



BENEMÉRITA UNIVERSIDAD AUTÓNOMA DE PUEBLA
FACULTAD DE CIENCIAS DE LA COMPUTACIÓN

***Reconocimiento de Patrones de Repetición en
Canciones usando Técnicas de Estadística
Bayesiana***

**Tesina para obtener el título de Licenciatura en
Ingeniería en Ciencias de la Computación**

PRESENTA:

Karla Karina Gómez Matías

Director de Tesis:

Dra. María Teresa Torrijos Muñoz

FEBRERO 2024



—

Resumen

El presente trabajo consiste en el desarrollo de una aplicación de clasificación de textos automático, que permite asignar la categoría a la oración ingresada de acuerdo con el texto de entrenamiento (letras de canciones), mediante el uso de una técnica de aprendizaje automático de la familia de los clasificadores Naive Bayes, usados para la clasificación de textos.

El teorema de Bayes, formulado por Thomas Bayes, calcula la probabilidad de que ocurra un evento en función del conocimiento previo de las condiciones relacionadas con un evento.

El clasificador Naive Bayes Simple considera la probabilidad de cada término dada la clase de forma binaria, es decir que la palabra está o no está en la clase, entonces la probabilidad condicional de la palabra, dada la clase puede no estar considerada. Con la finalidad de mejorar el clasificador se utilizó el algoritmo de clasificación Naive Bayes Multinomial, ya que considera la frecuencia de las palabras, para calcular la probabilidad condicional dada la clase, mejorando la asignación de las categorías del clasificador.

Palabras claves: clasificación de textos, aprendizaje automático, Teorema de Bayes, Naive Bayes Simple, Naive Bayes Multinomial.

Índice

Capítulo 1: Introducción.....	14
1.1 Conceptos Generales.....	14
1.2 Planteamiento del Problema	15
1.3 Objetivos Generales.....	15
1.4 Objetivos Específicos del Proyecto.	16
1.5 Justificación de la Investigación	16
Capítulo 2: Estado del Arte.....	17
2.1 Review on Data Science and Prediction	17
2.2 A survey study of success factors in data science projects.....	18
2.3 Classification and Analysis of Techniques and Tools for Data Visualization Teaching	19

2.4	Estudio de demanda y empleabilidad de la carrera de ingeniería en ciencia de datos en el Ecuador	20
2.5	Clasificación automática de papers de Ciencias de la Computación	21
Capítulo 3 Marco Teórico		23
3.1	Ciencia de Datos	23
3.2	Machine Learning	24
3.3	Clasificación de Textos	25
3.4	Procesamiento del Lenguaje Natural	25
3.4.1	Extracción de características o codificación de características	26
1.	Tokenization	26
2.	Stopwords	26
3.	Modelo BoW (Bag of Words)	26
3.5	Algoritmo por Clasificación	27

3.3.1 Naive de Bayes	27
1. Crear el modelo.....	27
2. Clasificar	28
3.3.2 Matemáticas del algoritmo Naive de Bayes	28
3.3.3 Independencia Naive Bayes Multinomial	30
Capítulo 4 Desarrollo del Clasificador	30
4.1 Método de clasificación jerárquica	31
4.2 Recopilación de datos	31
4.3 Clasificación de texto en el clasificador.....	31
4.4 Implementación	32
4.5 Procesamiento del algoritmo	32
a) Conjunto de datos.....	32
b) Filtrado y Limpieza de datos	35

✓	Tokenization.....	35
✓	Stopwords	43
✓	Bag of Words	44
c)	Análisis de los Datos.....	46
✓	Tabla de frecuencias.....	46
d)	Entrenamiento.....	53
✓	Crear Modelo	53
✓	Probabilidades a priori de cada clase/categoría.....	53
✓	Suavizado de Laplace.....	54
✓	Normalizar los datos	57
✓	Palabras desconocidas	57
✓	Clasificar	59
✓	Clasificar una oración nueva.....	59

✓ Naive Bayes.....	60
✓ Continuando con la clasificar de la oración anterior aplicando NBM..	62
Capítulo 5: Datos y Aprendizaje	67
5.2 Aprendizaje Supervisado.....	89
Capítulo 6: Evaluación del desempeño del clasificador.....	91
6.1 Matriz de confusión.....	95
Métrica para la matriz de confusión.....	95
6.2 Matriz de confusión del clasificador	97
6.3 Macro promedio	100
Capítulo 7: Conclusión	101
7.1 Trabajo a futuro realizar la clasificación de sentimientos.....	102
Bibliografía.....	103

Índice de tablas

Tabla 1 Ejemplos de aplicaciones de aprendizaje automático	24
Tabla 2 Total del conjunto de datos.....	32
Tabla 3 Títulos de las canciones (data set)	33
Tabla 4 Pruebas del clasificador.....	88
Tabla 5 Matriz de confusión.....	95
Tabla 6 Casos posibles	97
Tabla 7 Matriz de confusión NB.....	98
Tabla 8 Matriz de confusión NBM y la evaluación.....	99
Tabla 9 Valores de precisión y recall de las cinco categorías.	100

Índice de Figuras

Fig. 1 Arquitectura del Sistema	32
Fig. 2 Datos de entrenamiento (data set)	34
Fig. 3 Método separar palabras	43
Fig. 4 El switch convierte la categoría a mayúsculas	44
Fig. 5 Método ingresarPalabra separa la cadena en tokens y asigna un índice a cada una.	45
Fig. 6 Tabla de ocurrencia de las 9238 palabras del texto de entrenamiento	51
Fig. 7 Conteo de textos por categoría.	52
Fig. 8 En VerCategorias se imprime la cantidad de textos que tiene cada categoría, en totalTextos nos da el total de textos de entrenamiento.	52
Fig. 9. Calcular probabilidad por categoría.....	53
Fig. 10 Probabilidad a priori de las clases.....	53
Fig. 11. Probabilidad de la palabra climb sin suavizado.....	55

Fig. 12 Probabilidad de las palabras never climb sin suavizado	55
Fig. 13. Incrementa en uno las palabras blue y climb con suavizado.	56
Fig. 14. Ocurrencia de this y real por categoría.	57
Fig. 15 Realiza el suavizado de Laplace en los token´s this y real.	58
Fig. 16 Probabilidades del token this en las cinco categorías.	59
Fig. 17 Método clasificar	61
Fig. 18: Código para la impresión de los resultados	62
Fig. 19. Probabilidad de la oración ingresada para la categoría 0.	64
Fig. 20 Probabilidad de la oración ingresada para la categoría 1.	64
Fig. 21 Probabilidad de la oración ingresada para la categoría 2.	65
Fig. 22 Probabilidad de la oración ingresada para la categoría 3.	65
Fig. 23 Probabilidad de la oración ingresada para la categoría 4	66
Fig. 24. Resultado de la clasificación.	66

Fig. 25 Arquitectura del Sistema	89
Fig. 26. Dentro de los if anidados se ingresaron los valores verdaderos positivos de la tabla 3.	90
Fig. 27 Clase Archivo, método leer imprime las líneas del dataset.	90
Fig. 28 Frecuencia de los token´s de la oración	91
Fig. 29. Probabilidad de los token´s dentro de la categoría 3.	92
Fig. 30 Probabilidad de la categoría 3.....	92
Fig. 32 Probabilidades para la canción	93
Fig. 33 .Letra de la canción: I Will Be Hero by Tiesto.....	94

Índice de Listas

Lista 1 9238 palabras/token´s.....	42
Lista 2 Stop words	43

Capítulo 1: Introducción

La Ciencia de Datos analiza grandes cantidades de datos con base a diferentes modelos que aporten información, que en su estado normal no se ven a simple vista. El origen de los datos es diverso pueden ser datos tabulares, imágenes, sonidos, texto y video.

La Ciencia de Datos recopila, procesa y analiza los datos para encontrar soluciones eficaces, uno de los propósitos de este documento es conocer algunas investigaciones que se han realizado en torno a la ciencia de datos e implementarlos en una aplicación.

1.1 Conceptos Generales

La Ciencia de Datos es el proceso de construir un modelo representativo que se ajuste a los datos de observación, según el problema se pueden clasificar en tareas tales como clasificación, análisis de asociación, agrupamiento y regresión. Cada tarea de ciencia de datos utiliza algoritmos de aprendizaje específicos como arboles de decisión, redes neuronales artificiales, K- vecinos más cercanos (K-NN) entre otros.

1.2 Planteamiento del Problema

Día a día estamos en contacto con sistemas que recolectan información con el objetivo de modelar nuestro comportamiento, el convertir esa información en datos requiere del uso de herramientas estadísticas y computacionales siendo la Ciencia de Datos la encargada de ello.

Hoy en día las empresas comienzan a apostar por la Ciencia de Datos y el uso de tecnologías como el Big Data, Data Analytics y Machine Learning entre otras herramientas para rentabilizar los datos y predecir una variedad de situaciones comerciales.

