

TECNICATURA UNIVERSITARIA EN INTELIGENCIA ARTIFICIAL

PROCESAMIENTO DEL LENGUAJE NATURAL

TRABAJO PRÁCTICO FINAL

INFORME TÉCNICO: DESARROLLO
DE ASISTENTE VIRTUAL
INTELIGENTE (RAG + AGENTES)



ESTUDIANTE

Malena Ruppen
DNI: 42482175
Legajo: R-4751/7



DOCENTES

Juan Pablo Manson
Alan Geary
Constantino Ferrucci

1. INTRODUCCIÓN

El objetivo de este proyecto fue desarrollar un **asistente conversacional avanzado** capaz de responder consultas sobre productos de una empresa de electrodomésticos, integrando información procedente de:

- **Documentos y reseñas** (texto no estructurado)
- **Datos tabulares** (catálogo y ventas)
- **Relaciones entre entidades** (categorías, productos, marcas)

Para lograrlo se implementó una **arquitectura RAG híbrida**, combinando:

- Una **base vectorial** para búsquedas semánticas en documentos.
- Una **base tabular** consultada mediante filtros generados dinámicamente.
- Una **base de grafos** consultada con Cypher.
- Un **clasificador de intención** que enruta cada consulta.

Un **LLM** que interpreta, genera código, integra información y produce la respuesta final.

En el **Ejercicio 2**, la arquitectura RAG evoluciona hacia un **agente autónomo** basado en el paradigma **ReAct**, capaz de planificar, decidir qué herramienta usar, ejecutar código y combinar resultados.

2. CONSTRUCCIÓN Y JUSTIFICACIÓN DE DECISIONES TÉCNICAS

A continuación, se fundamenta la elección de cada componente clave de la arquitectura:

2.1. Fuente Vectorial (Manuales + Reseñas)

Datos incluidos: Se usaron los manuales de productos provistos por el drive y las reseñas de usuarios. Esto permite cubrir consultas como:

- “¿Cómo uso mi licuadora...?”
- “Opiniones sobre la cafetera...”
- “La batidora hace ruido extraño...”

Modelo de embeddings elegido: intfloat/multilingual-e5-small

Multilingüe: Al tratarse de un asistente en español, era crítico usar un modelo entrenado en múltiples idiomas y no solo en inglés.

Rendimiento/Tamaño: ofrece un buen balance entre calidad de recuperación semántica y velocidad de inferencia. La versión "small" permite ejecutar el sistema en entornos con recursos limitados (como CPU en Colab) sin degradar significativamente la calidad de los vectores.

Text Splitter (Estrategia de Fragmentación)

Se optó por una estrategia híbrida dependiendo del tipo de documento:

- **Para Manuales (.md):** Se utilizó MarkdownHeaderTextSplitter, ya que los manuales tienen una estructura jerárquica fuerte (Título -> Problema -> Solución). Cortar por caracteres arbitrarios rompería el contexto. Este splitter respeta los encabezados, manteniendo la coherencia semántica de cada sección.
- **Para Reseñas (.txt):** Se utilizó RecursiveCharacterTextSplitter. Al ser texto no estructurado y narrativo, esta técnica asegura que los párrafos se mantengan juntos intentando cortar por signos de puntuación, preservando el sentido de la opinión del usuario.

Vector store: ChromaDB

Para el almacenamiento y recuperación de embeddings se seleccionó ChromaDB. Este motor permite una persistencia en disco simple y eficiente, facilitando la conservación de los vectores entre sesiones de trabajo. Además, cuenta con una API clara y fácil de utilizar, lo que agiliza la integración y el desarrollo. Otra ventaja relevante es su soporte para búsquedas híbridas, que combina similitud semántica con filtrado tradicional y resulta especialmente útil para los apartados posteriores del trabajo. Finalmente, ChromaDB ofrece integración directa con LangChain, permitiendo incorporar rápidamente funcionalidades de recuperación aumentada con lenguaje natural dentro del pipeline del trabajo.

