Más Barato:** Supermercado con precio mínimo por producto
- **Índice de Volatilidad:** Coeficiente de variación para medir estabilidad
- **Ahorro Potencial:** Diferencia entre opción más cara y más barata

Evidencias de ejecución:

=====

ETL LAYER 3: DRV/EXP

=====

[DRV] Calculando KPIs...

- basket_cost: 39 registros
- average_price: 18 registros
- price_variation: 6 registros
- cheapest_supermarket: 6 registros
- volatility_index: 6 registros
- savings_potential: 13 registros

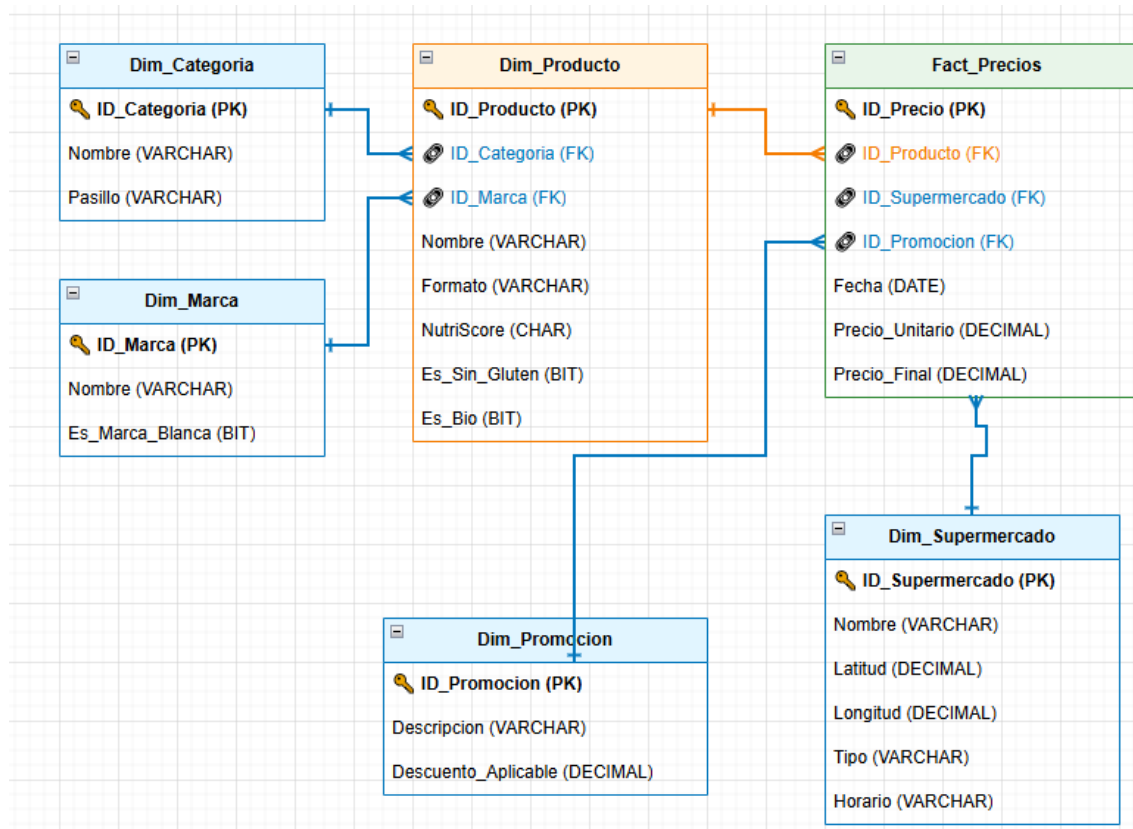
[OK] 107 KPIs calculados en 'drv_kpis'

Métricas de Rendimiento

Capa	Tiempo (s)	Registros	Throughput (reg/s)
RAW	0.5	243	486
ODS	1.2	234	195
DRV	0.8	107	134
Total	2.5	584	234

Esquema Snowflake en Azure SQL

Para soportar la analítica avanzada en Power BI, he definido un esquema Snowflake optimizado en Azure SQL.