Desafortunadamente muchas empresas están desperdiciando el valor de los datos por desconocimiento, burocracia o incapacidad de adaptación a las nuevas tecnologías, falta de organización o de profesionales capacitados en la materia.

1.3 Objetivos Generales.

Implementar el uso de un algoritmo Naive Bayes que permita identificar patrones de repetición en canciones, esto mediante el uso de algoritmo de clasificación de textos.

1.4 Objetivos Específicos del Proyecto.

- ✓ Establecer los procedimientos que se usaran para formular y evaluar el clasificador.
- ✓ Desarrollar el clasificador bayesiano multinomial.
- ✓ Realizar pruebas.
- ✓ Analizar y comparar los resultados de acuerdo con la precisión, recall y puntaje F1.

1.5 Justificación de la Investigación

El análisis de datos mejorar la toma de decisiones, proporcionar nuevas perspectivas, ayuda a combatir problemas de desarrollo social y de salud.

Hoy en día nuestro país se enfrenta a gran cantidad de retos y la facilidad o dificultad de su resolución surge en gran necesidad de que herramientas tenemos a nuestra disposición. Las decisiones basadas en datos representan una ventaja competitiva para cualquier ámbito.

Lo anterior motiva este trabajo, en el que se pretende realizar un análisis de letras de canciones de diferentes géneros musicales para realizar un clasificador que pueda reconocer el género/categoría de una canción basándose únicamente en la letra.

Capítulo 2: Estado del Arte

Se analizaron cinco artículos científicos relacionados con la investigación de la ciencia de datos para conocer las áreas de oportunidad

2.1 Review on Data Science and Prediction

En el artículo Revisión sobre la Ciencia de Datos y Predicción nos menciona que la ciencia de datos ha sido muy eficaz para combatir los desafíos que enfrenta el Big Data ya que se ha producido una expansión de datos gracias a los dispositivos inteligentes, las redes sociales y la web.

Casi todos los aspectos de la vida humana se interconectaron con grandes cantidades de información creando una necesidad para las personas que pudieran administrar y rastrearla recopilación de información. Desde el punto de vista de la ingeniería la increíble escala del Big Data ha demostrado que muchos modelos de base de datos están desactualizados, estos modelos se crearon para resumir datos y acceso rápido y no para el descubrimiento de conocimientos. Están optimizados para la consulta de los usuarios y no para el descubrimiento de varios patrones complejos.

La toma de decisiones de hoy se basa en Big Data, donde las computadoras actúan como mejores tomadores de decisiones que las personas, refiriéndose como mejores en el aspecto escalabilidad, precisión y costo.

En conclusión, el Big Data realiza la extracción de grandes volúmenes de conjuntos de datos mientras que la ciencia de datos utiliza algoritmos para entrenarlos y obtener información relevante.

2.2 A survey study of success factors in data science projects

En el artículo Un estudio de encuestas sobre los factores de éxito en proyectos de ciencias de datos se realizó una encuesta a 237 profesionales de la ciencia de datos sobre el uso de metodologías de gestión de proyectos para la ciencia de datos,

Las preguntas más relevantes fueron:

Q10 ¿Sueles seguir alguna metodología de proyectos de ciencia de datos?

Q11 ¿Seleccione la metodología que conoce y/o ha utilizado para sus proyectos de ciencia de datos?

Las metodologías más mencionadas son: CRISP-DM, Microsoft TDSP, ciclos de vida Agile DS Lifecycle, Domino DS Lifecycle, IBM FMDS, RAMSY.

Con respecto a las metodologías de ciencia de datos, Agile DS Lifecycle la más utilizado seguido de CRISP-DM y Microsoft TDSP. Solo el 25% de los encuestados afirman usan algún tipo de metodología de proyectos de ciencia de datos este bajo porcentaje muestra la falta de una metodología de procesos definidos para la gestión

de proyectos de ciencia de datos siendo una de las principales causas de las fallas actuales y de los desafíos de gestión por el contrario de los profesionales que se adhieren a una metodología de proyecto ponen un mayor énfasis en los riesgos y peligros potenciales del proyecto, el control de versiones, la canalización de implementación a producción y la seguridad y privacidad de los datos.

2.3 Classification and Analysis of Techniques and Tools for Data Visualization

Teaching

En artículo Clasificación y Análisis de Técnicas y Herramientas para la Enseñanza de la Visualización de Datos su principal objetivo es mejorar la enseñanza de la visualización de los datos.

La visualización de los datos es una forma de representar gráficamente la información y los datos, se pueden usar varios gráficos diferentes para fines operativos, como mejorar la eficiencia, monitorear procesos, estudiar la distribución geográfica de datos, buscar tendencia y relaciones, su principal función está en el análisis de datos.

Presentan una nueva clasificación de las técnicas gráficas de datos en tres grupos diferentes con atributos diferentes de los gráficos:

- Uso: se definen según el uso o tratamiento que hagan de los datos crudos o procesados (Exploratorio o Explicativo).

- Forma: se pueden definir de acuerdo con los diferentes elementos gráficos utilizados para construirlos que se pueden resumir en su forma siendo el mayor atributo utilizado en la clasificación de gráficos de datos, existen diferentes subtipos de grafos: Hachas, círculos, mapas, rectángulos, redes y dibujos.
- Movimiento se definen según su movimiento es decir si la pantalla tiene o no elementos que varíen según el tiempo (Estático y Dinámico).

Se realizó un análisis comparativo de las principales herramientas de visualización de datos, dando como resultado que las ventajas de utilizar software libre es que tiene una mayor variedad de visualización de datos y una mayor capacidad de personalización de esta. El software con licencia, su principal ventaja es la velocidad al crear cualquier implementación.

2.4 Estudio de demanda y empleabilidad de la carrera de ingeniería en ciencia de datos en el Ecuador

Mediante una encuesta aplicada a las instituciones educativas del centro del país Ecuador, utilizando recursos virtuales, se obtuvieron como respuesta que 4,210 estudiantes y 81 instituciones públicas y privadas obteniendo como resultado una tendencia favorable para justificar la creación de la carrera de Ingeniería en Ciencia de Datos en el Ecuador.

El instrumento utilizado para la recolección fue una encuesta que consta de 18 preguntas y fue generada en línea lo que permitió obtener información sobre la

necesidad e intereses de estudiar la carrera de Ciencia de Datos, permitiendo conocer los criterios valorados entre el perfil profesional versus la oferta académica del programa,

Las tendencias de interés en el área de ciencia de datos se obtiene un interés del 57.93% en instituciones públicas, frente a un 43.03% de instituciones privadas. En cuanto a nivel estudiantil se determina que un 82.54% de los encuestados le interesaría una carrera profesional en área de ciencia de datos.

La importancia de las Tecnologías de la información y el análisis de datos en esta era hacen necesario la creación de profesionales en ciencia de datos, que sean capaces de seleccionar, preparar, analizar, evaluar y comunicar cantidades masivas de datos de cualquier tipo, de manera ética y responsable para la toma de decisiones inteligentes y la resolución de problemas complejos en los sectores científicos, tecnológicos, empresariales y sociales

2.5 Clasificación automática de papers de Ciencias de la Computación

La tesis consiste en el desarrollo de un sistema de clasificación automática de papers (artículos científicos), que permiten asignar la categoría adecuada de acuerdo con el contenido de estas características puede ahorrar a los publicadores el tiempo significativo que implica atravesar un proceso de selección de categorías jerárquicas para clasificarlos correctamente, y de esta manera facilitar su búsqueda una vez que se han publicado en los repositorios bibliográficos online.

Para la implementación de este sistema se obtuvo una colección de artículos de todo tipo del sitio web de la Librería Digital de la ACM (Association for Computing Machinery) a partir de la cual se generaron los data sets de entrenamiento que se utilizarán junto con una herramienta llamada MALLET para determinar la categoría de papers que aún no han sido clasificados, estudiando el caso particular de las últimas taxonomías que presenta dicho sitio (versiones de 1998 y 2012).

La clasificación automática de textos siempre ha sido un tema de investigación importante desde la existencia de documentos digitales, ya que ésta permite manejar la enorme cantidad de información disponible en la web. Está basada principalmente en técnicas de machine learning, las cuales construyen automáticamente un clasificador que aprende las características de las categorías a partir de un conjunto preclasificado de documentos. Estas técnicas juegan un rol muy importante en la extracción y resumen de información, recuperación de texto, entre otros.

En este trabajo se ha hecho una revisión general de los problemas de clasificación jerárquica, y un subconjunto de sus posibles soluciones. Se ha diseñado para tal fin una aplicación que permita automatizar el proceso de clasificación de artículos científicos y facilitar la configuración de los parámetros de esta, para poder comunicarse de manera sencilla con herramientas de clasificación externa, utilizando para este caso particular la implementación en Java de MALLET, la cual posee eficientes metodologías para realizar este tipo de tareas.