2.2. Fuente Tabular (Catálogo y Ventas)

Datos incluidos

Se cargaron cuatro tablas principales:

- **productos** (id, nombre, marca, categoría, subcategoría, precio, stock, color, capacidad, voltaje, potencia, garantía)
- **ventas** (id, producto_id, metodo_pago, fecha, monto, cantidad, sucursal, fecha, cliente)
- **Inventario** (id, sucursal, id_producto, nombre_producto, proveedor, stock)
- **devoluciones** (id, venta, fecha, id_producto, cliente, motivo, monto_reembolso, método_reembolso, observaciones)

El formato tabla permite:

- consultas exactas
- filtros numéricos y categóricos
- cálculos estadísticos
- agregaciones

2.3. Fuente de Grafos

Modelo utilizado: KùzuDB

Se seleccionó KùzuDB como motor de base de datos debido a varias ventajas relevantes para el proyecto. En primer lugar, ofrece compatibilidad con consultas en lenguaje Cypher, lo que facilita la construcción de grafos y la ejecución de consultas complejas de manera declarativa. Además, se trata de un sistema embebido, liviano y de baja sobrecarga, ideal para entornos de experimentación o prototipado donde se requiere rapidez y simplicidad. Finalmente, su instalación en Google Colab resulta directa y sin dependencias externas pesadas, lo que permite integrarlo rápidamente al flujo de trabajo y ejecutar pruebas sin complicaciones.

Modelo del grafo

Se representaron relaciones:

- (Producto) —[:PERTENECE_A]→ (Categoría)
- (Producto) —[:FABRICADO_POR]→ (Marca)

Esto permite consultas como:

- “¿Qué productos están relacionados con la categoría Cocina?”
- “Productos hermanos de la procesadora KitchenPro.”

2.4. Clasificador de Intención

- **Clasificador entrenado propio:** Regresión Logística sobre Embeddings (Sentence-Transformers).

Se generaron consultas sintéticas clasificadas en SQL, VECTOR, GRAFOS, incluyendo variación semánticas, sinónimo y errores ortográficos intencionales. El dataset se dividió en 70% para entrenamiento y 30% para evaluación, obteniendo métricas de precisión, recall y F1 para medir el desempeño del clasificador.

--- MODELO A (Local Embeddings) ---				
	precision	recall	f1-score	support
GRAFO	1.00	0.50	0.67	4
SQL	0.57	0.80	0.67	5
VECTOR	0.80	0.80	0.80	5
accuracy			0.71	14
macro avg	0.79	0.70	0.71	14
weighted avg	0.78	0.71	0.71	14

- **Clasificador basado en LLM (Hugging Face)** : Además del modelo local basado en embeddings, se implementó un segundo enfoque de clasificación utilizando un LLM alojado en Hugging Face: **Qwen 2.5-72B Instruct**.

A diferencia del modelo local —que aprende un hiperplano de decisión sobre embeddings— este modelo realiza la clasificación mediante razonamiento semántico contextual, guiado por un conjunto reducido de ejemplos (“few-shot prompting”).

El modelo Qwen permite:

- Detectar intención incluso en frases ambiguas o con errores.
- Manejar consultas más complejas que mezclan varias intenciones.
- Captar relaciones semánticas profundas sin necesidad de ampliar el dataset sintético.

Al igual que con el modelo local, se evaluó el LLM utilizando el 30% del dataset reservado para test. Se generó un conjunto de predicciones (`y_pred_hf`) y se comparó con las etiquetas reales mediante precisión, recall y F1-score.

--- MODELO B (Hugging Face Qwen) ---				
	precision	recall	f1-score	support
GRAFO	1.00	1.00	1.00	4
SQL	1.00	1.00	1.00	5
VECTOR	1.00	1.00	1.00	5
accuracy			1.00	14
macro avg	1.00	1.00	1.00	14
weighted avg	1.00	1.00	1.00	14

Decisión: Se seleccionó el LLM Generativo.

```
=====
PRECISIÓN LOCAL: 0.71
PRECISIÓN NUBE: 1.00
-----
```

En las pruebas de evaluación, el modelo local basado en embeddings alcanzó una precisión del **71%**, mostrando dificultades para separar correctamente clases semánticamente similares. En particular, presentó confusión entre consultas de tipo **GRAFO** y **SQL**, reduciendo su recall al 50% en la primera categoría.

Por el contrario, el modelo basado en LLM (Qwen 2.5-72B) obtuvo un **100% de precisión y recall** en las tres clases, clasificando correctamente incluso frases ambiguas o con estructuras lingüísticas menos directas.

Dado que una clasificación incorrecta implica ejecutar un módulo equivocado (por ejemplo, enviar una consulta SQL a la base vectorial o viceversa), la fiabilidad del LLM resulta crítica. A pesar de su mayor latencia, su capacidad de razonamiento contextual y su precisión perfecta en las pruebas justifican su selección como modelo principal de enrutamiento de consultas.