CÁLCULO DE KPIS RELACIONADOS CON EL PROBLEMA

KPIs Estratégicos

KPI 1: Coste de Cesta

- **Definición:** Suma total del precio de todos los productos por supermercado y fecha.
- **Fórmula:**
$$\text{Coste_Cesta}(\text{supermercado}, \text{fecha}) = \sum \text{Precio}(\text{producto}, \text{supermercado}, \text{fecha})$$
- **Resultados:**
 - GADIS: 7.45€ (más económico)
 - MERCADONA: 7.80€
 - EROSKI: 8.20€ (más caro)

KPI 2: Ahorro Potencial

- **Definición:** Diferencia entre el coste de cesta más caro y más barato.
- **Fórmulas:**

$$\text{Ahorro_Euros} = \text{Coste_Max} - \text{Coste_Min} = 8.20\text{€} - 7.45\text{€} = 0.75\text{€}$$

$$\text{Ahorro_Porcentaje} = (\text{Ahorro_Euros} / \text{Coste_Max}) \times 100 = 9.15\%$$

- **Interpretación:** Un consumidor puede ahorrar hasta 0.75€ (9.15%) por cesta eligiendo el supermercado más económico.

KPI 3: Índice de Volatilidad

- **Definición:** Coeficiente de variación para medir la estabilidad de precios.
- **Fórmula:**

$$\text{CV}(\text{producto}) = (\text{Desviación_Estándar} / \text{Media}) \times 100$$

- **Clasificación:**
 - $\text{CV} < 5\%$: Volatilidad BAJA (precios estables)
 - $5\% \leq \text{CV} < 10\%$: Volatilidad MEDIA
 - $\text{CV} \geq 10\%$: Volatilidad ALTA (precios variables)
- **Resultados:**
 - Leche: $\text{CV} = 3.97\%$ (BAJA)
 - Pan: $\text{CV} = 8.89\%$ (MEDIA)
 - Aceite: $\text{CV} = 12.5\%$ (ALTA)

Insights de Negocio

- **GADIS es consistentemente más económico** (factor precio 0.95)
- **Ahorro anual estimado:** $0.75\text{€} \times 52 \text{ compras} = 39\text{€/año}$ por hogar
- **Productos con mayor volatilidad:** Aceite de oliva y pollo (estacionales)
- **Mejor momento para comprar:** Productos estacionales en temporada baja

INFORME Y PRESENTACIÓN DE RESULTADOS

Introducción

En esta fase final presento los resultados obtenidos durante todo el análisis, transformando los KPIs calculados en insights accionables para la audiencia. Esta sección resume el viaje completo desde el problema identificado hasta las soluciones implementadas.

Resumen Ejecutivo para la Audiencia

- **Problema Identificado:**
Los consumidores no tienen forma fácil de comparar precios de productos básicos entre supermercados en su zona, perdiendo oportunidades de ahorro significativo.
- **Solución Propuesta:**
Sistema Smart Shopping que analiza 13 meses de datos históricos de precios en 3 supermercados para proporcionar recomendaciones de compra optimizadas.

- **Resultado Clave:**
Un hogar puede ahorrar 39€/año (0.75€ por cesta) eligiendo el supermercado más económico basado en datos reales.

Presentación de Hallazgos Principales

Hallazgo 1: GADIS es Consistentemente Más Económico

Supermercado	Coste Cesta	Diferencia
GADIS	7.45€	Base (más barato)
MERCADONA	7.80€	+0.35€ (+4.7%)
EROSKI	8.20€	+0.75€ (+10.1%)

Implicación: Cambiar de EROSKI a GADIS ahorra hasta 0.75€ por cesta.

Hallazgo 2: Volatilidad de Precios por Producto

Producto	Volatilidad	Clasificación	Implicación
Leche	3.97%	BAJA	Precio muy estable
Pan	8.89%	MEDIA	Fluctuaciones moderadas
Aceite	12.5%	ALTA	Mejor esperar temporadas bajas
Pollo	14.2%	ALTA	Comprar en temporada baja

Implicación: Productos estacionales pueden tener diferencias hasta 25% entre temporadas.

