



UNIVERSIDAD DISTRITAL  
FRANCISCO JOSÉ DE CALDAS

---

# Extracción de Insights Clave para la Toma de Decisiones a partir de Comentarios Negativos de Detractores del Banco BBVA

---

Trabajo de Grado Modalidad Pasantía

**Autor:** Wilson Eduardo Jerez Hernández  
**Director:** Luis Alejandro Masmela Caíta  
**Profesional designado:** John Pablo Calvo López

Bogotá, DC  
Diciembre de 2024

---

## Resumen

Este trabajo presenta el desarrollo de un modelo de análisis de sentimientos basado en técnicas de Machine Learning, con el objetivo de clasificar y obtener insights a partir de comentarios negativos emitidos por los detractores del Banco BBVA. El propósito principal es identificar áreas de mejora en los servicios del banco para aumentar la satisfacción del cliente y reducir el número de detractores. Para lograrlo, se emplea procesamiento de lenguaje natural (NLP) sobre datos textuales obtenidos de redes sociales y encuestas internas.

## Palabras Clave

Machine Learning, Análisis de Sentimiento, NLP, BBVA, Clientes Detractores.

## Agradecimientos

Quiero agradecer profundamente a mi familia por su amor, paciencia y apoyo incondicional en cada paso de este camino.

A mis docentes, especialmente a mi director Alejandro Masmela, por su orientación y confianza durante todo este proceso.

Al semillero IPREA, gracias por darme la oportunidad de aprender y crecer en lo que más me apasiona.

A mi Sangha, por ser siempre una fuente de tranquilidad, sabiduría y luz en mi vida.

A John Calvo, gracias por tu guía, paciencia y enseñanzas. Tu apoyo fue clave para lograr el equilibrio entre mi crecimiento académico y personal.

Este trabajo es fruto del esfuerzo de todos, y les estaré siempre agradecido.

## Índice

<b>1. Introducción</b>	<b>4</b>
<b>2. Objetivos</b>	<b>5</b>
2.1. Objetivo General . . . . .	5
2.2. Objetivos Específicos . . . . .	5
<b>3. Marco Teórico</b>	<b>6</b>
3.1. BBVA Colombia: Una Entidad Financiera Líder . . . . .	6
3.1.1. Historia y Fundación . . . . .	6
3.1.2. Evolución y Fusión . . . . .	6
3.1.3. Adquisición de Granahorrar y Consolidación . . . . .	6
3.1.4. Innovación y Transformación Digital . . . . .	6
3.1.5. Inversión en Tecnología y Desarrollo Sostenible . . . . .	7
3.2. Análisis de Texto y Modelado de Tópicos para la Identificación de Áreas de Insatisfacción en Clientes Detractores . . . . .	7
3.2.1. Técnicas de Análisis de Texto . . . . .	7
3.2.2. Modelado de Tópicos con LDA . . . . .	10
3.2.3. Evaluación de Coherencia de Modelos . . . . .	10
3.2.4. Visualización de Resultados y Toma de Decisiones Estratégicas . . . . .	10
3.3. Descripción de los Datos . . . . .	10
<b>4. Creación y Desarrollo del Modelo</b>	<b>12</b>
4.1. Descripción del Modelo y Proceso de Implementación . . . . .	12
4.2. Pseudocódigo del Modelo . . . . .	14
<b>5. Resultados del Modelo de Tópicos</b>	<b>16</b>
5.1. Número Óptimo de Tópicos . . . . .	16
5.2. Distribución de los Temas Identificados . . . . .	16
5.3. Interpretación de los Tópicos . . . . .	17
5.4. Recomendaciones . . . . .	18

## ÍNDICE

---

<b>6. Anexos</b>	<b>19</b>
<b>7. Referencias</b>	<b>19</b>

### 1. Introducción

El sector bancario ha experimentado transformaciones profundas debido a la creciente digitalización, la evolución en las expectativas de los clientes y una competencia cada vez más feroz. En este entorno, la capacidad de captar y analizar la retroalimentación de los clientes se ha vuelto esencial para las entidades bancarias que desean adaptarse rápidamente a las demandas del mercado. Los comentarios negativos, en particular, representan una fuente valiosa de información, ya que pueden revelar problemas recurrentes o áreas de oportunidad que, si se abordan correctamente, podrían traducirse en una mayor satisfacción y fidelización de los clientes.

En este contexto, BBVA Colombia ha tomado la iniciativa de implementar un modelo de análisis de sentimientos utilizando técnicas avanzadas de Machine Learning. Este modelo se enfoca en la clasificación y análisis de los comentarios emitidos por los detractores, aquellos clientes que han expresado insatisfacción o frustración con los servicios ofrecidos. El objetivo principal de este enfoque es identificar de manera eficiente las áreas críticas de mejora, permitiendo a BBVA Colombia realizar ajustes estratégicos que mejoren la experiencia del cliente y reduzcan el índice de detractores.

Para lograrlo, se emplea el procesamiento de lenguaje natural (NLP), una rama de la inteligencia artificial que permite analizar y comprender grandes volúmenes de datos textuales de manera automatizada. Las fuentes de estos datos incluyen redes sociales, donde los clientes expresan sus opiniones de forma pública, y encuestas internas del banco, que proporcionan una visión más detallada y directa de las experiencias de los usuarios. A través de este análisis, BBVA Colombia busca convertir los comentarios negativos en oportunidades de mejora, fortaleciendo así su oferta de valor y consolidando su posición en un mercado altamente competitivo.