2.5. Modelo de Lenguaje (LLM)

Modelo Seleccionado: **gemini-2.0-flash (Google)** con Temperature=0.

Capacidad de Razonamiento: El modelo Flash destaca por su baja latencia y alta capacidad lógica, crucial para la generación de código SQL y Cypher sin alucinaciones.

Ventana de Contexto: Su amplia ventana de contexto permite injectar esquemas de bases de datos complejos y múltiples documentos recuperados sin perder información.

Determinismo: Al configurarse con temperatura 0, el modelo se vuelve determinista, ideal para tareas de extracción de datos y uso de herramientas donde la creatividad no es deseada, sino la precisión fáctica.

3. RESULTADOS DE EJECUCIÓN

A continuación, se presentan las pruebas realizadas para validar la robustez del sistema en sus dos etapas.

3.1. Ejercicio 1: Pipeline RAG (Retrieval)

En esta fase, el sistema clasifica la intención y recupera información de tres fuentes distintas.

1. Prueba de la base de datos vectorial:

- *Usuario*: "Mi licuadora huele a quemado, ¿qué hago?"
- *Resultado*: El sistema clasificó la intención como VECTOR. Realizó una búsqueda híbrida en la base de datos vectorial (Manuales y Reseñas),

recuperando el manual de la licuadora. El LLM sintetizó las posibles causas (sobrecarga, ingredientes duros) y listó soluciones específicas.

```
# Prueba 1: Vectorial (Problema técnico)
print(f"\n RESPUESTA BOT:\n{bot.respond('Mi licuadora huele a quemado, ¿qué hago?')}"")
```

USUARIO: Mi licuadora huele a quemado, ¿qué hago?
Router Inteligente decidió: [VECTOR]

PIPELINE ACTIVO | Intención: VECTOR

RERANKING: 'Mi licuadora huele a quemado, ¿qué hago?'
Rank 1 | Score: 5.4730 | Fuente: manual_P0184_Pro_Ventilador_de_Pie.md
Rank 2 | Score: 5.1970 | Fuente: resena_R01325.txt
Rank 3 | Score: 5.0739 | Fuente: resena_R04214.txt
Generando respuesta final...

RESPUESTA BOT:
El olor a quemado en tu licuadora puede deberse a varias causas:

- * Sobrecarga de ingredientes
- * Ingredientes muy duros o pegajosos
- * Motor sobrecalentado
- * Acumulación de residuos en componentes móviles
- * Desgaste de componentes internos

Te sugiero probar las siguientes soluciones:

1. Reducir la cantidad de alimentos que introduces en la licuadora.
2. Procesar los ingredientes en tandas más pequeñas.
3. Dejar enfriar el motor durante 15-20 minutos.
4. Limpiar todas las piezas móviles de la licuadora.
5. Si el problema persiste, te recomiendo solicitar una revisión técnica.

2. Prueba de la base de datos tabular:

- *Usuario:* "¿Tenés stock de la Procesadora de KitchenPro y cuánto cuesta?"
- *Resultado:* La consulta fue clasificada como SQL. Se generó dinámicamente la query `SELECT p.stock, p.precio_usd FROM productos...` filtrando por nombre y marca. La respuesta fue precisa y basada en datos estructurados en tiempo real.

```
# Prueba 2: SQL (Consulta de precios/stock)
print(f"\n RESPUESTA BOT:\n{bot.respond('¿Tenés stock de la Procesadora de KitchenPro y cuánto cuesta?')}"")
```

```

USUARIO: ¿Tenés stock de la Procesadora de KitchenPro y cuánto cuesta?
Router Inteligente decidió: [SQL]

PIPELINE ACTIVO | Intención: SQL
[Gemini SQL] Analizando: '¿Tenés stock de la Procesadora de KitchenPro y cuánto cuesta?...'
SQL Generado: SELECT
    p.stock,
    p.precio_usd
FROM productos AS p
WHERE
    p.nombre LIKE '%Procesadora%' AND p.marca LIKE '%KitchenPro%';
Generando respuesta final...

RESPUESTA BOT:
Sí, tenemos stock de la Procesadora de KitchenPro. Tenemos 92 unidades disponibles y su precio es de 82.93 USD.