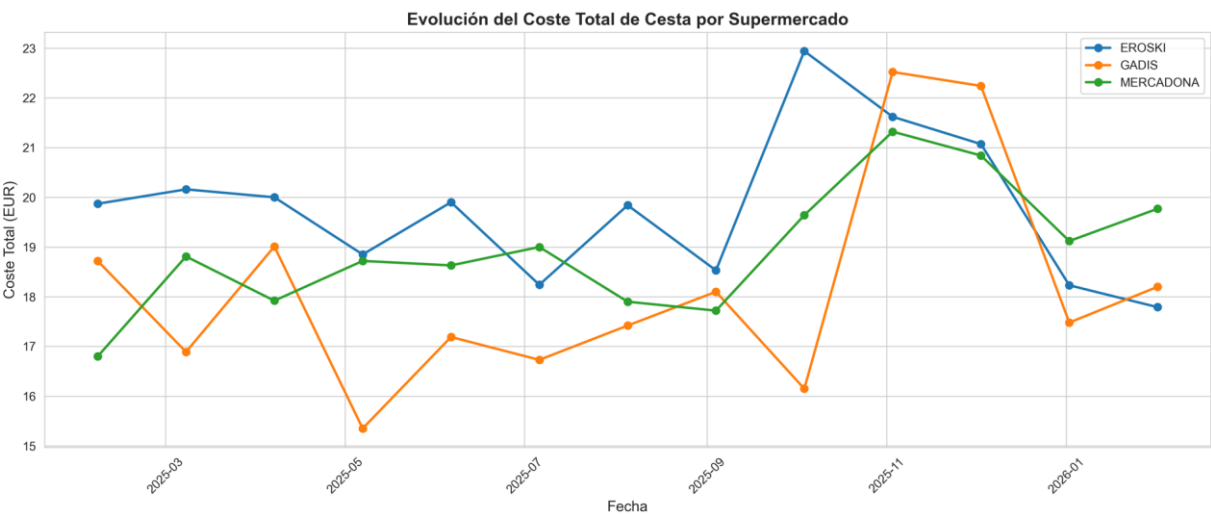
Hallazgo 3: Patrón de Estacionalidad