Capítulo 3 Marco Teórico

3.1 Ciencia de Datos

La Ciencia de Datos explora y analiza datos mediante tecnologías que provienen de las matemáticas, la estadística y la informática con el objetivo de obtener información para la toma de decisiones; entre las que encontramos el análisis exploratorio, el aprendizaje automático, el aprendizaje profundo el procesamiento del lenguaje natural, la visualización de datos y el diseño experimental.

Las dos tecnologías que más se utilizan dentro de la ciencia de datos son el Machine Learning y el Deep Learning ambas englobadas en el campo de la IA.

En ambos casos se busca la construcción de sistemas que sean capaces de aprender a resolver problemas a partir de conjuntos de datos que se entrenan de forma:

- ✓ supervisada: cuando los datos de entrenamiento están previamente etiquetados por humanos

- ✓ No supervisada cuando el conjunto de datos no está etiquetado.

Ejemplos de las tecnologías: sistemas de predicción ortográfica o traducción automática hasta los coches autónomos o los sistemas de visión artificial.

3.2 Machine Learning

El aprendizaje automático se dedica al estudio de agentes/algoritmos como: redes bayesianas, máquinas de vectores de soporte, análisis de clústeres, para analizar y procesar datos, que aprenden o evolucionan basados en su experiencia, consiste en dejar que los agentes descubran patrones recurrentes a partir de observaciones en los conjuntos de datos, estos datos pueden ser números, palabras, imágenes, estadísticas, etc.

Existe una amplia variedad de aplicaciones que utilizan agentes basados en aprendizaje entre ellas se encuentran:	
procesamiento de Lenguaje Natural	se utiliza para el análisis sintáctico y morfológico de los textos, extracción de Información, clasificación automática de documentos.
En el diagnóstico médico	según la historia clínica y los síntomas que presenta el paciente
En Biología	reconocimiento de tumores, patrones en cadenas de ADN.
En finanzas e industrias bancarias	para definir modelos de riesgo y fraude crediticio.
En el análisis de imágenes	detectar objetos dentro de una imagen como rostros, personas, etc.

Tabla 1 Ejemplos de aplicaciones de aprendizaje automático

3.3 Clasificación de Textos

La clasificación de textos consiste en determinar la categoría de un texto que establezca la clase a la que pertenece de acuerdo con ciertas características, facilitando la organización de la información

La decisión sobre que etiqueta debe asignar a las palabras se determina a partir de un modelo construido con el texto de entrenamiento.

3.4 Procesamiento del Lenguaje Natural

El procesamiento del lenguaje natural NLP permite que las computadoras interpreten el lenguaje humano.

El lenguaje puede expresarse por escrito (texto) u oralmente (voz). El NLP está más avanzado en el tratamiento de textos, ya que los datos son fáciles de encontrar en formato electrónico.

Las computadoras solo entienden de bytes y dígitos es por lo que el texto debe convertirse en números, específicamente en vectores de números.

3.4.1 Extracción de características o codificación de características

1. Tokenization

Divide el texto en pequeñas unidades como palabras o frases

2. Stop words

Las palabras generalmente se filtran antes de procesar un lenguaje natural se denominan stop words, son las palabras más comunes en cualquier idioma como artículos, preposiciones, pronombres y conjunciones, se filtran ya que no agregan mucha información al texto

3. Modelo BoW (Bag of Words)

Es un modelo que representa un texto como una bolsa de palabras, sin considerar su orden o relación entre ellas, extraer características del texto, es decir la frecuencia de palabras de un documento.

Representándose con el modelo vectorial en que cada documento es representado como un vector de dimensión igual al tamaño del vocabulario de la clase, y en que el valor de cada atributo corresponde a la frecuencia de la palabra en la clase correspondiente.

3.5 Algoritmo por Clasificación

Este algoritmo está basado en el teorema de Bayes y clasifica cada valor como independiente de cualquier otro. Permite predecir una categoría en función de un conjunto dado, utilizando la probabilidad.

3.3.1 Naive de Bayes

Es de los clasificadores más utilizados por su facilidad y rapidez. Constituye una técnica supervisada que necesita tener ejemplos clasificados para funcionar.

El funcionamiento del algoritmo destaca en dos partes:

1. Crear el modelo

- ✓ Calcula las probabilidades a priori de cada clase.
- ✓ Para cada clase, realiza un recuento de los valores de atributos que toma cada ejemplo. Se debe distribuir cada clase por separada para mayor eficiencia del algoritmo.
- ✓ Aplicar la Corrección de Laplace.
- ✓ Normalizar para obtener un rango de valores $[0,1]$.

2. Clasificar

1. Para cada clase disponible se determinan los valores de probabilidad de cada valor de los atributos del nuevo ejemplo
2. Aplicar la fórmula de Naive Bayes.

3.3.2 Matemáticas del algoritmo Naive de Bayes

Dado un vector de características $X = (x_1, x_2, \dots, x_n)$ y una variable de clase Y , el teorema de Bayes establece que:

$$P(y|X) = \frac{P(X|y) * P(y)}{P(X)} \quad (1)$$

Donde:

$P(y)$ es la probabilidad a priori, y es la probabilidad de que la hipótesis sea cierta (independiente de los datos).

$P(X|y)$ es el likelihood de una clase Y dados los datos X , probabilidad de los datos dado que la hipótesis era cierta

$P(X)$ evidencia es la probabilidad de los datos (independientemente de la hipótesis.)

$P(y|X)$ es la probabilidad posteriori de la hipótesis dado los datos, la distribución de probabilidad final para la clase.

Utilizando la regla de la cadena, la probabilidad $P(X|y)$ se puede descomponer como:

$$P(X|y) = P(x_1, x_2, \dots, x_n|y)$$

$$P(X|y) = P(x_1|x_2 \dots x_n, y) * P(x_2|x_3 \dots x_n, y) \dots P(x_n|y) \quad (2)$$

Debido a la suposición de independencia condicional de Naive, las probabilidades condicionales son independientes entre sí.

$$P(X|y) = P(x_1|y) * P(x_2|y) \dots P(x_n|y) \quad (3)$$

Entonces por independencia condicional:

$$P(y|X) = \frac{P(x_1|y) * P(x_2|y) \dots P(x_n|y) * P(y)}{P(x_1) * P(x_2) \dots P(X_n)} \quad (4)$$

Y como el denominador permanece constante para todos los valores, la probabilidad posteriori puede ser:

$$P(y|x_1, x_2, \dots, x_n) = P(y) \prod_{i=1}^n P(x_i|y) \quad (5)$$

El clasificador Naive Bayes combina este modelo con una regla de decisión. Una regla común es elegir la hipótesis que sea más probables: o la regla de decisión máxima a posteriori o MAP (máximum a posteriori hypothesis).

$$y = \operatorname{argmax}_y P(y) \prod_i P(x_i|y) \quad (6)$$

3.3.3 Independencia Naive Bayes Multinomial

La independencia condicional asume que las probabilidades de $P(x_i | y_j)$ son independientes de la clase Y .

$$P(x_1, x_2, \dots, x_n | y)$$

$$P(x_1, x_2, \dots, x_n | y) = P(x_1 | y) * P(x_2 | y) * P(x_3 | y) * \dots * P(x_n | y)$$

$$Y_{MAP} = \operatorname{argmax}_{y \in Y} P(x_1, x_2, \dots, x_n | y) P(y)$$

$$Y_{NB} = \operatorname{argmax}_{y \in Y} P(y) \prod_i^n P(x_i | y_j) \quad (7)$$

Los valores $P(x_i | y_j)$ se estiman con la frecuencia de los datos, no se hace búsqueda de hipótesis, sino que se cuenta la frecuencia de ocurrencias.

Capítulo 4 Desarrollo del Clasificador

Con el objetivo de identificar el género de la canción basándose en parte de la letra (frase, coro) de forma automática, se desarrolla una aplicación en el lenguaje de programación Java, el método para realizar la clasificación es por Naive Bayes Multinomial, dando como resultado el género musical u interprete.

4.1 Método de clasificación jerárquica

Se clasificaron las canciones por medio de una organización jerárquica teniendo 5 categorías, cada categoría contiene de 5 a 10 canciones.

4.2 Recopilación de datos

El primer paso fue obtener el conjunto de datos; todas las letras de las canciones se obtuvieron de la página web música.com.

4.3 Clasificación de texto en el clasificador

Dada una clase Y y un conjunto de palabras X de la canción a clasificar, se calcula la probabilidad de que X_i se clasifique dentro de la clase Y_i , con ello se conoce la probabilidad de X_i dada una clase y la probabilidad de la clase Y_i pero se necesita encontrar el valor máximo de la expresión para encontrar la clase en la que mejor se clasifique X_i .

La definición de las clases son las categorías/genero de las canciones y la clasificación se realiza por palabras.