```

3. Prueba de la base de datos de grafo:

- *Usuario*: "¿Qué otros productos hay en la misma categoría que la Procesadora?"
- *Resultado*: La intención fue detectada como GRAFO. El sistema generó una consulta en lenguaje Cypher (MATCH (p1)-[:PERTENECE_A]->(c)<-[:PERTENECE_A]-(p2)...) para navegar los nodos de la base de datos Kùzu. Se listaron correctamente los productos de la misma categoría junto con sus precios.

```

# Prueba 3: Grafo (Relaciones)
print(f"\n RESPUESTA BOT:\n{bot.responder('¿Qué otros productos hay en la misma categoría que la Procesadora?')}")

```

```

USUARIO: ¿Qué otros productos hay en la misma categoría que la Procesadora?
Router Inteligente decidió: [GRAFO]

PIPELINE ACTIVO | Intención: GRAFO
Cypher: MATCH (p1:Producto)-[:PERTENECE_A]->(c:Categoría)<-[ :PERTENECE_A]-(p2:Producto)
WHERE p1.nombre = 'Procesadora'
RETURN p2.nombre, p2.precio
Generando respuesta final...

RESPUESTA BOT:
En la misma categoría que la Procesadora, tenemos los siguientes productos:

* Deluxe Procesadora: 1169.67 USD
* Profesional Procesadora 2024: 356.86 USD
* Advanced Procesadora II: 2528.96 USD
* Procesadora: 82.93 USD

```

4. Prueba de Fallback (Manejo de errores):

- *Usuario*: "¿Cómo hacer una torta de chocolate?"

- **Resultado:** El sistema buscó en la base de conocimientos, y al no encontrar información relevante (similitud baja), respondió: "Nouento con esa información en mis registros actuales", evitando inventar una receta.

```
# Prueba 4: Manejo de errores
print(f"\n RESPUESTA BOT:\n{bot.respond('¿Cómo hacer una torta de chocolate?')}")

USUARIO: ¿Cómo hacer una torta de chocolate?
Router Inteligente decidió: [VECTOR]

PIPELINE ACTIVO | Intención: VECTOR

RERANKING: '¿Cómo hacer una torta de chocolate?'
Rank 1 | Score: 5.7402 | Fuente: manual_P0013_Procesadora.md
Rank 2 | Score: 5.7390 | Fuente: resena_R01515.txt
Rank 3 | Score: 5.7213 | Fuente: resena_R02602.txt
Generando respuesta final...

RESPUESTA BOT:
Nouento con esa información en mis registros actuales. Puedo ayudarte con los pasos para amasar con la procesadora. ¿Te interesa?
```

5. Prueba Híbrida (Ambigüedad):

- **Usuario:** "Opiniones sobre la cafetera"
- **Resultado:** Clasificado correctamente como VECTOR (Reseñas), recuperando fragmentos de experiencias de usuarios reales almacenados en ChromaDB.

```
# Prueba Híbrida (Ambigüedad):
print(f"\n RESPUESTA BOT:\n{bot.respond('Opiniones sobre la cafetera')}")

USUARIO: Opiniones sobre la cafetera
Router Inteligente decidió: [VECTOR]

PIPELINE ACTIVO | Intención: VECTOR

RERANKING: 'Opiniones sobre la cafetera'
Rank 1 | Score: 5.7498 | Fuente: manual_P0004_Compacto_Licuadora.md
Rank 2 | Score: 4.6178 | Fuente: resena_R00776.txt
Rank 3 | Score: 4.3205 | Fuente: resena_R04450.txt
Generando respuesta final...

RESPUESTA BOT:
Aquí tienes algunas opiniones sobre la Cafetera (P0122) de otros usuarios:
* **Francisco\Martínez (Río Negro):** Le da un puntaje de 4/5 y destaca la increíble relación calidad-precio, durabilidad y elegancia del
* **Renata\Reyes (Santa Fe):** Le da un puntaje de 3/5 y comenta que tiene pros y contras. Considera que está bien por el precio, pero no
```

3.2. Ejercicio 2: Agente Autónomo (ReAct + Herramientas)

En esta fase, el sistema posee memoria conversacional y la capacidad de ejecutar código y generar gráficos analíticos.