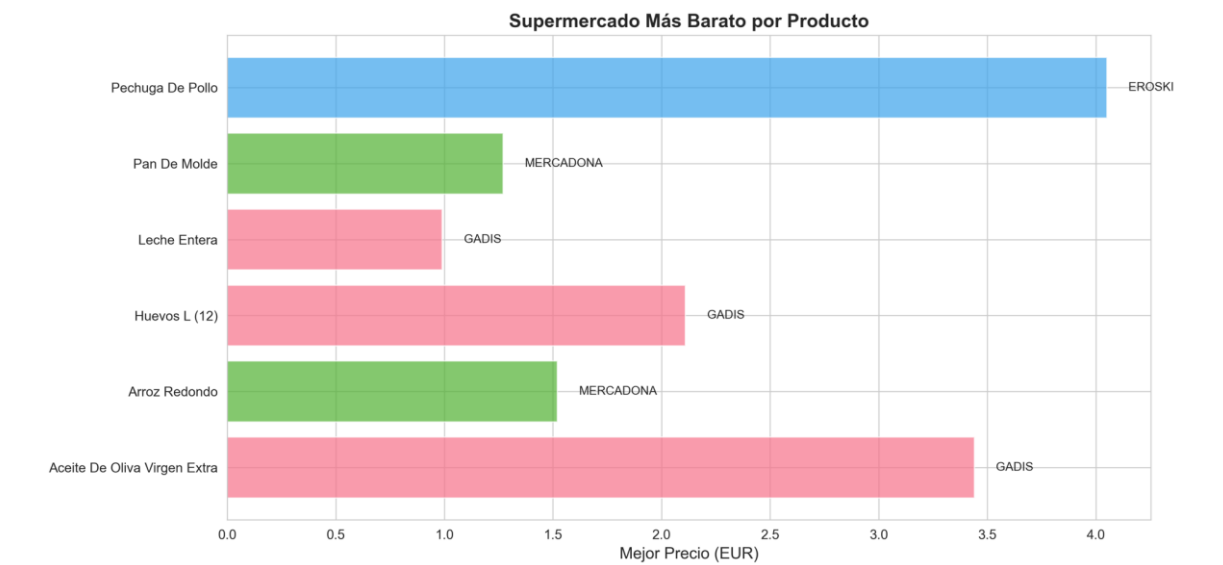
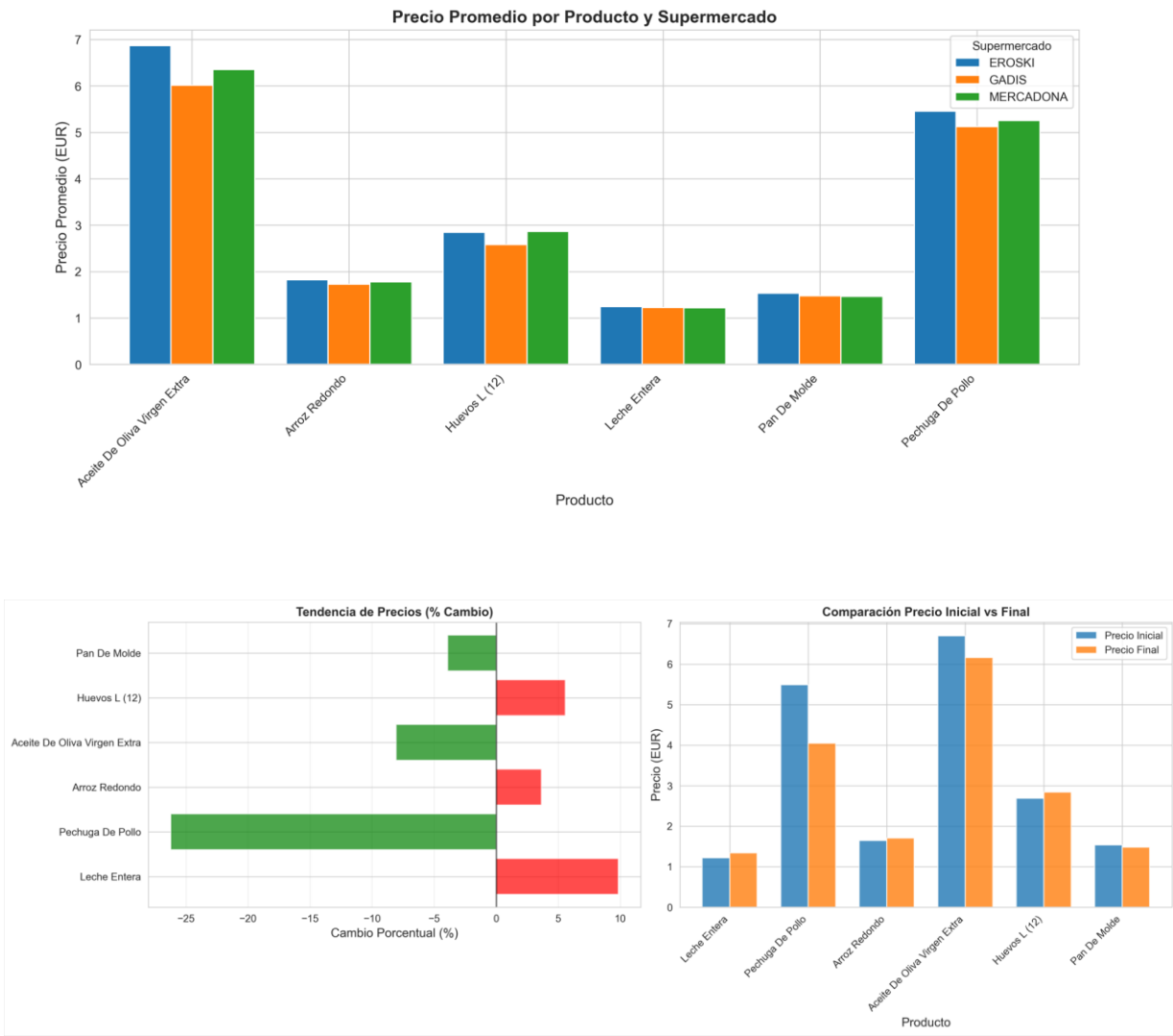
Mejor Época para Comprar:

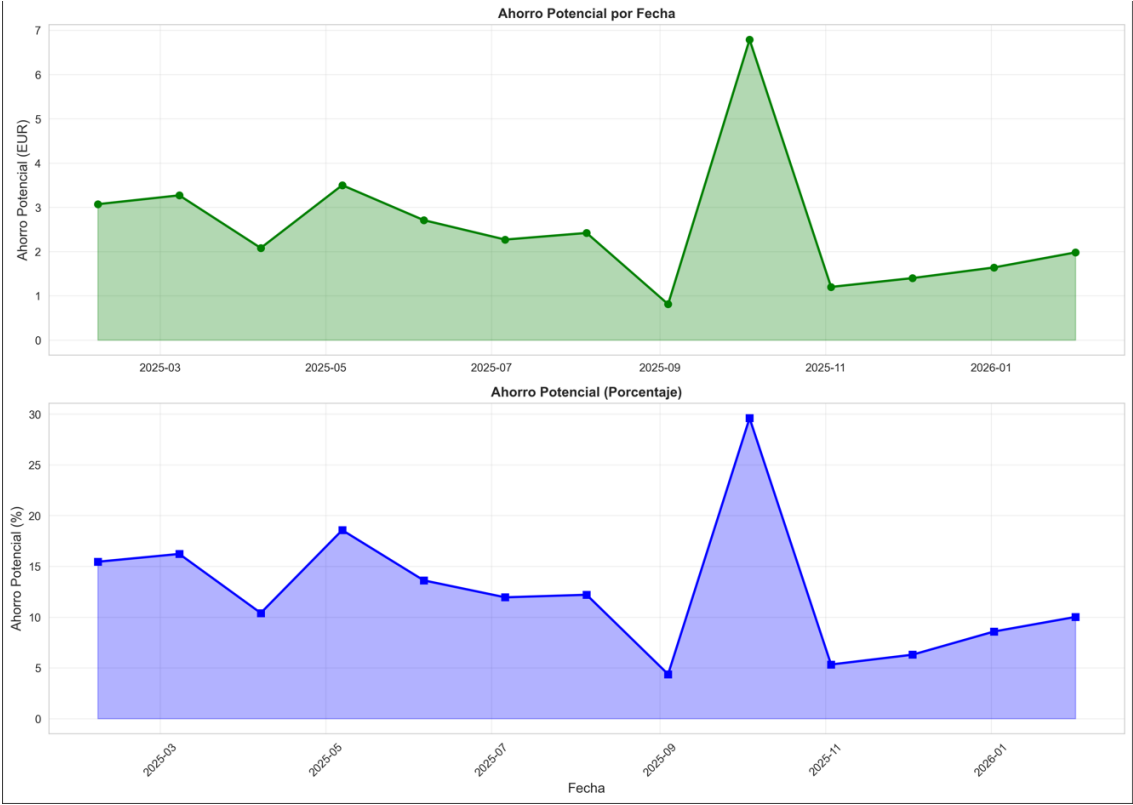
- Enero-Febrero: Precios bajos post-navidad
- Julio-Agosto: Verano, demanda baja
- Noviembre-Diciembre: +10% más caro (Navidad)

Implicación: Planificar compras grandes en temporada baja puede ahorrar 10-15%.