4.4 Implementación

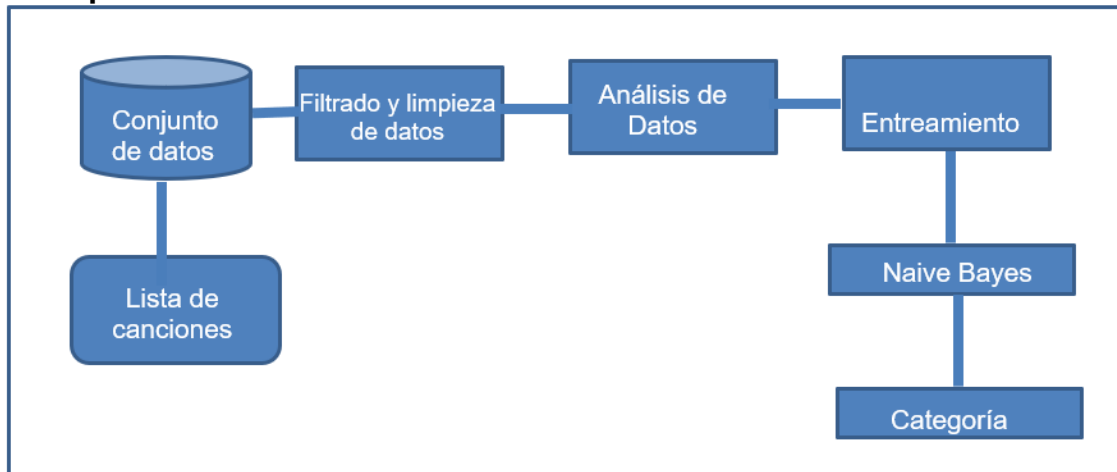


Fig. 1 Arquitectura del Sistema

4.5 Procesamiento del algoritmo

a) Conjunto de datos

Al ser una clasificación supervisada, los datos ya están previamente clasificados manualmente en 5 categorías, contando con un total de 9238 palabras, en la tabla 2 se muestra la cantidad de palabras por categorías.

Categoría	Conjunto de datos.
Queen/ Rock	1412
Eminem/ Hip Hop	4259
Daft Punk/ Electrónica	1915
Tiesto/ <u>Electropop</u>	1033
Other	619
TOTAL	9238

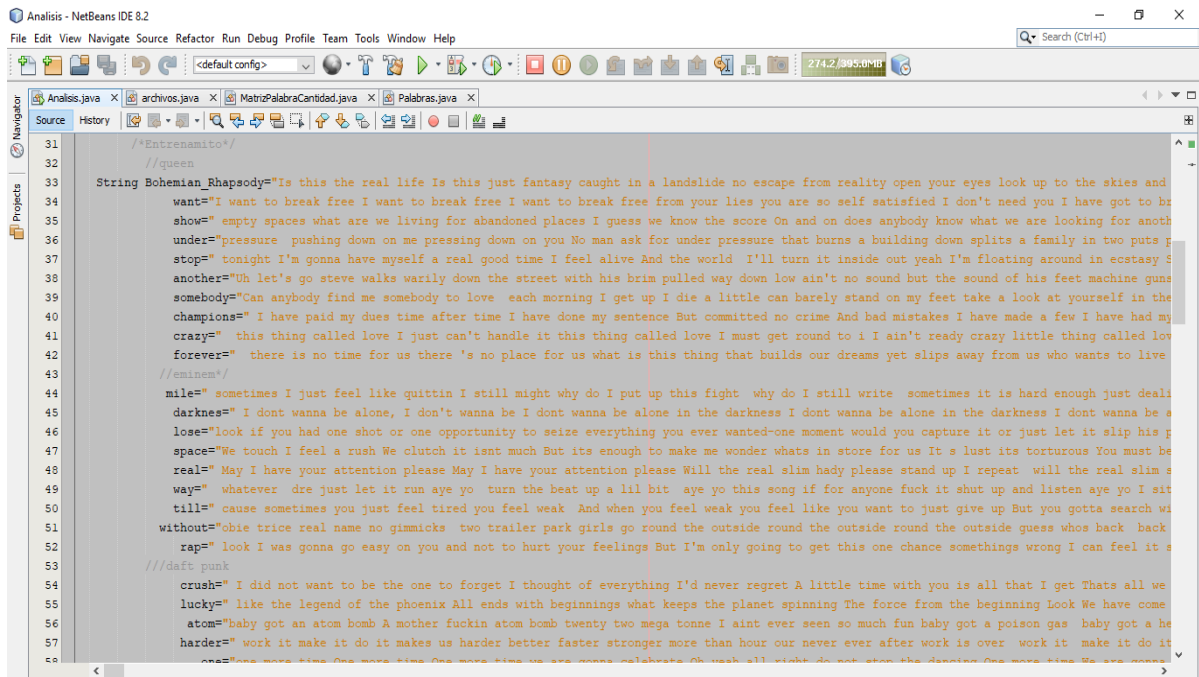
Tabla 2 Total del conjunto de datos.

Para el entrenamiento del clasificador se ingresan las letras de las canciones que se muestran en la tabla 3.

Queen/ Rock	Eminem/ Hip Hop	Daft Punk/ Electrónica	Tiesto/Electropop	Other
Bohemian rhapsody	8 mile	Instant crush	Century	BabaO'Riley
I want to break free	Darkness	Get lucky	Escapeme	I love rock 'n' roll
The show must go on	Lose yourself	Atom bomb	Feel it in my bones	Rockin' in the free world
Under pressure	Space bound	Harder, better, faster, stronger	Here on earth	Dream on
Don't stop me now	The real Slim shady	One more time	It's not the things you say	Born to be wild
Another one bites the dust	The way I am	Lose yourself to dance	I will be here	
Some body to love	Till collapse	Technologic	Red lights	
We are the champions	Without me	Face to face	Ritual	
Crazy Little thing called love	Rap god	Give life back to music	You are my diamond	
Who wants to live forever		Fragments of time		

Tabla 3 Títulos de las canciones (data set)

Los datos de entrenamiento están dentro de la clase **Análisis**, declarados como variables de tipo **String**, son un total de 43 canciones.



```
31  /*Entrenamito*/
32  //queen
33  String Bohemian_Rhapsody="Is this the real life Is this just fantasy caught in a landslide no escape from reality open your eyes look up to the skies and
34  want="I want to break free I want to break free I want to break free from your lies you are so self satisfied I don't need you I have got to br
35  show=" empty spaces what are we living for abandoned places I guess we know the score On and on does anybody know what we are looking for anoth
36  under="pressure pushing down on me pressing down on you No man ask for under pressure that burns a building down splits a family in two puts p
37  stop=" tonight I'm gonna have myself a real good time I feel alive And the world I'll turn it inside out yeah I'm floating around in ecstasy S
38  another="Uh let's go steve walks warily down the street with his brim pulled way down low ain't no sound but the sound of his feet machine guns
39  somebody="Can anybody find me somebody to love each morning I get up I die a little can barely stand on my feet take a look at yourself in the
40  champions=" I have paid my dues time after time I have done my sentence But committed no crime And had mistakes I have made a few I have had my
41  crazy=" this thing called love I just can't handle it this thing called love I must get round to i I ain't ready crazy little thing called lov
42  forever=" there is no time for us there 's no place for us what is this thing that builds our dreams yet slips away from us who wants to live
43
44  //eminem*
45  mile=" sometimes I just feel like quittin I still might why do I put up this fight why do I still write sometimes it is hard enough just deali
46  darknes=" I dont wanna be alone, I don't wanna be I dont wanna be alone in the darkness I dont wanna be alone in the darkness I dont wanna be a
47  lose="look if you had one shot or one opportunity to seize everything you ever wanted-one moment would you capture it or just let it slip his p
48  space="We touch I feel a rush We clutch it isnt much But its enough to make me wonder whats in store for us It s lust its torturous You must be
49  real=" May I have your attention please May I have your attention please Will the real slim hady please stand up I repeat will the real slim s
50  way=" whatever dre just let it run aye yo turn the beat up a lil bit aye yo this song if for anyone fuck it shut up and listen aye yo I sit
51  till=" cause sometimes you just feel tired you feel weak And when you feel weak you feel like you want to just give up But you gotta search wi
52  without="obie trice real name no gimmicks two trailer park girls go round the outside round the outside round the outside guess whos back back
53  rap=" look I was gonna go easy on you and not to hurt your feelings But I'm only going to get this one chance somethings wrong I can feel it s
54
55  //daft punk
56  crush=" I did not want to be the one to forget I thought of everything I'd never regret A little time with you is all that I get Thats all we
57  lucky=" like the legend of the phoenix All ends with beginnings what keeps the planet spinning The force from the beginning Look We have come
58  atom="baby got an atom bomb A mother fuckin atom bomb twenty two mega tonne I aint ever seen so much fun baby got a poison gas baby got a he
59  harder=" work it make it do it makes us harder better faster stronger more than hour our never ever after work is over work it make it do it
60  one="One more time One more time One more time we are gonna celebrate Oh yeah all right do not stop the dancing One more time We are gonna
```

Fig. 2 Datos de entrenamiento (data set)

En la clase **Palabras** se crea un **arrayList** llamado **palabrasL** (lista de palabras) en el cual se va añadiendo cada una de las palabras en el método **IngresarPalabra**.