## 2. Objetivos

### 2.1. Objetivo General

Desarrollar un modelo de análisis de sentimiento basado en técnicas de Machine Learning que permita extraer insights clave a partir de los comentarios negativos emitidos por los detractores del Banco BBVA Colombia, con el fin de proporcionar información estratégica que facilite la toma de decisiones en la alta gerencia y contribuya a la mejora continua de los servicios ofrecidos.

### 2.2. Objetivos Específicos

1. **Transformación de comentarios en representaciones vectoriales:** Se implementará el modelo **TF-IDF Vectorizer** para convertir los comentarios negativos en representaciones vectoriales. Este proceso permitirá identificar las palabras y términos clave que están relacionados con las principales áreas de insatisfacción de los clientes. Esto servirá como una base sólida para el análisis posterior.
2. **Identificación de temas mediante Topic Modeling (LDA):** Aplicaremos técnicas de **Topic Modeling**, utilizando **Latent Dirichlet Allocation (LDA)**, para identificar los temas más recurrentes en los comentarios negativos. De esta forma, los comentarios se agruparán en tópicos relevantes, facilitando una visión clara de las áreas donde se concentran los principales problemas.
3. **Evaluación de la coherencia del modelo:** Para seleccionar el modelo óptimo, se calcularán los valores de **coherencia** para distintos números de temas. Este paso garantiza que los temas extraídos sean consistentes y representen adecuadamente las preocupaciones de los clientes.
4. **Análisis de áreas críticas de insatisfacción:** A partir de los temas identificados, se realizará un análisis profundo para determinar las áreas del producto o servicio que presentan mayores niveles de insatisfacción. Este análisis proporcionará **recomendaciones claras** para mejorar el servicio y reducir las quejas en los puntos críticos.
5. **Visualización de resultados:** Se generarán **visualizaciones** que muestren la distribución de los temas en los comentarios, ofreciendo una comprensión intuitiva de la prevalencia de cada tema y su relación con la insatisfacción de los clientes. Estas visualizaciones serán clave para **comunicar los resultados** a la alta gerencia de manera efectiva.
6. **Soporte para la toma de decisiones estratégicas:** Finalmente, los **insights** obtenidos del análisis de temas y la distribución de los comentarios proporcionarán **información clave** para orientar las acciones correctivas y preventivas. El objetivo es mejorar la experiencia del cliente y reducir el número de detractores.

### **3. Marco Teórico**

#### **3.1. BBVA Colombia: Una Entidad Financiera Líder**

BBVA Colombia es una destacada institución bancaria que forma parte del Grupo BBVA, uno de los conglomerados financieros más grandes del mundo. A lo largo de su historia, la entidad ha experimentado una notable evolución, marcada por adquisiciones estratégicas, innovaciones tecnológicas y un fuerte compromiso con el desarrollo sostenible[2].

##### **3.1.1. Historia y Fundación**

BBVA Colombia fue fundada en 1956 bajo el nombre de Banco Ganadero, una entidad financiera de economía mixta cuyo objetivo inicial era fomentar el desarrollo de la industria agropecuaria en Colombia. En 1996, el Banco Bilbao Vizcaya (BBV) adquirió el 34,7% de las acciones del Banco Ganadero, iniciando su expansión en el país y fortaleciendo su presencia en el sector financiero colombiano. [2]

##### **3.1.2. Evolución y Fusión**

En 1998, BBV incrementó su participación en el Banco Ganadero adquiriendo un 15% adicional de sus acciones, lo que le permitió tomar el control total de la entidad. A partir de ese momento, la institución pasó a llamarse BBV Banco Ganadero. En 1999, la fusión entre el Banco Bilbao Vizcaya y Argentaria dio lugar a BBVA, y en 2004, el nombre de la filial colombiana fue formalmente cambiado a BBVA Colombia, consolidando su identidad bajo la marca global. [2]

##### **3.1.3. Adquisición de Granahorrar y Consolidación**

Uno de los hitos más importantes en la historia de BBVA Colombia fue la adquisición de Granahorrar en 2005, una entidad financiera estatal enfocada en el mercado hipotecario. Un año después, las dos instituciones se fusionaron bajo la marca BBVA Colombia, lo que permitió consolidar su liderazgo en el sector hipotecario y ampliar significativamente su presencia en el país. [2]

##### **3.1.4. Innovación y Transformación Digital**

BBVA Colombia ha sido pionera en la transformación digital del sector bancario en el país. A través de diversas iniciativas, la entidad ha promovido el uso de canales digitales para mejorar la experiencia de sus clientes y ofrecer servicios más eficientes. Un ejemplo destacado es la campaña Uga Uga, diseñada

para impulsar el cambio hacia la banca digital, la cual fue galardonada por su impacto positivo en la adopción de la tecnología por parte de los colombianos. [3]

### 3.1.5. Inversión en Tecnología y Desarrollo Sostenible

En 2023, BBVA Colombia realizó una inversión récord de más de \$235.000 millones en tecnología, con el fin de renovar sus aplicaciones móviles y expandir su oferta de productos 100 % digitales. Esta apuesta por la innovación digital ha fortalecido la competitividad del banco en un mercado cada vez más digitalizado. Además, ese mismo año, BBVA Colombia registró un crecimiento del 98 % en financiación sostenible, movilizando \$6,7 billones para proyectos con impacto social y medioambiental, reafirmando su compromiso con el desarrollo sostenible y el bienestar de las comunidades. [4]

## 3.2. Análisis de Texto y Modelado de Tópicos para la Identificación de Áreas de Insatisfacción en Clientes Detractores

El análisis de texto y el modelado de tópicos son herramientas clave para comprender las preocupaciones de los clientes detractores. Este marco teórico describe los conceptos principales utilizados en este estudio, con un enfoque en cómo estas técnicas permiten identificar áreas de insatisfacción y mejorar la toma de decisiones estratégicas.

