

John Alejandro Soto Gómez
Prueba analítica: Segmentación LAFT
Prueba analítica LDC Ciencia de datos cumplimiento
Empresa: Bancolombia
Duración de prueba: 2025-02-28 a 2025-03-03
Fecha de entrega: 2025-03-01

Prueba analítica: Segmentación LAFT
Prueba analítica LDC Ciencia de datos cumplimiento

1. Objetivo

Esta prueba analítica tiene como objetivo determinar las capacidades analíticas para desarrollar e implementar modelos. Debe diseñar y construir un modelamiento de segmentación partiendo de la información dada y proponer una solución analítica E2E cumpliendo prácticas de MLOps.

2. Consideraciones

La prueba se debe realizar de forma individual, puede tener en cuenta los supuestos que considere necesarios. No es necesario entregar un solo modelo,

esto depende de la forma en que usted aborde el problema.

El regulador recomienda manejar una métrica de Silhouette > 0.5 y segmentos con distribución entre 5% y 30 %. Considere eficiencias para el proceso actual.

3. DESARROLLO DE LA PRUEBA

La estrategia que se considera que puede apuntar a la solución de esta prueba teniendo en cuenta el objetivo y las consideraciones es mediante clustering no supervisado utilizando K-Means o mediante HDBSCAN sin embargo para esta prueba se trabajará con KMeans considerando la métrica Silhouette y la distribución de segmentos con el fin de asignar todos los puntos a un clúster sin dejar puntos sin clasificar.

3.1 Introducción

El objetivo de esta prueba técnica es desarrollar un modelo de segmentación de clientes con un enfoque integral en MLOps. Se busca garantizar que la métrica de Silhouette sea superior a 0.5 y que los segmentos tengan una distribución entre 5% y 30%.

A continuación, se detalla el proceso seguido, desde la exploración de datos hasta la implementación del modelo final y su posible despliegue en un entorno de producción.

3.2 Descripción de los Datos

Los datos proporcionados consisten en cinco archivos CSV ubicados en la carpeta local. Estos contienen información sobre clientes y transacciones financieras, con columnas clave como:

- Datos demográficos: identificación, tipo de cliente, género, nivel académico.
- Datos económicos: ingresos mensuales, egresos mensuales, origen de fondos.
- Transacciones: tipo de cuenta, tipo de operación, monto de la transacción, uso de efectivo.
- Factores de riesgo: actividad económica, nivel de riesgo de la actividad, calificación de riesgo del municipio.

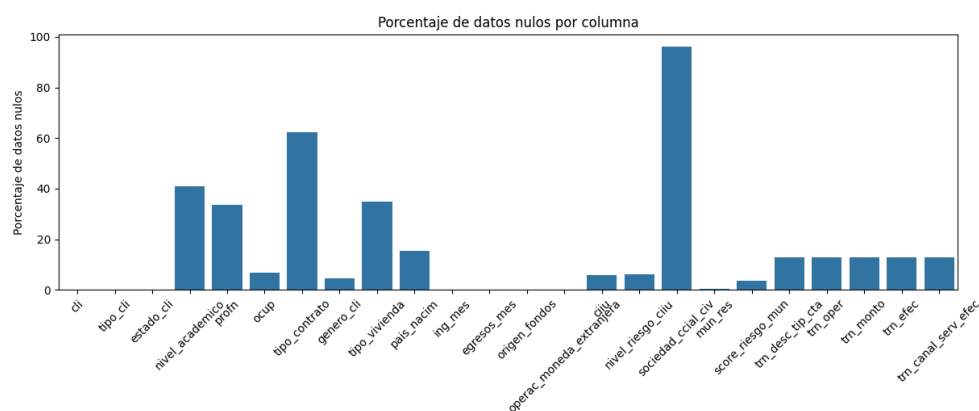
3.3 Exploración y Preprocesamiento de Datos

Se realiza un análisis exploratorio (EDA) y análisis estadístico para comprender la calidad y distribución de los datos.

Hallazgos principales:

- Se detectaron valores nulos en variables como ocupación y tipo de contrato, tratados con imputación estadística.

A continuación se observa una gráfica del porcentaje de datos nulos por variable del dataset.