El método **verPalabras** imprime el **arrayList** **palabrasL**, que es el conjunto por considerar.

b) Filtrado y Limpieza de datos

✓ Tokenization

La tokenización se realiza dentro del método **separarpalabras** con ayuda de un **Apalabras** (Arreglo de palabras) y la función Split separa la palabra de la cadena, después de un espacio o de algunos de estos caracteres [, @ '!] +

obteniendo las palabras/token's.

*****	81-sends	103-free	155-mindless	207-with
0-this	82-shivers	104-from	156-crime	208-grin
1-real	83-down	105-your	157-behind	209-giving
2-life	84-spine	106-lies	158-curtain	210-bill
3-just	85-body	107-self	159-pantomime	211-overkill
4-fantasy	86-aching	108-satisfied	160-hold	212-have
5-caught	87-goodbye	109-need	161-line	213-find
6-landslide	88-everybody	110-have	162-want	214-carry
7-escape	89-gotta	111-knows	163-take	215-pressure
8-from	90-leave	112-fallen	164-anymoreThe	216-pushing
9-reality	91-behind	113-love	165-show	217-down
10-open	92-face	114-first	166-must	218-pressing
11-your	93-truth	115-time	167-yeah	219-under
12-eyes	94-wanna	116-this	168-Inside	220-that
13-look	95-sometimes	117-know	169-heart	221-burns
14-skies	96-wish	118-real	170-breaking	222-building
15-poor	97-never	119-yeah	171-make	223-splits
16-need	98-been	120-strange	172-flaking	224-family
17-sympathy	99-born	121-true	173-smile	225-puts
18-because	100-silhouette	122-over	174-still	226-people
19-easy	101-scaramouche	123-like	175-stays	227-streets
20-come	102-will	124-sure	176-whatever	228-That
21-little	103-fandango	125-when	177-happens	229-okay
22-high	104-thunderbolt	126-walk	178-leave	230-terror
23-anyway	105-lightning	127-that	179-chance	231-knowing
24-wind	106-very	128-door	180-heartache	232-what
25-blows	107-frightening	129-baby	181-failed	233-this
26-does	108-galileo	130-life	182-romance	234-world
27-really	109-figaro	131-still	183-learning	235-about
28-matter	110-magnifico	132-goes	184-warmer	236-watching
29-mama	111-nobody	133-used	185-soon	237-some
30-killed	112-loves	134-living	186-turning	238-good
31-against	113-family	135-without	187-round	239-friends
32-head	114-spare	136-side	188-corner	240-screaming
33-pulled	115-monstrosity	137-live	189-outside	241-pray
34-trigger	116-bismillah	138-alone	190-dawn	242-tomorrow
35-dead	117-mamma	139-make	191-inside	243-gets
36-begun	118-beelzebub	140-heve	192-dark	244-higher
37-have	119-devil	141-empty	193-aching	245-chipping
38-gone	120-aside	142-spaces	194-free	246-around
39-thrown	121-think	143-what	195-soul	247-kick
40-away	122-stone	144-living	196-painted	248-brains
41-mean	123-spit	145-abandoned	197-like	249-floor
42-make	124-love	146-places	198-wings	250-these
43-back	125-baby	147-guess	199-butterflies	251-days
44-again	126-right	148-know	200-fairytales	252-never
45-time	127-outta	149-score	201-yesterday	253-rains
46-tomorrow	128-here	150-does	202-will	254-pours
47-carry	129-yeah	151-anybody	203-grow	255-screaming:
48-nothing	130-anyone	152-looking	204-never	256-high
49-matters	131-want	153-another	205-friends	257-turned
50-late	132-break	154-hero	206-face	---

258-away	319-racing	378-doorway	437-just	496-done
259-from	320-passing	379-bullets	438-relief	497-sentence
260-like	321-Lady	380-beat	439-work	498-committed
261-blind	322-Godiva	381-yeah	440-hard	499-crime
262-fence	323-there	382-another	441-works	500-mistakes
263-work	324-stopping	383-bites	442-everyday	501-made
264-keep	325-burning	384-dust	443-life	502-share
265-coming	326-hundred	385-gone	444-till	503-sand
266-with	327-degrees	386-gonna	445-ache	504-kicked
267-love	328-that	387-think	446-bones	505-face
268-slashed	329-they	388-going	447-home	506-come
269-torn	330-call	389-along	448-earned	507-through
270-Love	331-Mister	390-without	449-goes	508-mean
271-insanity	332-Fahrenheit	391-when	450-down	509-champions
272-laughs	333-travelling	392-took	451-knees	510-friends
273-breaking	334-speed	393-everything	452-start	511-keep
274-give	335-light	394-that	453-pray	512-fighting
275-ourselves	336-wanna	395-kicked	454-praise	513-till
276-more	337-make	396-happy	455-tears	514-losers
277-chance	338-supersonic	397-satisfied	456-from	515-cause
278-cause	339-such	398-long	457-eyes	516-world
279-such	340-just	399-stand	458-please	517-taken
280-fashioning	341-give	400-heat	459-wants	518-bows
281-word	342-(cause	401-look	460-help	519-curtain
282-dares	343-havin	402-take	461-every	520-calls
283-care	344-rocket	403-bite	462-everybody	521-brought
284-edge	345-ship	404-shoot	463-they	522-fame
285-night	346-Mars	405-there	464-goin	523-fortune
286-loves	347-collision	406-plenty	465-crazy	524-everything
287-change	348-course	407-ways	466-water	525-that
288-caring	349-satellite	408-hurt	467-brain	526-goes
289-last	350-control	409-bring	468-common	527-with
290-dance	351-machine	410-ground	469-sense	528-thank
291-tonight	352-ready	411-cheat	470-nobody	529-been
292-gonna	353-reload	412-treat	471-left	530-roses
293-have	354-atom	413-leave	472-believe	531-pleasure
294-myself	355-bomb	414-standing	473-yeah	532-cruise
295-real	356-about	415-repeating	474-someone	533-consider
296-good	357-explode	416-anybody	475-feel	534-challenge
297-time	358-woman	417-find	476-rhythm	535-before
298-feel	359-ball	418-somebody	477-keep	536-whole
299-alive	360-stere	419-love	478-losing	537-human
300-world	361-walks	420-each	479-beat	538-race
301-turn	362-warily	421-morning	480-alright	539-gonna
302-inside	363-down	422-little	481-gonna	540-lose
303-yeah	364-street	423-barely	482-face	541-this
304-floating	365-with	424-stand	483-defeat	542-thing
305-around	366-brim	425-feet	484-gotta	543-called
306-ecstasy	367-pulled	426-take	485-this	544-love
307-stop	368-sound	427-look	486-prison	545-just
308-cause	369-feet	428-yourself	487-cell	546-handle
309-having	370-machine	429-mirror	488-some	547-must
310-shooting	371-guns	430-Lord	489-free	548-round
311-star	372-ready	431-what	490-anywhere	549-ready
312-leaping	373-this	432-doing	491-have	550-crazy
313-through	374-hanging	433-have	492-paid	551-little
314-like	375-edge	434-spent	493-dues	552-(this
315-tiger	376-your	435-years	494-time	553-thing)
316-defying	377-seat	436-believing	495-after	554-cries