### 3.2.1. Técnicas de Análisis de Texto

El análisis de texto es el proceso de extraer información significativa de grandes volúmenes de datos no estructurados, como los comentarios de los clientes. Las técnicas clave que se implementan en este estudio incluyen el **TF-IDF Vectorizer**, **stemming** y **lematización**, cada una de las cuales se describe a continuación [8].

**TF-IDF Vectorizer** ( $\text{TF-IDF}(t, d, D)$ ) : Es una técnica de procesamiento de texto que transforma los comentarios en representaciones vectoriales numéricas, permitiendo analizar la frecuencia de términos. TF-IDF (*Term Frequency-Inverse Document Frequency*) pondera la importancia de una palabra según su frecuencia en un documento, pero también teniendo en cuenta cuántos documentos contienen esa palabra.

La fórmula de TF-IDF para una palabra  $t$  en un documento  $d$  dentro de un conjunto de documentos  $D$  es:

$$\text{TF-IDF}(t, d, D) = \text{TF}(t, d) \times \text{IDF}(t, D)$$



Donde:

- $TF(t, d)$  es la *frecuencia del término*  $t$  en el documento  $d$ . Se calcula como el número de veces que aparece el término  $t$  dividido por el número total de términos en el documento.

$$TF(t, d) = \frac{\text{número de veces que el término } t \text{ aparece en el documento } d}{\text{Total de términos en el documento } d}$$

- $IDF(t, D)$  es la *frecuencia inversa del documento*, y se calcula como:

$$IDF(t, D) = \ln \left( \frac{\text{Número total de documentos en el corpus } D}{\text{número de documentos que contienen el término } t} \right)$$

**Nota:** Se puede sumar un 1 al numerador y al denominador para evitar divisiones entre 0.

Este factor pondera menos las palabras que son comunes en muchos documentos.

Por ejemplo, consideremos dos documentos (comentarios):

- Documento 1: “*El producto llegó tarde*”
- Documento 2: “*El servicio es deficiente*”

Después de eliminar las *stopwords* como “el” y aplicar lematización, obtenemos:

- Documento 1: “*producto*”, “*llegar*”, “*tarde*”
- Documento 2: “*servicio*”, “*deficiente*”

El vocabulario total es: {“*producto*”, “*llegar*”, “*tarde*”, “*servicio*”, “*deficiente*”}.

Calculamos la TF (Frecuencia de Término) para cada palabra en cada documento:

Para Documento 1:

$$TF(\text{término}, \text{Doc1}) = \frac{\text{Frecuencia del término en Doc1}}{\text{Total de términos en Doc1}} = \frac{1}{3} \approx 0,333$$

Para Documento 2:

$$TF(\text{término}, \text{Doc2}) = \frac{1}{2} = 0,5$$

Calculamos la IDF (Frecuencia Inversa de Documento) para cada término:

$$IDF(\text{término}, D) = \ln \left( \frac{\text{Número total de documentos}}{\text{Número de documentos que contienen el término}} \right)$$

Aplicando esto:

- Para “producto”:  $IDF = \ln\left(\frac{2}{1}\right) = 0,693$
- Para “llegar”:  $IDF = \ln\left(\frac{2}{1}\right) = 0,693$
- Para “tarde”:  $IDF = \ln\left(\frac{2}{1}\right) = 0,693$
- Para “servicio”:  $IDF = \ln\left(\frac{2}{1}\right) = 0,693$
- Para “deficiente”:  $IDF = \ln\left(\frac{2}{1}\right) = 0,693$

Calculamos el TF-IDF multiplicando TF por IDF para cada término en cada documento:

Para Documento 1:

- $TF-IDF(\text{producto}, \text{Doc1}) = 0,333 \times 0,693 \approx 0,231$
- $TF-IDF(\text{llegar}, \text{Doc1}) = 0,333 \times 0,693 \approx 0,231$
- $TF-IDF(\text{tarde}, \text{Doc1}) = 0,333 \times 0,693 \approx 0,231$

Para Documento 2:

- $TF-IDF(\text{servicio}, \text{Doc2}) = 0,5 \times 0,693 \approx 0,347$
- $TF-IDF(\text{deficiente}, \text{Doc2}) = 0,5 \times 0,693 \approx 0,347$

Los vectores TF-IDF resultantes son:

- Documento 1: [0,231 (producto), 0,231 (llegar), 0,231 (tarde), 0, 0]
- Documento 2: [0, 0, 0, 0,347 (servicio), 0,347 (deficiente)]

Estos vectores numéricos representan la importancia de cada término en el contexto del corpus y son utilizados por algoritmos de Machine Learning para el análisis y clasificación de textos.

**Stemming** : El stemming es un proceso para reducir las palabras a su raíz o forma base. En español, por ejemplo, las palabras “*amando*”, “*amar*” y “*amaba*” se reducirían a la raíz “*am*”, lo que ayuda a simplificar el análisis eliminando variaciones gramaticales que no añaden significado adicional. Aunque el stemming puede resultar útil, no siempre conserva el significado exacto de las palabras, ya que recorta las terminaciones sin considerar el contexto.

**Lematización** : A diferencia del stemming, la lematización transforma una palabra en su forma base (o lema) teniendo en cuenta el contexto y la gramática. Por ejemplo, la palabra “*corriendo*” se lematiza a “*correr*”. En el análisis de texto, la lematización es fundamental para agrupar términos relacionados, mejorando la precisión de los resultados.