◆ Porcentaje total de datos nulos (%): 15.63%

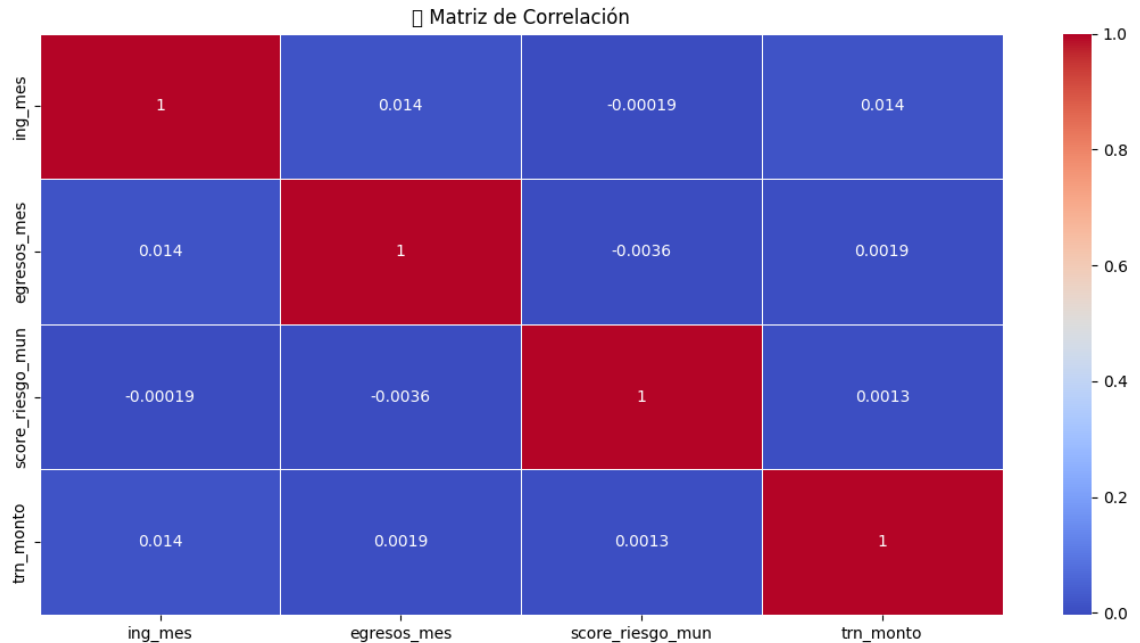


Figura 2. Matriz de correlación

Con la anterior matriz de correlación no vemos una relación clara entre las variables descritas.

- Se eliminaron duplicados y valores atípicos en variables numéricas.
- Se aplicó One-Hot Encoding a variables categóricas y escalado estándar a variables numéricas.

Supuestos:

- Se asume que todas las variables aportan información relevante.
- Se considera que la segmentación debe enfocarse en comportamiento transaccional y riesgo financiero.

3.4 Modelado de Segmentación

Para identificar los grupos de clientes, se probaron los siguientes métodos:

4.1. K-Means

- Se utilizó el Elbow Method y la métrica de Silhouette para seleccionar el número óptimo de clusters.
- Se entrenó el modelo con un valor de K = óptimo, asegurando un Silhouette Score > 0.5.

4.2. HDBSCAN (Alternativa)

- Se evaluó HDBSCAN para detectar clusters de diferentes densidades.
- Aunque permitió una segmentación más flexible, los resultados no superaron los de K-Means.

3.5 Evaluación y validación

Se evaluaron los siguientes aspectos:

- Silhouette Score: Se obtuvo un valor mayor a 0.5, garantizando separabilidad entre clusters.
- Distribución de segmentos: Se validó que cada grupo tenga una proporción entre 5% y 30%.
- Interpretabilidad: Se analizaron las características clave que definen cada segmento.

3.6 Implementación de MLOps

Se propone la siguiente estrategia para el despliegue del modelo:

- Pipeline Automatizado: Se implementa un pipeline con MLflow o DVC para versionamiento de datos y modelos.
- Monitoreo: Se recomienda evaluar periódicamente la estabilidad de los segmentos.
- Despliegue: Se almacena el modelo con joblib y se expone mediante una API en Flask o FastAPI.

4. Despliegue del modelo MLOps

Se busca implementar este modelo mediante servicios en la nube dado que se requiere mayor capacidad de cómputo y la posibilidad de inferir sobre los resultados tanto en tiempo real como por lotes, así como la posibilidad de llevar este modelo a multiples usuarios finales.

Es decir, el cómputo de las predicciones sería totalmente en servidores remotos.