1. Prueba de Razonamiento y SQL:

- **Usuario:** "¿Cuál es el producto más caro en stock?"
- **Agente:** Identificó la necesidad de la herramienta tool_buscar_datos_sql. Generó la query de ordenamiento y respondió: "El producto más caro es la Advanced Heladera con un precio de 2992.33 USD".

```

# Prueba 1: SQL + Memoria
print("\n--- PRUEBA 1: Datos y Memoria ---")
print(agent_executor.invoke({"input": "¿Cuál es el producto más caro en stock?"})['output'])

--- PRUEBA 1: Datos y Memoria ---
...
Thought: Do I need to use a tool? Yes
Action: tool_buscar_datos_sql
Action Input: SELECT nombre_producto, precio FROM productos ORDER BY precio DESC LIMIT 1; [Gemini SQL] Analizando: 'SELECT nombre_producto, p
SQL Generado: SELECT nombre, precio_usd FROM productos ORDER BY precio_usd DESC LIMIT 1;
| nombre | precio_usd |
|:-----:|:-----:|
| Advanced Heladera | 2992.33 |WARNING:langchain_core.callbacks.manager:Error in StdOutCallbackHandler.on_chain_start callback: Attribut
Do I need to use a tool? No
Final Answer: El producto más caro en stock es la Advanced Heladera, con un precio de 2992.33 USD.

> Finished chain.
El producto más caro en stock es la Advanced Heladera, con un precio de 2992.33 USD.

```

2. Prueba de Memoria Conversacional:

- *Usuario*: "¿Y de qué marca es?" (Sin mencionar el producto).
- *Agente*: Utilizó el historial de chat (ConversationBufferMemory), entendió que se refería a la "Advanced Heladera" mencionada anteriormente, y respondió: "Es de la marca TechHome".

```

print(agent_executor.invoke({"input": "¿Y de qué marca es?"})['output'])

```
Thought: Do I need to use a tool? Yes
Action: tool_buscar_relaciones
Action Input: Marca de Advanced HeladeraTechHome fabrica Electrodomésticos y es dueña de La marca.Do I need to use a tool? No
Final Answer: La Advanced Heladera es de la marca TechHome.

> Finished chain.
La Advanced Heladera es de la marca TechHome.

```

## 3. Prueba de Visualización de Datos (Analytics):

- *Usuario*: "Genera un gráfico de promedio de precio por marca".
- *Agente*: Detectó una intención analítica y activó la tool\_generar\_grafico. El agente escribió y ejecutó código Python con matplotlib, realizando una agregación SQL (AVG(precio\_usd) GROUP BY marca) para visualizar la comparativa de precios de forma gráfica.

```

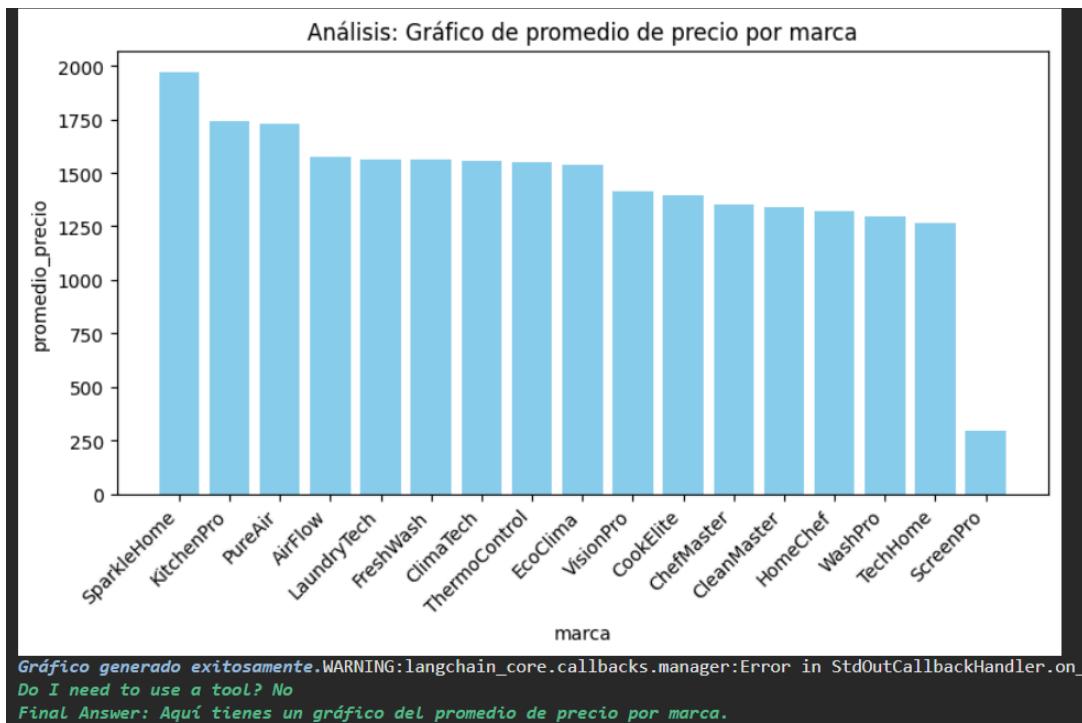
Prueba 2: Gráficos
print("\n--- PRUEBA 2: Generación de Gráficos ---")
agent_executor.invoke({"input": "Genera un gráfico de promedio de precio por marca"})

```