### 3.2.2. Modelado de Tópicos con LDA

El **Latent Dirichlet Allocation (LDA)** es una técnica de modelado de tópicos que permite identificar los temas más comunes en un conjunto de documentos. El LDA asume que cada documento es una mezcla de varios temas y que cada tema está compuesto por un conjunto de palabras clave. Por ejemplo, en los comentarios de clientes, podríamos identificar dos temas recurrentes: “tiempos de entrega” y “calidad del producto”, a partir de términos como “*tarde*”, “*demora*” (para el primer tema) y “*defectuoso*”, “*malo*” (para el segundo).

El proceso de LDA implica asignar probabilidades a palabras y temas en un conjunto de documentos, lo que facilita el agrupamiento de comentarios en tópicos relevantes. Este enfoque es esencial para encontrar patrones comunes de insatisfacción entre los clientes detractores.[9]

### 3.2.3. Evaluación de Coherencia de Modelos

Para garantizar que los temas identificados sean útiles y consistentes, se calcula la coherencia de los modelos generados. La coherencia mide la interpretabilidad de los temas, asegurando que las palabras agrupadas dentro de cada tema tengan un significado coherente entre sí. Por ejemplo, un tema con palabras como “*defectuoso*”, “*malo*”, “*roto*” tendría una coherencia alta porque las palabras están estrechamente relacionadas con problemas de calidad del producto.[10]

### 3.2.4. Visualización de Resultados y Toma de Decisiones Estratégicas

Una vez identificados los temas clave, se generan visualizaciones que permiten comprender la distribución de los mismos en los comentarios. Estas visualizaciones ofrecen una representación gráfica de la prevalencia de cada tema, facilitando la toma de decisiones estratégicas.

Por ejemplo, si un tema relacionado con “demoras en la entrega” es predominante, esto sugiere un área de oportunidad para mejorar la logística de la empresa. La interpretación de estos resultados es fundamental para proponer soluciones que mejoren la experiencia del cliente y reduzcan el número de detractores.

## 3.3. Descripción de los Datos

Este estudio se basa en un conjunto de 398 registros de quejas de clientes de una entidad bancaria ficticia. Es importante aclarar que, por motivos de seguridad, no se utilizan datos reales de bancos como BBVA Colombia. En su lugar, se optó por descargar una base de datos libre y pública [7], que contiene datos ficticios para fines de análisis y estudio.

Cada registro incluye información detallada sobre la interacción del cliente con el banco, su queja específica, el tiempo de resolución y una medida posterior de satisfacción. Estos datos son esenciales para identificar áreas de mejora en la atención al cliente y optimizar los servicios que ofrece la institución.

Cada queja está asociada a un producto o servicio, como *“sucursal web”*, *“cajero automático”* o *“tarjeta de crédito o débito”*. Esta categorización permite un análisis preciso de las áreas que generan mayor frustración en los clientes, facilitando la focalización de los esfuerzos de mejora en los productos que presentan más problemas. Por ejemplo, las quejas frecuentes relacionadas con las *“tarjetas de crédito o débito”* podrían señalar la necesidad de revisar los procesos de emisión o reemplazo de tarjetas.

Uno de los indicadores clave es la columna *“Days\_To\_Resolve”*, que mide en días el tiempo que tarda en resolverse cada queja. Un tiempo de resolución prolongado puede agravar la insatisfacción del cliente y afectar negativamente la percepción de la calidad del servicio. Reducir este tiempo podría ser un factor crítico para mejorar la satisfacción general del cliente.

El conjunto de datos también incluye el *“NPS Response”* (Net Promoter Score)[5], un índice que mide la disposición de los clientes a recomendar los servicios del banco. Este indicador es fundamental para evaluar el éxito de las estrategias de atención al cliente. Un puntaje bajo (de 0 a 6, clientes detractores) puede revelar problemas importantes en la calidad del servicio o en la resolución de quejas, mientras que un puntaje alto (9 o 10, promotores) indica que el cliente quedó satisfecho con su experiencia, a pesar de haber presentado una queja.

Por último, las descripciones textuales de las quejas, incluidas en la columna *“Complaint”*, representan un recurso valioso para aplicar técnicas de análisis de texto, como TF-IDF y modelado de tópicos. Estas técnicas permiten identificar palabras clave y temas recurrentes que impactan la satisfacción del cliente, revelando patrones o problemas sistémicos que no son evidentes a simple vista pero que requieren atención urgente.

## 4. Creación y Desarrollo del Modelo

El desarrollo del modelo para la identificación de temas en los comentarios de clientes se realizó con ayuda del lenguaje de programación *Python* utilizando diversas herramientas y bibliotecas de procesamiento de lenguaje natural, como *NLTK*, *spaCy*, *scikit-learn*, y *Gensim*.

A continuación, se detalla el proceso de creación del modelo en etapas bien definidas, desde la preprocesamiento de los datos hasta la selección del número óptimo de temas utilizando el algoritmo **Latent Dirichlet Allocation** (LDA).