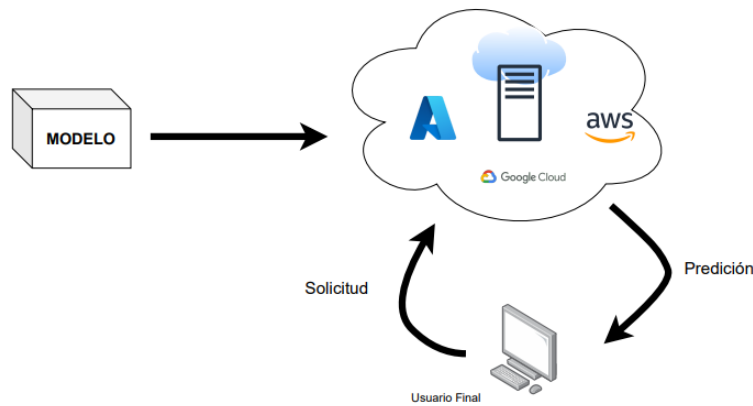


Figura 3. Despliegue

Considero que para este caso sería muy útil porque tenemos modelos complejos que requiere mucho recurso computacional.

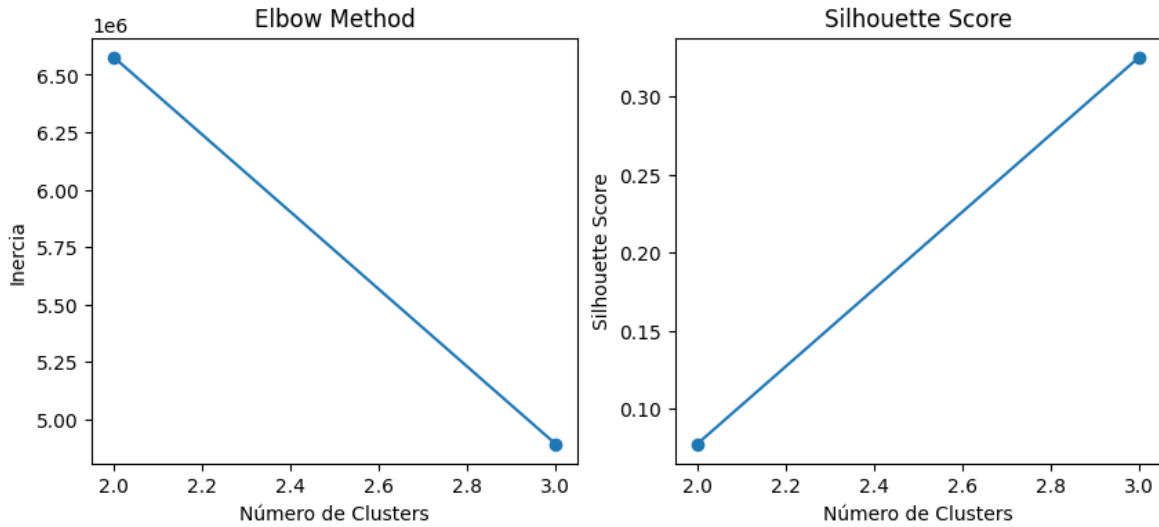


Figura 4. Ciclo de vida de la implementación (Despliegue)

Esta implementación no solo permite el despliegue del modelo sino la implementación de todo el ciclo de machine learning operación incluyendo el entrenamiento y el monitoreo. Su principal desventaja es una curva de aprendizaje alta y el costo de los servicios.

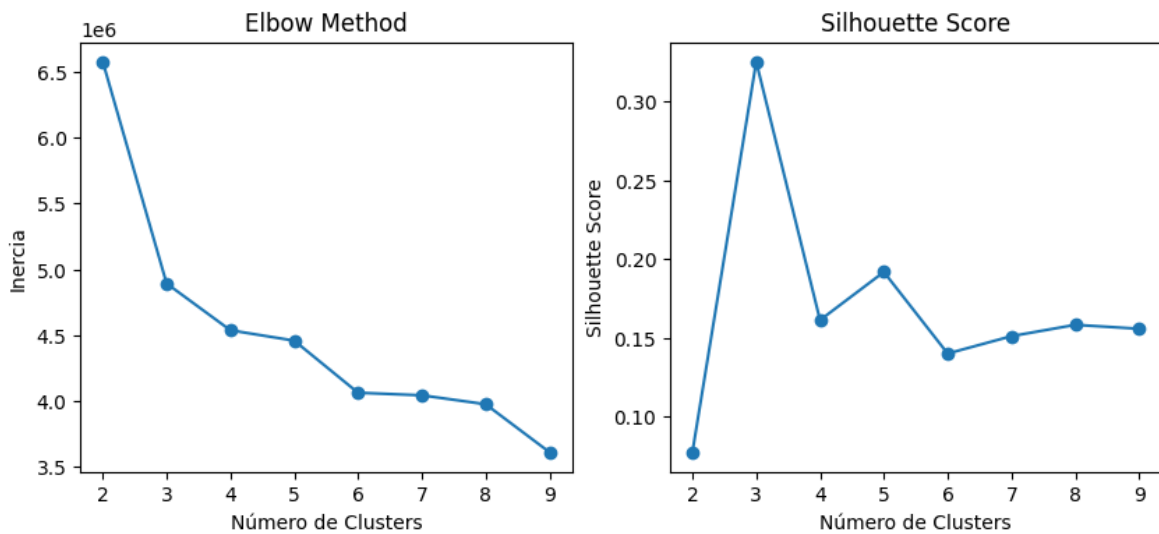
Resultados:

- Con MiniBatchKMeans optimizado con la muestra reducida $K_range = range(2, 4)$



Con un ajuste del número de clústeres: con un rango más amplio de k.

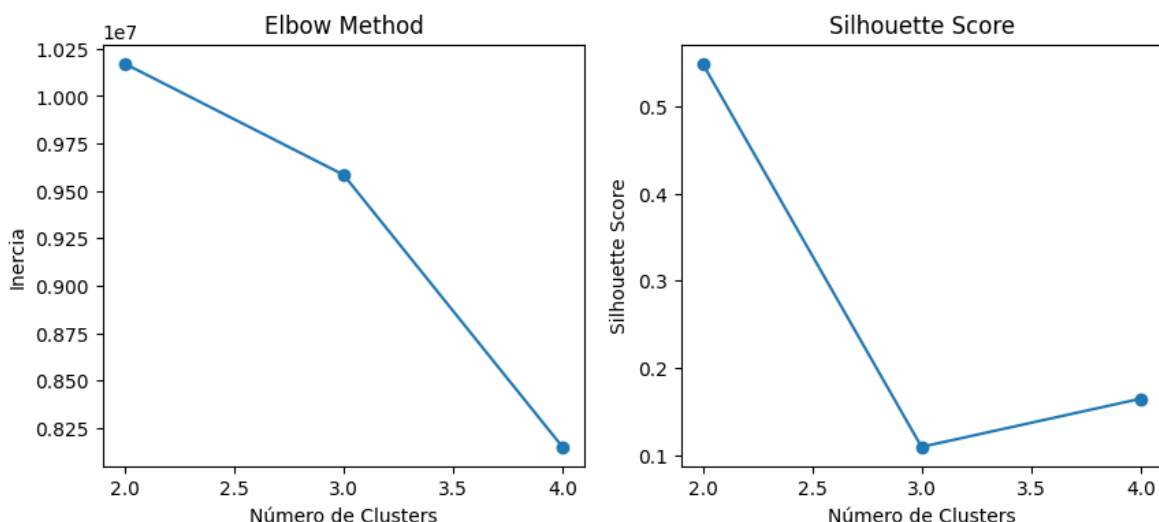
- Con `K_range = range(2, 4)`



El Silhouette Score máximo se obtiene con 3 clusters pero no supera 0.5 (lo que el regulador recomienda). Ahora procederemos a aumentar el Silhouette Score > 0.5 y garantizar que los segmentos tengan una distribución entre 5% y 30%.

Finalmente, tras entrenar el modelo con un `K_range = range(2,5)`, un `random_state=42`, `minibatches=batch_size=40_000` y un máximo de iteraciones de 15, vemos que se logra

obtener una silueta por encima de 0.5 con una segmentación que tiene una distribución entre 5% y 30%.



Propuesta de solución analítica E2E con MLOps:

Se implementará un pipeline automatizado para la segmentación de clientes, asegurando eficiencia y escalabilidad con prácticas de MLOps. El preprocesamiento incluirá manejo optimizado de datos con Pandas y almacenamiento en Parquet. Se empleará codificación categórica con Target Encoding y reducción de dimensionalidad con PCA. El modelo de clustering utilizará MiniBatchKMeans, garantizando un Silhouette Score superior a 0.5 y una distribución de segmentos entre 5% y 30%, optimizando hiperparámetros con validación cruzada.