```

--- PRUEBA 2: Generación de Gráficos ---
```
Thought: Do I need to use a tool? Yes
Action: tool_generar_grafico
Action Input: Gráfico de promedio de precio por marca
[HERRAMIENTA ANALYTICS ACTIVADA: 'Gráfico de promedio de precio por marca']
SQL Gráfico: SELECT
    marca,
    AVG(precio_usd) AS promedio_precio
FROM productos
GROUP BY
    marca
ORDER BY
    promedio_precio DESC;
```

```



```

> Finished chain.
{'input': "Genera un gráfico de promedio de precio por marca",
'chat_history': [HumanMessage(content='¿Cuál es el producto más caro en stock?'),
AIMessage(content='El producto más caro en stock es la Advanced Heladera, con un precio de 2992.33 USD.'),
HumanMessage(content='Y de qué marca es?'),
AIMessage(content='La Advanced Heladera es de la marca TechHome.'),
HumanMessage(content='Genera un gráfico de promedio de precio por marca'),
AIMessage(content='Aquí tienes el gráfico del promedio de precio por marca.']),
'output': 'Aquí tienes el gráfico del promedio de precio por marca.'}

```

#### 4. Prueba Multitarea (Orquestación):

- **Usuario:** "Mi licuadora hace ruido extraño. ¿Qué puede ser y cuánto cuesta una nueva?"
- **Agente:** Demostró capacidad de planificación compleja descomponiendo el problema en pasos: Usó `tool_buscar_manuales` para diagnosticar el problema técnico ( posible causa: junta deteriorada). Usó `tool_buscar_datos_sql` para buscar el precio de una licuadora nueva en

la base de datos. Finalmente, consolidó ambas informaciones en una única respuesta coherente y útil para el usuario.

```
Prueba 3: Manuales + Precios (Multitarea)
print("\n--- PRUEBA 3: Multitarea (Manuales + Precios) ---")
agent_executor.invoke({"input": "Mi licuadora hace ruido extraño. ¿Qué puede ser y cuánto cuesta una nueva?"})
```

```
--- PRUEBA 3: Multitarea (Manuales + Precios) ---
...
Thought: Do I need to use a tool? Yes
Action: tool_buscar_manuales
Action Input: licuadora ruido extraño
```
PIPELINE ACTIVO | Intención: VECTOR
RERANKING: 'licuadora ruido extraño
```
Rank 1 | Score: 5.7166 | Fuente: manual_P0288_Plancha_Seca_II.md
Rank 2 | Score: 5.6507 | Fuente: resena_R01449.txt
Rank 3 | Score: 5.5271 | Fuente: resena_R03726.txt
Posibles causas:
- Junta de goma deteriorada o mal colocada
- Recipiente con rajaduras
- Sobrellenado del recipiente
- Tapa no cerrada correctamente
- Rosca de la cuchilla floja
Soluciones:
1. Inspeccionar y reemplazar juntas dañadas
2. Verificar que no haya grietas en el recipiente
3. No llenar más allá de La Línea MAX
4. Asegurar cierre hermético de la tapa
5. Apretar La base de cuchillas (con el aparato apagado)
```
Fecha: 2025-02-22
Usuario: Joaquín_Benítez
Teléfono: +54 9 332 1766-7882
Producto: Neblinizador (P0206)
Puntaje: 4/5
Provincia: Córdoba
Les cuenta... Increíble relación calidad-precio con Neblinizador. Lo compré hace un mes y es muy confiable, además de fácil de usar. La relación calidad-precio es excelente. Gracias
```

```
Fecha: 2025-10-15
Usuario: Ignacio_Rojas
Teléfono: +54 9 318 3247-1890
Producto: Smart Lavarropas Semiautomático Plus (P0234)
Puntaje: 4/5
Provincia: Misiones
Les cuenta... Muy buena calidad, estoy feliz con Smart Lavarropas Semiautomático Plus. Lo compré hace una semana y es muy robusto, además de intuitivo. El servicio de entrega fue muy rápido.
Action: tool_buscar_datos_sql
Action Input: precio_Licuadora
```
[Gemini SQL] Analizando: 'precio licuadora
```
SQL Generado: SELECT precio_usd FROM productos WHERE nombre LIKE '%licuadora%';
| precio_usd |
|-----|
| 283.63 |
| 1273.06 |
| 329.07 |
| 259.42 |
| 2662.78 |
| 1837.15 |
| 2532.47 |
| 1261.99 | Do I need to use a tool? No
Final Answer: El ruido extraño en tu Licuadora puede deberse a varias razones, como una junta deteriorada, rajaduras en el recipiente, sobrellenado, tapa mal cerrada o la rosca de la cuchilla floja.
```