### 4.1. Descripción del Modelo y Proceso de Implementación

El flujo de trabajo seguido para desarrollar el modelo incluye los siguientes pasos:

1. **Limpieza de datos:** El proceso de limpieza de datos es una de las fases más críticas y que mayor tiempo consume en el análisis de comentarios. En general, los datos recopilados, como los de encuestas o comentarios en redes sociales, suelen estar llenos de errores ortográficos, gramaticales y, en ocasiones, incoherencias, principalmente por la rapidez con la que los usuarios completan las encuestas o publican sus opiniones. En este caso, la base de datos descargada ya estaba limpia, lo cual es poco común. Sin embargo, dado que no podemos utilizar datos bancarios sensibles, se detallará el proceso de limpieza tal como se aplicaría en situaciones habituales utilizando bibliotecas de procesamiento de lenguaje natural (NLP) en *Python*.
  - **Corrección ortográfica y gramatical:** El primer paso es revisar los datos para corregir errores. Para ello, se utilizarán bibliotecas de NLP como *NLTK* y *spaCy*. Por ejemplo, si un cliente escribe “el serviciio es malo” o “no me guto la atención”, estos se corrigen a “el servicio es malo” y “no me gustó la atención”, respectivamente. Esta etapa también incluye el uso de expresiones regulares (con la biblioteca *re*) para identificar y corregir palabras mal escritas.
  - **Unificación de idioma:** Los comentarios pueden estar en varios idiomas, por lo que, tras las correcciones, se traduce todo a un solo idioma. Para este análisis, se optó por el español, aprovechando que muchas herramientas como *spaCy* tienen soporte en este idioma. Esto asegura uniformidad en el análisis y facilita el uso de técnicas de procesamiento de texto en un solo lenguaje.
  - **Eliminación de ruido:** Si los comentarios provienen de redes sociales, se eliminan elementos irrelevantes como hashtags, emoticonos o stickers, que no aportan valor al análisis. Para ello, se emplea la biblioteca *re* para eliminar caracteres especiales. Por ejemplo, un comentario como “#MalServicio 😞 ” se reduce a “Mal Servicio” eliminando hashtags y emoticonos.

- **Limpieza adicional:** Se eliminan números, puntuación innecesaria y espacios en blanco adicionales para asegurar consistencia en los datos. Este paso utiliza expresiones regulares (biblioteca *re*) para eliminar caracteres no alfabéticos y homogenizar los datos, facilitando la extracción de significado relevante.
- **Normalización:** Por último, se normalizan los datos convirtiendo todo a minúsculas y homogenizando los formatos, lo cual asegura consistencia en el texto. La biblioteca *unicodedata* es útil para eliminar acentos y caracteres especiales, garantizando un análisis uniforme.
- **Lematización y Stemming:** Se utiliza *spaCy* para lematizar, es decir, reducir las palabras a su forma base (por ejemplo, “corriendo” a “correr”), y el stemmer *SnowballStemmer* de *NLTK* para reducir palabras a sus raíces (por ejemplo, “correr” y “corre” a “corr”). Esto mejora la precisión y la simplicidad en el análisis posterior.

Este proceso garantiza que los datos estén en condiciones óptimas para su posterior análisis y modelado.

2. **Preprocesamiento del Texto:** Esta etapa comienza con la eliminación de palabras irrelevantes, conocidas como *stopwords*, que son palabras comunes que no aportan información significativa, como artículos, preposiciones y conjunciones. Ejemplos de *stopwords* en español incluyen “el”, “la”, “y”, “de”, “en”, “con” y “para”. A continuación, se aplican técnicas de lematización y *stemming* para reducir las palabras a su forma base o raíz. La lematización utiliza herramientas como *spaCy* para transformar cada palabra a su lema, teniendo en cuenta su contexto gramatical. El *stemming*, por otro lado, con el apoyo de *NLTK*, reduce las palabras a sus raíces, sin considerar el contexto, eliminando únicamente sufijos y terminaciones comunes. Estas técnicas ayudan a homogenizar el texto y mejorar la precisión del análisis de temas.
3. **Creación del Corpus:** Tras el preprocesamiento de los textos, se procede a la creación del *corpus* y del diccionario de términos. Primero, los comentarios se tokenizan, es decir, se dividen en palabras individuales (tokens) y se eliminan las *stopwords* utilizando la biblioteca *spaCy*. Por ejemplo, la frase “El servicio fue excelente y rápido” se tokeniza y elimina las *stopwords* para quedar como: “servicio”, “excelente”, “rápido”. Luego, se genera un diccionario que contiene cada palabra única presente en los comentarios junto con su frecuencia de aparición. Este diccionario es la base para transformar el texto en su representación numérica.
4. **Aplicación de TF-IDF:** Para transformar los comentarios en representaciones vectoriales, se aplicó la técnica *TF-IDF* (Term Frequency-Inverse Document Frequency), que pondera la importancia de cada palabra en función de su frecuencia dentro de un comentario específico y su aparición general en todo el conjunto de datos. La frecuencia de término (TF) mide cuántas

veces aparece una palabra en un comentario, mientras que la frecuencia inversa de documentos (IDF) reduce el peso de palabras comunes que aparecen en la mayoría de los comentarios, destacando así términos más relevantes.

Por ejemplo, la frase “El servicio fue excelente y rápido” se procesa eliminando *stopwords*, quedando como: “servicio”, “excelente”, “rápido”. Luego de aplicar TF-IDF, podría representarse numéricamente como un vector [0.3, 0.5, 0.3], donde cada valor refleja la importancia de las palabras “servicio”, “excelente” y “rápido” en el contexto del *corpus*. Esta técnica permite que el modelo priorice palabras significativas en lugar de aquellas que son meramente funcionales, facilitando un análisis más preciso de temas y patrones en los comentarios.