555-like	614-aside	673-slam	732-pants	791-hear
556-baby	615-dares	674-shut	733-chase	792-homey
557-cradle	616-love	675-whole	734-gotta	793-whenever
558-night	617-when	676-manhoods	735-move	794-daug
559-swings	618-must	677-been	736-asap	795-bailin
560-jives	619-touch	678-stripped	737-mommas	796-trailer
561-shakes	620-tears	679-have	738-poor	797-tomorrow
562-over	621-with	680-vicked	739-little	798-love
563-jelly	622-your	681-must	740-baby	799-kiss
564-fish	623-lips	682-split	741-sister	800-goodbye
565-kinda	624-fingertips	683-fuck	742-dont	801-whenever
566-there	625-have	684-shit	743-understand	802-need
567-goes	626-today	685-goin	744-sits	803-there
568-knows	627-sometimes	686-home	745-front	804-second
569-rock	628-just	687-world	746-buries	805-blow
570-roll	629-feel	688-shoulder	747-nose	806-everything
571-drives	630-like	689-back	748-colors	807-work
572-gives	631-quitin	690-mile	749-until	808-live
573-cold	632-still	691-road	750-crayon	809-didnt
574-fever	633-might	692-Chorus	751-gets	810-wouldnt
575-leaves	634-this	693-make	752-dull	811-deal
576-cool	635-fight	694-plan	753-hand	812-wasn
577-sweat	636-write	695-time	754-while	813-skillest
578-gotta	637-hard	696-stand	755-brother	814-borderline
579-relax	638-enough	697-travel	756-mother	815-Detroit
580-track	639-dealin	698-land	757-Aint	816-city
581-take	640-with	699-take	758-tellin	817-limits
582-back	641-real	700-matters	759-really	818-different
583-seat	642-life	701-into	760-goes	819-certain
584-Hitchhike	643-wanna	702-hands	761-head	820-significance
585-long	644-jump	703-once	762-wish	821-certificate
586-ride	645-stage	704-over	763-could	822-authenticity
587-motor	646-kill	705-tracks	764-daddy	823-even
588-bike	647-mics	706-never	765-neither	824-credibility
589-until	648-show	707-lock	766-keep	825-seen
590-hitchhike	649-these	708-gone	767-runnin	826-heard
591-Freddie	650-people	709-know	768-wanted	827-smelled
592-yeah	651-what	710-where	769-upset	828-incredible
593-there	652-level	711-sorry	770-cause	829-upon
594-time	653-skills	712-momma	771-blew	830-pedestal
595-place	654-white	713-grown	772-grew	831-unsigned
596-what	655-hate	714-alcme	773-grow	832-havin
597-this	656-somethin	715-follow	774-nuts	833-rough
598-thing	657-aint	716-footsteps	775-Dont	834-porch
599-that	658-right	717-making	776-step	835-friends
600-builds	659-brake	718-only	777-pressures	836-kick
601-dreams	660-lights	719-that	778-much	837-dumb
602-slips	661-case	720-escape	779-whats	838-rhymes
603-away	662-fright	721-from	780-best	839-serve
604-from	663-drawin	722-eminem	781-wont	840-lunchline
605-wants	664-blank	723-walkin	782-tell	841-when
606-live	665-fault	724-train	783-moment	842-comes
607-forever	666-great	725-tryin	784-pray	843-crunch
608-chance	667-then	726-regain	785-please	844-punchlines
609-decided	668-falls	727-spirit	786-beggin	845-bust
610-world	669-insides	728-fore	787-pigeon	846-flow
611-only	670-crawl	729-same	788-holed	847-another
612-sweet	671-clam	730-crap	789-regular	848-crab
613-moment	672-wham	731-plant	790-hope	

848-crab	907-suddenly	966-pacin	1024-wait	1083-alcohol
849-bucket	908-surge	967-room	1025-crowd	1084-breath
850-luck	909-burst	968-valium	1026-thought	1085-reach
851-Rabbit	910-energy	969-then	1027-shit	1086-scope
852-maybe	911-occured	970-chase	1028-sold	1087-blackin
853-outlet	912-free	971-with	1029-only	1088-meds
854-startin	913-leaders	972-boone	1030-opening	1089-then
855-doubt	914-three	973-little	1031-early	1090-benzodiazepines
856-feelin	915-third	974-taste	1032-overreact	1091-gone
857-skeptic	916-scared	975-maybe	1033-told	1092-magazines
858-hang	917-bird	976-take	1034-relax	1093-sprawled
859-clothes	918-cross	977-smoos	1035-just	1094-floor
860-about	919-median	978-tear	1036-hope	1095-media
861-salvatic	920-curb	979-stage	1037-packed	1096-goin
862-Army	921-verbs	980-fuck	1038-before	1097-yeah
863-salvage	922-blur	981-colt	1039-they	1098-people
864-outfit	923-dont	982-need	1040-fill	1099-start
865-cold	924-vanna	983-somethin	1041-each	1100-time
866-plus	925-alone	984-stronger	1042-thatd	1101-10:05
867-stuck	926-darkness	985-caps	1043-totally	1102-starts
868-battlin	927-anymore	986-better	1044-wack	1103-already
869-mode	928-hello	987-vodka	1045-murder	1104-sweatin
870-defense	929-friend	988-round	1046-nobodys	1105-locked
871-thing	930-here	989-after	1047-what	1106-rapid
872-want	931-again	990-gettin	1048-nobody	1107-fire
873-pity	932-cant	991-loaded	1049-shows	1108-spittin
874-dark	933-this	992-thats	1050-panic	1109-concert-goers
875-bein	934-hole	993-shots	1051-mode	1110-scopes
876-pulled	935-like	994-double	1052-bout	1111-sniper
877-apart	936-walls	995-entendre	1053-snap	1112-vision
878-each	937-cloain	996-starin	1054-motherfuckin	1113-surprise
879-lirbs	938-help	997-service	1055-wacko	1114-nowhere
880-skin	939-feel	998-menu	1056-second	1115-slide
882-doin	940-these	999-benzo	1057-cancel	1116-clip
883-stove	941-curtains	1000-hear	1058-fans	1117-inside
884-explode	942-open	1001-music	1059-below	1118-leanin
885-kettle	943-something	1002-continue	1060-rush	1119-going
886-mouth	944-pulls	1003-crescendo	1061-entrance	1120-Kaiser
887-overloads	945-closed	1004-whole	1062-plan	1121-sore
888-learned	946-feels	1005-fuckin	1063-wreck	1122-finger
889-turn	947-loathing	1006-venue	1064-cameras	1123-trigger
890-takes	948-vegas	1007-from	1065-directions	1124-licensed
891-burned	949-haven	1008-window	1066-press	1125-owner
892-fallin	950-vaguest	1009-when	1067-apeshit	1126-prior
893-next	951-lost	1010-know	1068-bananas	1127-convictions
894-meet	952-make	1011-youre	1069-networks	1128-says
895-girl	953-small	1012-schizo	1070-commando	1129-skys
896-longer	954-wager	1013-cause	1071-extra	1130-limit
897-play	955-tomorrows	1014-peakin	1072-clips	1131-supplies
898-stupid	956-paper	1015-curtain	1073-ammo	1132-infinite
899-immature	957-would	1016-hotel	1074-hecklers	1133-strapped
900-every	958-odds	1017-loud	1075-armed	1134-soldier
901-ingredient	959-favor	1018-almost	1076-teeth	1135-hopping
902-courage	960-much	1019-though	1077-nother	1136-over
903-already	961-father	1020-sound	1078-Valium	1137-climbing
904-beat	962-think	1021-should	1079-fall	1138-fences
905-words	963-that	1022-ready	1080-ground	1139-some
906-urge	964-knew	1023-show	1081-crawl	1140-john
	965-keep	1024-wait	1082-dresser	1141-travolta

2092-constant	2151-admired	2210-cause	2269-filling
2093-lyrical	2152-wished	2211-sometimes	2270-minimal
2094-content	2153-fired	2212-just	2271-swap
2095-guilty	2154-drop	2213-feel	2272-millions
2096-conscience	2155-from	2214-tired	2273-listeners
2097-gotten	2156-label	2215-weak	2274-coming
2098-such	2157-stop	2216-when	2275-with
2099-rotten	2158-fables	2217-like	2276-gonna
2100-responses	2159-gonna	2218-want	2277-fear
2101-controversy	2160-able	2219-give	2278-shoved
2102-circles	2161-name	2220-gotta	2279-spirit
2103-seems	2162-pigeon	2221-search	2280-lives
2104-media	2163-hold	2222-within	2281-hear
2105-immediately	2164-some	2223-find	2282-lyrics
2106-points	2165-poppy	2224-that	2283-shock
2107-finger	2166-sensation	2225-inner	2284-miracle
2108-point	2167-caught	2226-strength	2285-product
2109-index	2168-rotating	2227-pull	2286-fissing
2110-pinky	2169-rock	2228-shit	2287-shizzle
2111-ring	2170-roll	2229-motivation	2288-whizzle
2112-thumb	2171-stations	2230-quitter	2289-this
2113-give	2172-patience	2231-matter	2290-plot
2114-bullshit	2173-deal	2232-fall	2291-listen
2115-they	2174-these	2233-flat	2292-bizzles
2116-pull	2175-cocky	2234-your	2293-forgot
2117-full	2176-casions	2235-face	2294-slizzle
2118-dudes	2177-wigger	2236-collapse	2295-does
2119-bullied	2178-tries	2237-spilling	2296-fuck
2120-shoes	2179-black	2238-these	2297-roof
2121-school	2180-talk	2239-raps	2298-comes
2122-blame	2181-accent	2240-long	2299-till
2123-marylin	2182-grab	2241-drop	2300-lights
2124-heroin	2183-ball	2242-youll	2301-legs
2125-where	2184-always	2243-never	2302-cant
2126-verse	2185-keep	2244-killing	2303-shut
2127-parents	2186-askin	2245-then	2304-mouth
2128-look	2187-same	2246-stop	2305-smoke
2129-middle	2188-questions	2247-pinning	2306-clears
2130-america	2189-hood	2248-them	2307-high
2131-tragedy	2190-grew	2249-hiphop	2308-perhaps
2132-upper	2191-till	2250-eminem	2309-bone
2133-class	2192-grabbin	2251-subliminal	2310-music
2134-city	2193-hair	2252-thoughts	2311-magic
2135-havin	2194-tearing	2253-sending	2312-there
2136-happenin	2195-drivin	2254-women	2313-certain
2137-attack	2196-crazy	2255-caught	2314-feeling
2138-eminem	2197-take	2256-webs	2315-apit
2139-glad	2198-racin	2257-spin	2316-people
2140-feed	2199-pacing	2258-hock	2317-moment
2141-fuel	2200-stand	2259-venom	2318-every
2142-That	2201-thankful	2260-adrenaline	2319-single
2143-need	2202-every	2261-shots	2320-minute
2144-fire	2203-bathroom	2262-penicillin	2321-spittin
2145-burn	2204-without	2263-could	2322-trying
2146-burnin	2205-someone	2264-illin	2323-hold
2147-returned	2206-standing	2265-amoxicillin	2324-onto
2148-wasn	2207-sign	2266-real	2325-again
2149-sick	2208-autograph	2267-enough	2326-while
2150-bein	2209-asshole	2268-criminal	2327-much