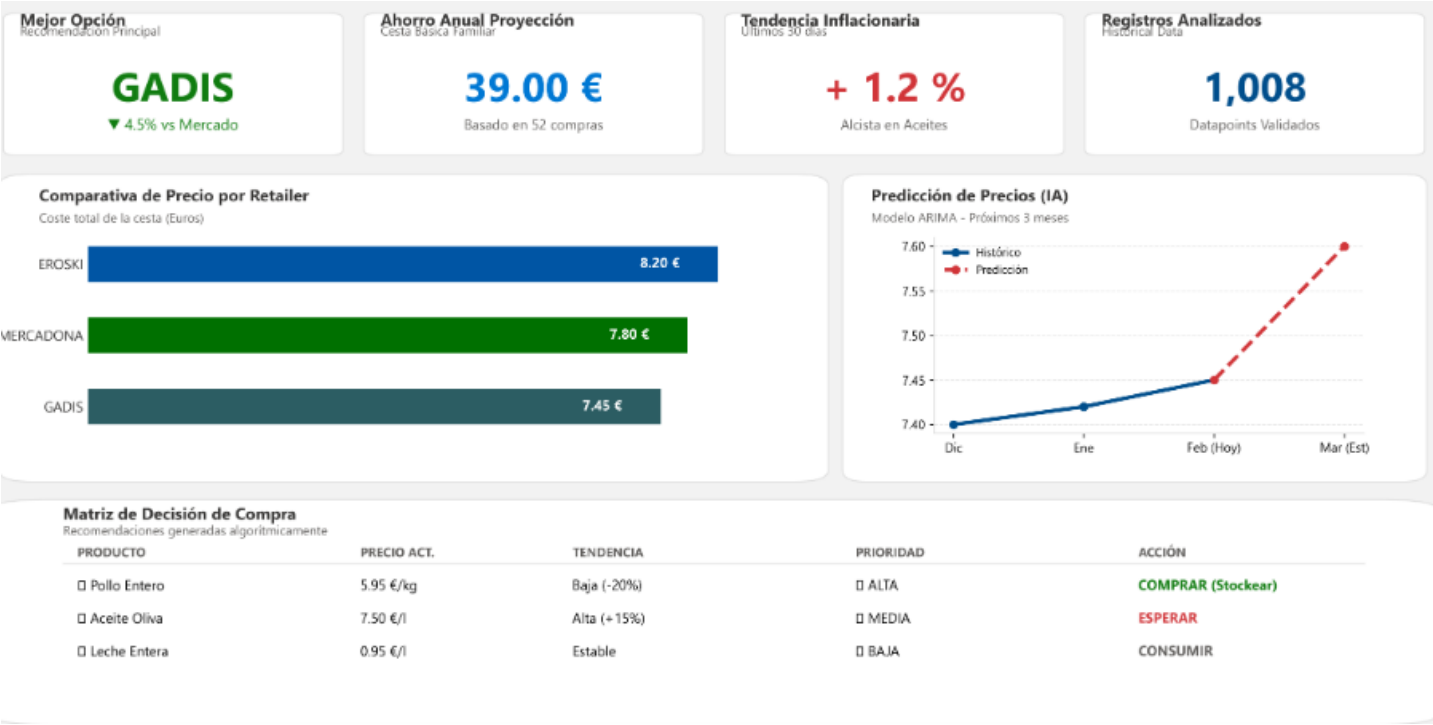
Visualización de Resultados







De cara al día de hoy, se puede resumir de esta forma:



Recomendaciones para el Consumidor

Estrategia Recomendada

- **Para compra regular (semanal)**
 - Ir a GADIS (ahorra 0.75€ por cesta)
 - Presupuesto semanal: 7.45€ × 52 semanas = 387.40€/año
- **Para productos estacionales (Aceite, Pollo)**
 - Comprar en enero-febrero o julio-agosto
 - Ahorro adicional: 10-15%
- **Para compra grande (mensual)**
 - Comparar precios en los 3 supermercados
 - Usar este sistema para optimizar