5. **Modelado de Tópicos con LDA:** Para identificar los temas subyacentes en los comentarios, se utilizó el modelo de *Latent Dirichlet Allocation* (LDA), implementado mediante la biblioteca Gensim de Python. Este algoritmo asume que cada comentario es una combinación de varios temas y que cada tema está representado por un conjunto de palabras clave con ciertas probabilidades de aparición. A partir de esta premisa, el modelo clasifica automáticamente los comentarios en grupos temáticos relevantes.

El número óptimo de temas fue determinado *evaluando la coherencia* de cada modelo generado con diferentes cantidades de temas, utilizando también la biblioteca Gensim. La coherencia mide cuán interpretables y consistentes son las palabras agrupadas en cada tema, lo cual ayuda a identificar la estructura subyacente de los comentarios. Por ejemplo, después de entrenar el modelo, algunos temas identificados podrían estar relacionados con “tiempo de espera” o “calidad del servicio”, destacando palabras como “espera”, “rápido” y “demora” para el primer tema, y “calidad”, “malo” y “atención” para el segundo.

Este enfoque permite al modelo LDA descomponer el conjunto de comentarios en temas representativos, facilitando la comprensión de los principales problemas o áreas de insatisfacción expresados por los clientes. Así, se puede analizar con mayor claridad cuáles son los factores recurrentes de insatisfacción y priorizar las acciones correctivas en las áreas que requieren atención.

6. **Visualización de Resultados:** Finalmente, se generaron visualizaciones para mostrar la distribución de los temas en los comentarios, lo cual permite identificar las áreas más recurrentes de insatisfacción entre los clientes.

## 4.2. Pseudocódigo del Modelo

A continuación, se presenta el pseudocódigo que describe el proceso de desarrollo del modelo:

1. Descargar y configurar las herramientas de procesamiento de texto:

- Instalar las bibliotecas necesarias: `kneed`, `spacy`, y `googletrans`.
  - Descargar el modelo de lenguaje en español de `spaCy`.
  - Importar las stopwords en español utilizando `NLTK`.
  - Configurar el lematizador en español con `spaCy`.
  - Configurar el stemmer utilizando `SnowballStemmer` de `NLTK`.
2. Preparación y traducción de los datos:
- Cargar el archivo de datos de comentarios.
  - Traducir las columnas relevantes al español utilizando la biblioteca `googletrans`.
  - Eliminar columnas innecesarias del conjunto de datos para optimizar el análisis.
3. Preprocesar los comentarios de la categoría 'otro':
- Filtrar los comentarios por la categoría específica.
  - Aplicar lematización y stemming en cada palabra de los comentarios:
    - Lematizar cada palabra si no es una palabra vacía o un signo de puntuación, utilizando `spacy`.
    - Aplicar stemming a las palabras resultantes con `SnowballStemmer`.
4. Crear un diccionario y corpus para el modelo LDA:
- Tokenizar y limpiar los comentarios para generar una lista de listas de palabras.
  - Crear un diccionario a partir de los textos preprocesados con `Gensim`.
  - Convertir los textos en un formato numérico (corpus) basado en el diccionario.
5. Aplicar TF-IDF para representar los textos como vectores:
- Utilizar TF-IDF para ponderar las palabras según su frecuencia e importancia en el texto.
  - Establecer un número máximo de palabras a incluir en el análisis para evitar términos raros.
6. Desarrollar el modelo LDA:
- Para un rango de posibles números de temas (por ejemplo, de 3 a 6):
    - Entrenar el modelo LDA utilizando `Gensim` con el número de temas seleccionado.
    - Obtener los términos más representativos de cada tema.
    - Calcular la coherencia del modelo usando la métrica '`c_v`' de `Gensim`.
  - Seleccionar el modelo con el mayor valor de coherencia.
7. Visualización de Resultados:
- Generar gráficos de coherencia en función del número de temas.
  - Visualizar la distribución de los temas en los comentarios.
  - Etiquetar y mostrar las principales palabras de cada tema en gráficos para facilitar la interpretación.



## 5. Resultados del Modelo de Tópicos

El análisis de los comentarios en la categoría “otro” utilizando el modelo **Latent Dirichlet Allocation (LDA)** permitió identificar los cinco principales temas recurrentes en las quejas de los clientes. La selección de esta categoría fue estratégica, ya que contenía el mayor número de comentarios, lo que la hace ideal para analizar las principales áreas de insatisfacción.

### 5.1. Número Óptimo de Tópicos

Como se muestra en la Figura 1, la evaluación de coherencia del modelo arrojó que el número óptimo de temas es **5**, ya que con este número se maximiza la coherencia de los temas identificados.

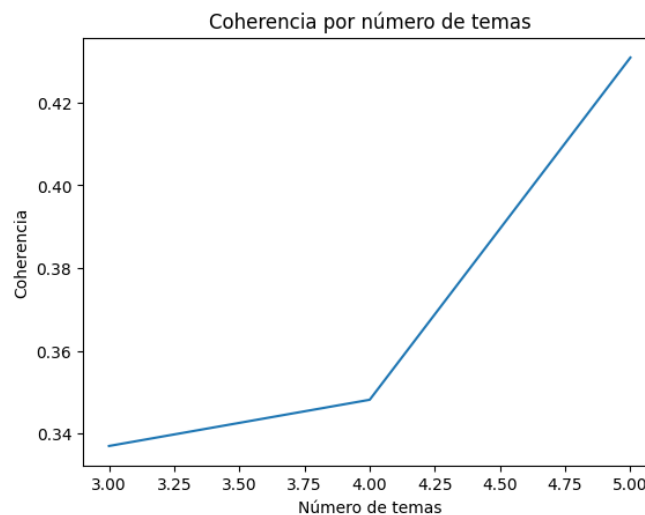


Figura 1: Coherencia por número de temas.