3032-special	3091-hurtin	3150-cake	3209-every
3033-cable	3092-many	3151-mistake	3210-kinda
3034-channel	3093-murder	3152-fatal	3211-counted
3035-went	3094-prove	3153-need	3212-being
3036-radio	3095-nice	3154-overseas	3213-friend
3037-station	3096-sacrifice	3155-vacation	3214-give
3038-very	3097-virgins	3156-trip	3215-away
3039-next	3098-flunkie	3157-broad	3216-about
3040-kill	3099-pill	3158-fall	3217-what
3041-lyrics	3100-junky	3159-retard	3218-really
3042-supersonic	3101-accolades	3160-want	3219-know
3043-speed	3102-skills	3161-forget	3220-where
3044-sama	3103-brung	3162-thought	3221-chained
3045-lamaa	3104-Full	3163-everything	3222-myself
3046-duma	3105-myself	3164-never	3223-unlocks
3047-assuming	3106-bully	3165-regret	3224-like
3048-human	3107-mind	3166-little	3225-door
3049-superhuman	3108-million	3167-time	3226-some
3050-innovative	3109-leagues	3168-with	3227-more
3051-rubber	3110-above	3169-that	3228-master
3052-anything	3111-speak	3170-Thats	3229-they
3053-ricocheting	3112-tongues	3171-need	3230-wanted
3054-glue	3113-tongue	3172-because	3231-someone
3055-devastating	3114-cheek	3173-take	3232-locked
3056-more	3115-drunk	3174-thing	3233-just
3057-than	3116-satan	3175-same	3234-summer
3058-demonstrating	3117-wheel	3176-when	3235-memory
3059-give	3118-asleep	3177-your	3236-dies
3060-audience	3115-seat	3178-round	3237-worked
3061-levitating	3120-bumping	3179-believe	3238-long
3062-fading	3121-heavy	3180-lips	3239-hard
3063-haters	3123-funky	3181-ground	3240-sees
3064-forever	3124-something	3182-wanna	3241-right
3065-waiting	3125-tugging	3183-place	3242-through
3066-fell	3126-struggling	3184-roche	3243-easy
3067-celebrating	3127-angels	3185-gives	3244-lies
3068-motivated	3128-fight	3186-anymore	3245-cracks
3069-elevating	3129-devils	3187-once	3246-road
3070-elevator	3130-asking	3188-lock	3247-would
3071-mainstream	3131-eliminate	3189-made	3248-disguise
3072-thats	3132-women	3190-offer	3249-runs
3073-jealous	3133-hate	3191-then	3250-scissor
3074-confuse	3134-consideration	3192-this	3251-seem
3075-found	3135-bitter	3193-picture	3252-wall
3076-hella	3136-hatred	3194-kids	3253-cannot
3077-fuse	3137-patient	3195-head	3254-break
3078-rock	3138-sympathetic	3196-hear	3255-down
3079-shock	3139-situation	3197-last	3256-else
3080-throw	3140-understand	3198-said	3257-fall
3081-lose	3141-discrimination	3199-listened	3258-thousand
3082-songs	3142-life	3200-problems	3259-lonely
3083-dont	3143-handing	3201-listen	3260-stars
3084-words	3144-lemons	3202-mine	3261-hiding
3085-occurs	3145-lemonade	3203-will	3262-cold
3086-ripping	3146-cant	3204-alone	3263-sing
3087-verses	3147-batter	3205-again	3264-doesn
3088-diverse	3148-supposed	3206-cause	3265-understand
3089-curtains	3149-bake	3207-does	3266-upset
3090-inadvertently	3150-cake	3208-happen	3267-swimming

3964-find	4023-long	4082-woman	4141-dusk
3966-just	4024-till	4083-night	4142-dawn
3966-yeah	4025-with	4084-with	4143-everybodys
3967-here	4026-yeah	4085-baby	4144-their
3968-fields	4027-would	4086-hand	4145-does
3969-fight	4028-singing	4087-under	4146-life
3970-meals	4029-love	4088-light	4147-know
3971-back	4030-rock	4089-near	4148-nobody
3972-into	4031-roll	4090-garbage	4149-knows
3973-living	4032-another	4091-pute	4150-where
3974-need	4033-dime	4092-away	4151-comes
3975-prove	4034-jukebox	4093-gone	4152-goes
3976-right	4035-baby	4094-hates	4153-lose
3977-forgiven	4036-come	4095-life	4154-half
3978-raise	4037-take	4096-what	4155-lives
3979-your	4038-your	4097-done	4156-books
3980-only	4039-time	4098-more	4157-written
3981-teenage	4040-dance	4099-that	4158-pages
3982-wasteland	4041-smiled	4100-will	4159-lived
3983-sally	4042-asked	4101-never	4160-learned
3984-take	4043-name	4102-school	4161-from
3985-hand	4044-that	4103-fall	4162-fools
3986-travel	4045-matter	4104-love	4163-sages
3987-south	4046-said	4105-cool	4164-true
3988-cross	4047-cause	4106-thousand	4165-things
3989-land	4048-same	4107-points	4166-come
3990-fire	4049-home	4108-homeless	4167-back
3991-lock	4050-where	4109-kinder	4168-sing
3992-past	4051-alone	4110-gentler	4169-with
3993-shoulder	4052-next	4111-machine	4170-year
3994-exodus	4053-were	4112-department	4171-laughter
3995-happy	4054-moving	4113-stores	4172-tear
3996-ones	4055-there	4114-toilet	4173-just
3997-near	4056-colors	4115-paper	4174-today
3998-together	4057-street	4116-styrofoam	4175-maybe
3999-before	4058-white	4117-boxes	4176-tomorrow
4000-much	4059-blue	4118-osone	4177-good
4001-older	4060-people	4119-layer	4178-lord
4002-they	4061-shufflin	4120-says	4179-will
4003-wasted	4062-their	4121-hope	4180-take
4004-dancing	4063-feet	4122-alive	4181-away
4005-there	4064-sleepin	4123-fuel	4182-dream
4006-record	4065-shoes	4124-burn	4183-yourself
4007-machine	4066-warning	4125-roads	4184-until
4008-knew	4067-sign	4126-drive	4185-your
4009-must	4068-road	4127-every	4186-your
4010-have	4069-ahead	4128-time	4187-motor
4011-been	4070-sayin	4129-that	4188-runnin
4012-about	4071-better	4130-look	4189-head
4013-seventeen	4072-dead	4131-mirror	4190-highway
4014-beat	4073-feel	4132-these	4191-looking
4015-going	4074-like	4133-lines	4192-adventure
4016-strong	4075-satan	4134-face	4193-whatever
4017-playing	4076-them	4135-gettin	4194-comes
4018-favorite	4077-forget	4136-clearer	4195-yeah
4019-song	4078-keep	4137-past	4196-darlin
4020-could	4079-suckin	4138-gone	4197-gonna
4021-tell	4080-free	4139-went	4198-make
4022-wouldn	4081-world	4140-like	4199-happen