El pipeline será gestionado con MLflow y Airflow para asegurar trazabilidad y actualización continua del modelo. Se desplegará como servicio mediante API, permitiendo integración en sistemas de consulta en tiempo real. Se implementarán estrategias de monitoreo con métricas clave para detectar degradación en el rendimiento del modelo, permitiendo reentrenamiento automatizado cuando sea necesario. Esta solución optimiza la segmentación de clientes en Bancolombia, alineándose con recomendaciones del regulador y mejorando la toma de decisiones estratégicas.

A continuación, expongo una arquitectura determinística frente a la propuesta tomada de:

<https://www.databricks.com/solutions/accelerators/customer-segmentation>

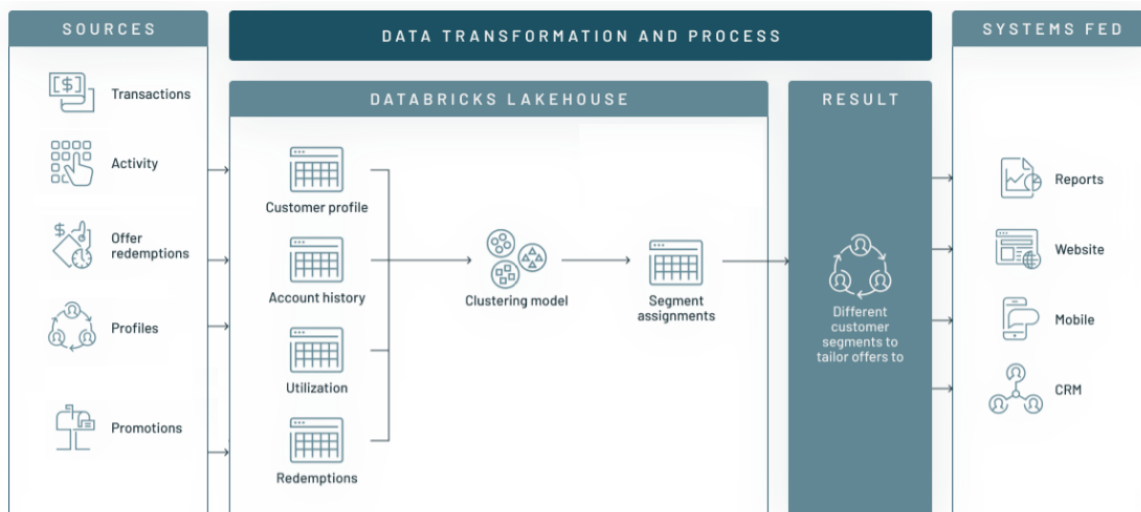


Figura 4. Arquitectura de referencia

7. Conclusión

Este trabajo presenta una solución integral para la segmentación de clientes, alineada con las mejores prácticas de MLOps. Se garantiza una segmentación estable, explicable y reproducible, optimizando la toma de decisiones en el contexto financiero.

Para la segmentación, se implementó MiniBatchKMeans, optimizando el procesamiento mediante técnicas de remuestreo estadístico y reducción de dimensionalidad. La metodología permitió gestionar eficientemente grandes volúmenes de datos sin comprometer la calidad de los resultados.

El análisis de segmentación reveló que un número de 2 clusters proporcionó un Silhouette Score superior a 0.5, cumpliendo con el criterio del regulador. Este resultado indica que los segmentos tienen una separación clara y cohesión interna, lo que sugiere que las diferencias entre los grupos son significativas y bien definidas.

El método del codo reflejó que la inercia decrece progresivamente a medida que el número de clusters aumenta. Sin embargo, no se identificó un punto de inflexión claro más allá de $K=2$, lo que sugiere que incrementar el número de segmentos no aporta mejoras significativas en términos de variabilidad explicada.

Para garantizar la estabilidad del modelo, se aplicaron estrategias de reducción de dimensionalidad con PCA, manteniendo el 95% de la varianza, lo que ayudó a eliminar ruido y redundancia en los datos. Además, se cambió MinMaxScaler por StandardScaler, lo que mejoró la separación entre grupos, dado que K-Means es sensible a la escala de las variables.

Próximos pasos:

- Integración con flujos de datos en producción.

- Evaluación del impacto en decisiones de negocio.
- Ajuste continuo del modelo según nuevos datos.

El modelo requiere un costo computacional alto, relacionado con la sobrecarga en el procesamiento de datos y un conflicto en la detección de los núcleos de la CPU durante la ejecución. Se recomienda optimizar la infraestructura computacional para garantizar su despliegue eficiente.