### 5.2. Distribución de los Temas Identificados

En la Figura 2 se presenta la distribución de los temas en los comentarios analizados. El tema 1 es el más frecuente, con 19 comentarios, seguido de los temas 4 y 5 con 17 y 16 comentarios, respectivamente.

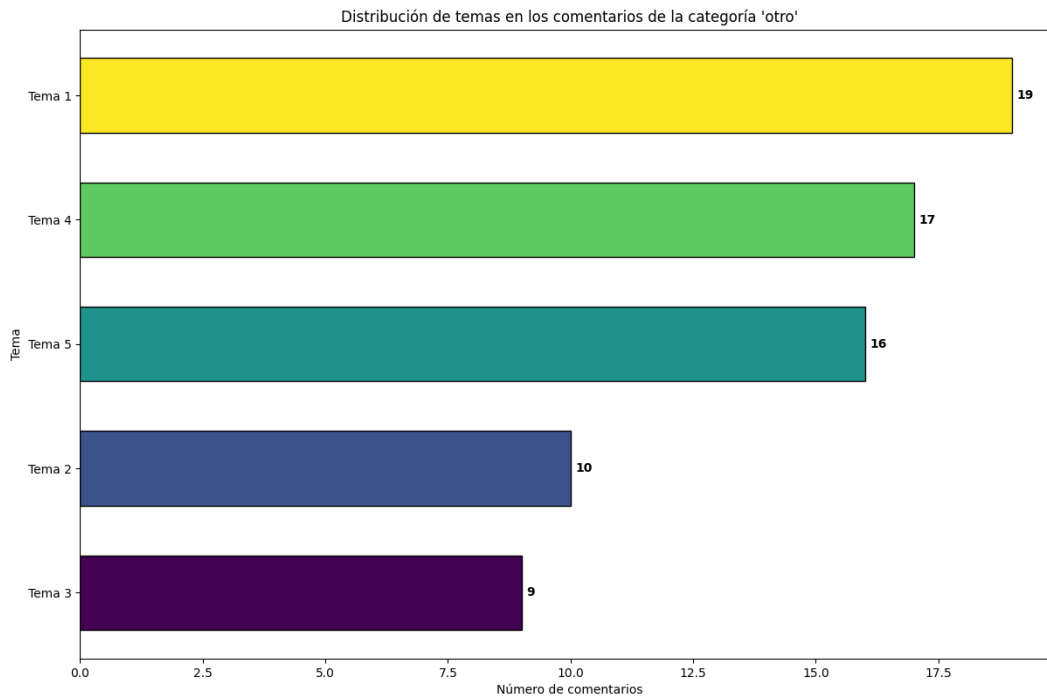


Figura 2: Distribución de los temas en los comentarios de la categoría “otro”.

### 5.3. Interpretación de los Tópicos

A continuación, se presenta una interpretación de los cinco temas más recurrentes encontrados en los comentarios:

- **Tema 1:** Este tópico refleja una fuerte insatisfacción con el servicio al cliente. Palabras como “*decepción*”, “*inaceptable*”, “*respuesta*” indican que los clientes no están recibiendo la asistencia adecuada o en el tiempo esperado.
- **Tema 2:** Las quejas en este tópico giran en torno a la actitud y comportamiento de los representantes del banco. Términos como “*grosería*”, “*relación*”, “*insatisfactorio*” sugieren que muchos clientes perciben interacciones negativas o poco profesionales.
- **Tema 3:** Este tema está relacionado con el tiempo de espera y la rapidez en la atención. Palabras como “*espera*”, “*minutos*”, “*rápido*” muestran que la demora en los procesos es una fuente recurrente de quejas entre los clientes.
- **Tema 4:** Este tópico se centra en los problemas relacionados con la interacción en sucursales

y con el personal del banco. Palabras como *“personal”*, *“sucursal”*, *“inquietud”* sugieren que hay deficiencias en el servicio personalizado y en la resolución de dudas o problemas.

- **Tema 5:** Este tema hace referencia a problemas técnicos con las aplicaciones bancarias y las transacciones móviles. Términos como *“aplicación”*, *“transacción”*, *“fallo”* indican que los usuarios experimentan problemas recurrentes al intentar realizar transacciones a través de la aplicación del banco.

En resumen, el modelo de LDA ha permitido identificar áreas clave donde los clientes muestran insatisfacción, tanto en la calidad del servicio al cliente como en problemas técnicos con las plataformas digitales. Estos hallazgos proporcionan una base sólida para diseñar estrategias de mejora enfocadas en reducir las quejas y mejorar la experiencia del cliente.

### 5.4. Recomendaciones

Con base en los resultados obtenidos y el análisis realizado, se proponen las siguientes recomendaciones para BBVA y futuros estudios relacionados con el análisis de comentarios de clientes:

- **Expansión del análisis a otros canales de retroalimentación:** Si bien este estudio se enfocó en comentarios textuales de detractores provenientes de encuestas y redes sociales, sería valioso ampliar el análisis a otros canales de interacción con el cliente, como chats en línea, llamadas al centro de atención y correos electrónicos. Esto permitiría captar una visión más integral de la experiencia del cliente, ya que diferentes canales pueden revelar problemas o inquietudes adicionales.
- **Integración de encuestas de satisfacción instantáneas:** Implementar encuestas de satisfacción en tiempo real, justo después de una interacción con el cliente (por ejemplo, tras el uso de la aplicación móvil o luego de completar una transacción), podría proporcionar datos inmediatos y más precisos sobre la satisfacción del cliente. Esto ayudaría a BBVA a identificar y corregir problemas de forma proactiva.
- **Desarrollo de modelos predictivos basados en NPS:** Se recomienda utilizar los resultados del análisis de los comentarios junto con el Net Promoter Score (NPS) para desarrollar modelos predictivos que permitan anticipar el comportamiento de los clientes. Estos modelos podrían ayudar a predecir la probabilidad de que un cliente se convierta en detractor o promotor, y permitirían tomar medidas preventivas para reducir la tasa de insatisfacción.
- **Exploración de técnicas de análisis más avanzadas:** Aunque el modelo LDA ha sido útil para identificar los temas principales, sería interesante explorar técnicas más avanzadas como modelos basados en embeddings (por ejemplo, BERT o Word2Vec) que capturan mejor el contexto

semántico de los comentarios. Estos modelos podrían proporcionar una visión más detallada de las preocupaciones de los clientes y ayudar a mejorar la clasificación de los comentarios.

- **Automatización de la retroalimentación:** Implementar sistemas automatizados que analicen los comentarios de forma continua, proporcionando reportes periódicos a la alta gerencia. Esta automatización permitiría una respuesta más rápida a los problemas emergentes y contribuiría a un ciclo de mejora continua.
- **Futuros estudios académicos y colaboraciones:** Para mejorar la robustez de los modelos desarrollados, sería beneficioso fomentar colaboraciones con instituciones académicas, que podrían explorar enfoques innovadores en la mejora del análisis de sentimiento y la identificación de patrones en grandes volúmenes de datos textuales. Además, futuros estudios podrían enfocarse en comparar la efectividad de diferentes algoritmos de análisis de texto en el contexto bancario.

En resumen, estas recomendaciones buscan ampliar el alcance de las herramientas de análisis desarrolladas en este trabajo y proponen futuros estudios para seguir mejorando la experiencia del cliente y la toma de decisiones en BBVA, con el objetivo final de reducir el número de detractores y mejorar la satisfacción general.

## 6. Anexos

Los anexos de este trabajo incluyen el código desarrollado en Google Colab: [https://colab.research.google.com/drive/1AVMlirkqufPE-BpA9\\_ZLStWFVA-XSetJ?usp=sharing](https://colab.research.google.com/drive/1AVMlirkqufPE-BpA9_ZLStWFVA-XSetJ?usp=sharing) y la Base de datos: <https://pub.towardsai.net/bank-complaints-fictional-data-b885cc907b7d>

## 7. Referencias

### Referencias

- [1] Nelson Salgado Reyes, Graciela Elizabeth Trujillo Moreno, «Análisis de sentimientos en datos de redes sociales: aplicación de técnicas de procesamiento de lenguaje natural y machine learning para analizar opiniones y sentimientos en datos de redes sociales,» 2024. URL: <https://dspace.itsjapon.edu.ec/xmlui/handle/123456789/4606>.
- [2] BBVA, «BBVA Colombia, 60 años de historia,» BBVA.com, 2023. URL: <https://www.bbva.com/es/co/bbva-colombia-60-anos-historia/>.

## REFERENCIAS

---

- [3] BBVA, «BBVA Colombia lidera la transformación digital del sector,» BBVA.com, 2023. URL: <https://www.bbva.com/es/bbva-colombia-lidera-la-transformacion-digital-del-sector/>.
- [4] BBVA, «El 2023 fue el año de hacer inversiones récord en Colombia: Mario Pardo, presidente de BBVA en Colombia,» BBVA.com, 2024. URL: <https://www.bbva.com/es/co/economia-y-finanzas/el-2023-fue-el-ano-de-hacer-inversio>
- [5] Hotjar, «Análisis del NPS: 5 formas de entender tu resultado NPS,» Hotjar.com, 2023. URL: <https://www.hotjar.com/es/net-promoter-score/analisis/>.
- [6] Elastic, «¿Qué es el análisis de sentimiento? Una guía completa del análisis de sentimiento,» Elastic.co, 2023. URL: <https://www.elastic.co/es/what-is/sentiment-analysis>.
- [7] Towards AI, «Bank Complaints Fictional Data,» Towards AI, 2023. URL: <https://pub.towardsai.net/bank-complaints-fictional-data-b885cc907b7d>.
- [8] Nazar Anchorena, C., «Extracción de patrones en las reseñas sobre celulares mediante el modelado de temas y el análisis de sentimientos,» Universidad Torcuato Di Tella, 2022. URL: [https://repositorio.utdt.edu/bitstream/handle/20.500.13098/11861/MiM\\_Nazar%20Anchorena\\_2022.pdf?isAllowed=y&sequence=1](https://repositorio.utdt.edu/bitstream/handle/20.500.13098/11861/MiM_Nazar%20Anchorena_2022.pdf?isAllowed=y&sequence=1).
- [9] David M. Blei, Andrew Y. Ng, Michael I. Jordan, «Latent Dirichlet Allocation,» Journal of Machine Learning Research, 3(2003), pp. 993-1022. URL: <https://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf>.
- [10] Florian Rosner, Andreas Hinneburg, Michael Röder, Matthias Netting, Andreas Both, «Evaluating topic coherence measures,» NIPS 2013 Workshop on Topic Models: Computation, Application, and Evaluation, 2013. URL: [https://mimno.infosci.cornell.edu/nips2013ws/nips2013tm\\_submission\\_7.pdf](https://mimno.infosci.cornell.edu/nips2013ws/nips2013tm_submission_7.pdf).