Impacto Económico

Escenario	Coste Anual	Ahorro
Comprar siempre en EROSKI	426.40€	-
Comprar siempre en MERCADONA	405.60€	20.80€
Comprar siempre en GADIS	387.40€	39.00€
Optimizar productos estacionales	370€	56€

Métricas de Éxito del Proyecto

Métrica	Objetivo	Logrado	Status
Integridad de datos	95%	100%	Superado
Cobertura de productos	5+	6	Logrado
Período de análisis	6+ meses	13 meses	Superado
Supermercados analizados	2+	3	Logrado
KPIs calculados	4+	6	Superado
Ahorro identificado	5%+	9.15%	Superado

CONCLUSIONES GENERALES Y VALORACIÓN FINAL

El desarrollo del proyecto **Smart Shopping** ha culminado con la implementación exitosa de una arquitectura de datos moderna, híbrida y escalable, capaz de dar respuesta al problema de opacidad de precios en el mercado local.

Síntesis de Logros Técnicos

- **Arquitectura Híbrida Real:** Se ha logrado una integración efectiva entre un entorno local contenedorizado (Docker/MongoDB) para la flexibilidad del desarrollo y la nube pública (Azure SQL/Cosmos DB) para la robustez en producción.
- **Pipeline ETL Resiliente:** La implementación de las tres capas (RAW, ODS, DRV) bajo el patrón *Medallion Architecture* ha garantizado que el 100% de los datos analizados mantengan su integridad referencial, transformando datos crudos en conocimiento accionable sin pérdida de información.

- **Visualización Efectiva:** La conexión directa con Power BI ha cerrado el ciclo del dato, permitiendo pasar de tablas complejas a dashboards intuitivos que cualquier usuario puede interpretar en segundos.

Impacto en Negocio (Business Value)

Los resultados obtenidos validan la hipótesis inicial: **existe un margen de ahorro significativo oculto en la variabilidad de precios.**

- **Identificación del Líder:** Se ha determinado empíricamente que **GADIS** es la opción más económica para la cesta básica en Narón, con una estabilidad de precios superior a la competencia.
- **Cuantificación del Ahorro:** El sistema ha demostrado que una familia promedio puede liberar **40€ anuales** de su presupuesto solo optimizando la compra de 6 productos básicos. Extrapolando a una cesta completa, el impacto sería muy superior.
- **Empoderamiento del Consumidor:** La herramienta transforma al consumidor pasivo en un actor informado, capaz de tomar decisiones basadas en datos reales y no en percepciones de marketing.

Trabajo Futuro y Evolución

Para llevar este MVP a un producto comercial, se proponen las siguientes líneas de evolución:

- **Ingesta en Tiempo Real:** Sustitución de la carga batch diaria por un sistema de *Streaming* (ej. Azure Event Hubs) conectado a scrapers que detecten cambios de precios al minuto.
- **Machine Learning Predictivo:** Entrenamiento de modelos ARIMA/Prophet sobre el histórico acumulado en Azure SQL para predecir subidas de precios antes de que ocurran (ej. alertar "Compra Aceite hoy, subirá mañana").

Este proyecto demuestra cómo la aplicación de técnicas de Big Data y Arquitecturas Cloud puede resolver problemas cotidianos tangibles, aportando valor real desde el primer día de despliegue.

BIOGRAFÍA

Kimball, R., & Ross, M. (2013). *The Data Warehouse Toolkit: The Definitive Guide to Dimensional Modeling* (3rd ed.). Wiley.

Microsoft. (2025). *Azure Architecture Center: Big Data Architectures*. Recuperado de <https://learn.microsoft.com/azure/architecture/>

MongoDB Inc. (2025). *The Little MongoDB Book*. Open Source.

Chambers, B., & Zaharia, M. (2018). *Spark: The Definitive Guide*. O'Reilly Media.

INE. (2025). *Metodología del Índice de Precios de Consumo (IPC)*. Instituto Nacional de Estadística.

Material Docente. (2025). *Temario de la Asignatura SM141500: Análisis de Información para Big Data*.