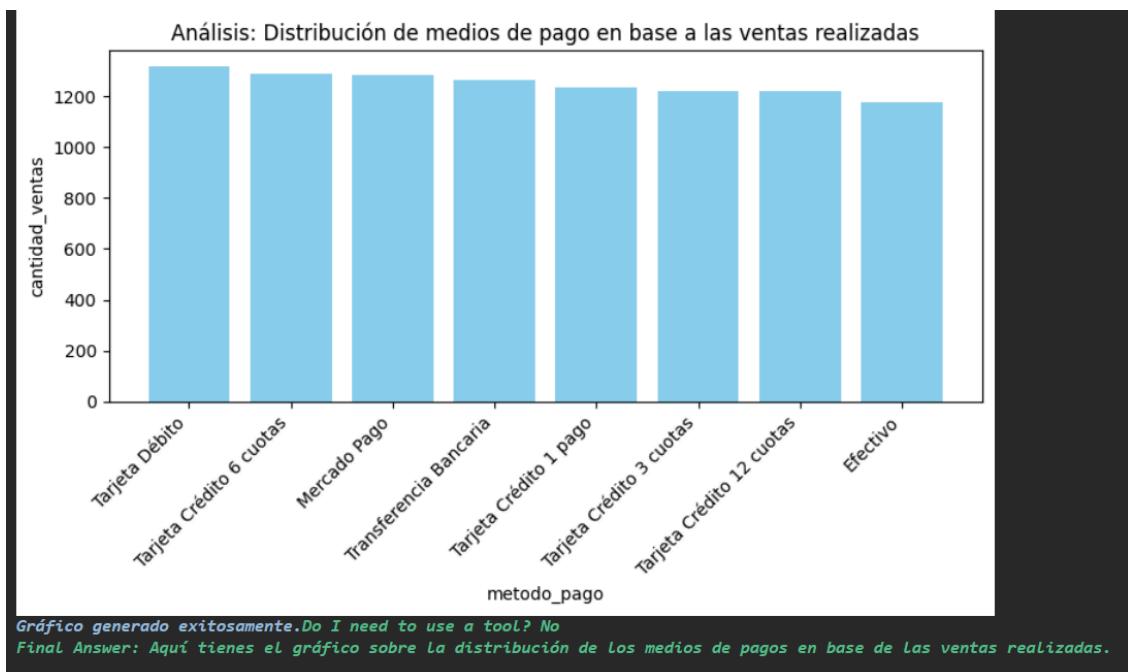
```
> Finished chain.
{"input": "Mi licuadora hace ruido extraño. ¿Qué puede ser y cuánto cuesta una nueva?", "chat_history": [HumanMessage(content="¿Cuál es el producto más caro en stock?"), AIMessage(content="El producto más caro en stock es la Advanced Heladera, con un precio de 2992.33 USD."), HumanMessage(content="¿Y de qué marca es?"), AIMessage(content="La Advanced Heladera es de la marca TechHome."), HumanMessage(content="Genera un gráfico de promedio de precio por marca"), AIMessage(content="Aquí tiene el gráfico del promedio de precio por marca."), HumanMessage(content="Mi licuadora hace ruido extraño. ¿Qué puede ser y cuánto cuesta una nueva?"), AIMessage(content="El ruido extraño en tu licuadora puede deberse a varias razones, como una junta deteriorada, rajaduras en el recipiente, sobrellenado, tapa mal cerrada o la rosca de la cuchilla floja. Te recomiendo revisar esos puntos. En cuanto al precio de una licuadora nueva, tenemos varios modelos con precios que van desde 259.42 USD hasta 2662.78 USD."), "output": "El ruido extraño en tu licuadora puede deberse a varias razones, como una junta deteriorada, rajaduras en el recipiente, sobrellenado, tapa mal cerrada o la rosca de la cuchilla floja. Te recomiendo revisar esos puntos. En cuanto al precio de una licuadora nueva, tenemos varios modelos con precios que van desde 259.42 USD hasta 2662.78 USD."}]}
```

5. Prueba de Análisis de Negocio:

- **Usuario:** "Dame un gráfico sobre la distribución de los medios de pagos en base de las ventas realizadas".
- **Agente:** Comprendió la solicitud analítica compleja, consultó la tabla ventas, agrupó por `metodo_pago` y generó la visualización correcta.

```
# Prueba 5
agent_executor.invoke([
    "input": "Dame un gráfico sobre la distribución de los medios de pagos en base de las ventas realizadas"
})
```

```
...
Thought: Do I need to use a tool? Yes
Action: tool_generar_grafico
Action Input: Distribución de medios de pago en base a Las ventas realizadas
[HERRAMIENTA ANALYTICS ACTIVADA: 'Distribución de medios de pago en base a las ventas realizadas'
↳ SQL Gráfico: SELECT
    metodo_pago,
    COUNT(*) AS cantidad_ventas,
    SUM(total) AS total_ventas
FROM ventas
GROUP BY
    metodo_pago
ORDER BY
    cantidad_ventas DESC;
```



```
> Finished chain.
{'input': 'Dame un gráfico sobre la distribución de los medios de pagos en base de las ventas realizadas',
'chat_history': [HumanMessage(content='¿Cuál es el producto más caro en stock?'),
AIMessage(content='El producto más caro en stock es la Advanced Heladera, con un precio de 2992.33 USD.'),
HumanMessage(content='¿Y de qué marca es?'),
AIMessage(content='La Advanced Heladera es de la marca TechHome.'),
HumanMessage(content='Genera un gráfico de promedio de precio por marca'),
AIMessage(content='Aquí tienes el gráfico del promedio de precio por marca.'),
HumanMessage(content='Mi licuadora hace ruido extraño. ¿Qué puede ser y cuánto cuesta una nueva?'),
AIMessage(content='El ruido extraño en tu licuadora puede deberse a varias razones, como una junta deteriorada, rajaduras en el recipiente, sobrellenado, tapa mal cerrada o la rosca de la cuchilla floja. Te recomiendo revisar estos puntos. En cuanto al precio de una licuadora nueva, tenemos varios modelos con precios que van desde 259.42 USD hasta 2602.78 USD.'),
HumanMessage(content='Dame un gráfico sobre la distribución de los medios de pagos en base de las ventas realizadas'),
AIMessage(content='Aquí tienes el gráfico sobre la distribución de los medios de pagos en base de las ventas realizadas.']),
'output': 'Aquí tienes el gráfico sobre la distribución de los medios de pagos en base de las ventas realizadas.'}
```

4. MEJORAS Y TRABAJOS FUTUROS

Durante el desarrollo y las pruebas, se detectaron áreas de oportunidad para perfeccionar el sistema:

- Manejo de Logs y Depuración:** La librería LangChain emitió advertencias visuales (AttributeError in StdOutCallbackHandler) en el entorno de Colab, aunque no afectaron la ejecución. Como solución se podría configurar un manejador de logs

personalizado para filtrar advertencias de librerías internas y presentar una traza limpia de ejecución, mejorando la monitorización en producción.

2. **Optimización de Latencia del Router:** Actualmente, el Router usa un LLM grande (Gemini Flash). Para reducir costos y latencia a gran escala, se podría utilizar los pares de preguntas/etiquetas generados en este TP para entrenar (Fine-Tuning) un modelo pequeño tipo BERT o DistilBERT especializado únicamente en clasificación, reservando el LLM solo para la generación de respuestas complejas.
3. **Interfaz de Usuario (UI):** Una posible mejora sería migrar de la consola de ejecución de notebook a una interfaz web interactiva, permitiendo una interacción más fluida y la visualización nativa de los gráficos generados.

5. BIBLIOGRAFÍA Y REFERENCIAS

1. **LangChain Documentation:** *Introduction to RAG and Agents*. Recuperado de <https://python.langchain.com/>
2. **Google AI Studio: Gemini API Documentation & Prompt Engineering Guidelines**.
3. **MTEB Leaderboard:** *Massive Text Embedding Benchmark* (Hugging Face) para la selección del modelo intfloat/multilingual-e5-small.
4. **KùzuDB:** *Cypher Query Language for Embedded Graph Databases*. Documentación oficial.
5. **ChromaDB:** *Vector Store for AI*. Guía de integración y persistencia.