4494-round	5097-like	7067-read	9105-good
4495-outside	5098-have	7068-tune	9106-lord
4496-round	5099-well	7069-print	9107-will
4497-outside	5100-truthful	7070-scan	9108-will
4498-guess	5101-blueprints	7071-send	9108-take
4499-whos	5102-simply	7072-rename	9109-away
4500-back	5103-rage	7073-touch	9110-sing
4501-back	5104-youthful	7074-bring	9111-with
4502-again	5105-exuberance	7075-watch	9112-sing
4503-shadys	5106-everybody	7076-turn	9113-year
4504-back	5107-loves	7077-leave	9114-sing
4505-tell	5108-root	7078-start	9115-laughter
4506-friend	5109-nuisance	7079-format	9116-sing
4507-guess	5110-earth	7080-break	9117-tear
4508-whoss	5111-like	7081-trash	9118-sing
4509-back	5112-asteroid	7082-change	9119-with
4510-guess	5113-nothing	7083-mail	9120-just
4511-whos	5114-shoot	7084-upgrade	9121-today
4512-back	5115-moon	7085-charge	9122-maybe
4513-guess	5116-since	7086-point	9123-tomorrow
4514-whos	5117-taken	7087-zoom	9124-good
4515-back	5118-school	7088-press	9125-lord
4516-guess	5119-with	7089-snap	9126-will
4517-whos	5120-this	7090-work	9127-take
4518-back	5121-music	7091-quick	9128-away
4519-guess	5122-cause	7092-erase	9129-your
4520-whos	5123-vehicle	7093-write	9130-motor
4521-back	5124-rhyme	7094-paste	9131-runnin
4522-guess	5125-lead	7095-save	9132-head
4523-whos	5126-school	7096-load	9133-highway
4524-back	5127-full	7097-check	9134-locking
4525-guess	5128-students	7098-quit	9135-adventure
4526-whos	5129-product	7099-rewrite	9136-whatever
4527-back	5130-rakin	7100-plug	9137-comes
4528-ovulating	6184-lives	7101-nlav	9138-yeah
4529-know	6185-baby	8472-always	9218-once
4530-cheney	6186-crystal	8473-ritual	9219-explode
4531-your	6187-ball	8474-always	9220-into
4532-husbands	6188-baby	8475-ritual	9221-space
4533-heart	6189-does	8476-come	9222-like
4534-problem	6190-care	8477-come	9223-true
4535-complicati	6191-baby	8478-come	9224-natures
4536-wont	6192-having	8480-come	9225-child
4537-they	6193-much	8481-babe	9226-were
4538-tried	6194-shit	8482-come	9227-born
4539-shut	6195-kickin	8483-come	9228-born
4540-down	6196-mother	8484-come	9229-wild
4541-feels	6197-fuckin	8485-come	9230-climb
4542-empty	6198-atom	8486-come	9231-high
4543-without	6199-bomb	8487-babe	9232-never
4544-come	6200-baby	8488-love	9233-wanna
4600-your	6201-fleet	8489-love	9234-born
4601-lips	6202-submarine	8490-love	9235-wild
		8491-babe	9236-born
			9237-wild

Lista 1 9238 palabras/token's

✓ Stopwords

Con el método **separarpalabras** también se realiza el stopwords, ignora todas las palabras del texto de entrenamiento que sean menores a 4 caracteres

11-no	998-in	2369-I	3618-To	4470-be	5604-to
12-I	999-of	2370-I	3619-up	4471-is	5605-be
13-am	1000-tv	2371-am	3620-we	4472-I	5606-to
14-go	1001-in	2372-In	3621-So	4473-is	5607-be
15-to	1002-in	2373-I	3622-to	4474-in	5608-on
16-me	1003-no	2374-am	3623-is	4475-It	5609-In
17-To	1004-on	2375-I	3624-up	4476-is	5610-it
18-me	1005-in	2376-t	3625-to	4477-ll	5611-in
19-a	1006-I	2377-I	3626-I	4478-to	5612-a
20-a	1007-be	2378-am	3627-m	4479-I	5613-of
21-my	1008-of	2379-I	3628-up	4480-I	5614-at
22-he	1009-us	2380-m	3629-to	4481-to	5615-a
23-is	1010-I	2381-so	3630-is	4482-I	5616-we
24-I	1011-I	2382-of	3631-up	4483-I	5617-to
25-it	1012-so	2383-I	3632-I	4484-in	5618-be
26-to	1013-I	2384-I	3633-m	4485-my	5619-we
27-If	1014-I	2385-or	3634-up	4486-I	5620-so
28-I	1015-t	2386-my	3635-to	4487-it	5621-I
29-am	1016-up	2387-I	3636-We	4488-a	5622-to
30-on	1017-It	2388-m	3637-up	4489-it	5623-be
31-on	1018-is	2389-be	3638-to	4490-to	5624-to
32-As	1019-I	2390-to	3639-we	4491-it	5625-be

Lista 2 Stop words

Con la función `length ()` se obtiene la longitud del token, y un `if` condicional de que el token sea mayor a 3 caracteres, se agrega `Apalabras` a `palabraL_temp` con la función `add ()`.

```
public ArrayList<String> separarPalabras(String texto,String catego ){
    ArrayList<String> palabrasL_temp;
    String[] Apalabras=texto.split("[, ?.@ ' !]+"); //separa la palabra de la cade
    palabrasL_temp = new ArrayList();
    for(int i=0; i<Apalabras.length;i++){
        //ingresa la palabras a la lista de palabras
        if (Apalabras[i].length() > 3 ) { // El contains()método comprueba si una cad
            //verifica que no se repita la palabra
            palabrasL_temp.add(Apalabras[i]);
        }
    }
}
```

Fig. 3 Método separarpalabras

✓ Bag of Words

Utiliza un vector de palabras en el cual se lleva un recuento de cada una de ella y su posición en el vector; se clasifica con un switch la categoría del token´s/palabras dentro de la matriz `cpcate` (contador de palabras por categoría).

```
54 switch (catego.toUpperCase()) { //toUpperCase()método co
55
56     case "QUEEN":
57         cpcate.ingresarPalabra(Apalabras[i], 0 );
58
59         count[0]++;
60
61         break;
62     case "EMINEM":
63         cpcate.ingresarPalabra(Apalabras[i], 1);
64         count[1]++;
65         break;
66     case "DAFT PUNK":
67         cpcate.ingresarPalabra(Apalabras[i], 2);
68         count[2]++;
69         break;
70     case "TIESTO":
71         cpcate.ingresarPalabra(Apalabras[i], 3);
72         count[3]++;
73         break;
74     case "OTHER":
75         cpcate.ingresarPalabra(Apalabras[i], 4);
76         count[4]++;
77         break;
78     default:
```

Fig. 4 El switch convierte la categoría a mayúsculas

Esto se realiza en la clase **MatrizPalabraCantidad** con el ArreyList `palabra` dentro del método **ingresarPalabra**, con un if y la función **contains** verifica que la palabra/token no esté dentro de la cadena `palabr`, si es verdadera la agrega a la cadena `palabr` con la función **add** indicándole la posición o índice para agregarla, cuando la palabra ya se encuentra en la cadena `palabr` entonces ubica su posición con la variable `dato` y la función **indexOf** que devuelve la posición del token para contar la frecuencia de ocurrencia de cada palabra (token) por categoría, se utiliza para realizar la tabla de frecuencias.

```
Analisis - NetBeans IDE 8.2
Edit View Navigate Source Refactor Run Debug Profile Team Tools Window Help
- default config -
StartPage x Analisis.java x archivos.java x MetodoPalabraCantidad.java x Palabras.java x
Source History
34 public void ingresarPalabra(String palabra, int categoria) { // agrega las palabras según categorías dadas
35
36     if (!palabra.contains(" ")) { //El usuario escribió correctamente y, una vez más, valida una palabra
37         palabra.add(indice, palabra); // una vez más valida un elemento al final de la lista y se le da
38         switch (categoria) {
39             case 0:
40                 queen[indice]**; // el índice para de acuerdo al género.
41                 // System.out.println("palabra: "+palabra+" "+ queen[indice]);
42                 break;
43             case 1:
44                 eminem[indice]**;
45                 // System.out.println("palabra: "+palabra+" "+ eminem[indice]);
46                 break;
47             case 2:
48                 daftpunk[indice]**;
49                 break;
50             case 3:
51                 tiesto[indice]**;
52                 break;
53             case 4:
54                 other[indice]**;
55                 break;
56             default:
57                 break;
58         }
59         indice++;
60
61     } else {
62         int dato=palabra.indexOf(" "); // devuelve la posición, mediante un número entero
63         switch (categoria) {
64             case 0:
65                 queen[dato]**;
66                 // System.out.println("palabra: "+palabra+" encontrada: "+ dato);
67                 break;
68             case 1:
69                 eminem[dato]**;
70                 break;
71             case 2:
72                 daftpunk[dato]**;
73                 break;
74             case 3:
75                 tiesto[dato]**;
76                 break;
77             case 4:
78                 other[dato]**;
79                 break;
80             default:
81                 break;
82         }
83
84     }
85 }
86 }
```

Fig. 5 Método ingresarPalabra separa la cadena en token's y asigna un índice a cada una.

c) Análisis de los Datos

✓ Tabla de frecuencias

Con el método **ver_palabraCateCantidad** (ver palabra por categoría y cantidad) se imprimen los datos del método **ingresarPalabra** para observar la frecuencia de las palabras en cada categoría.

```
*****Tabla de frecuencias *****
0-Palabra- this
  Queen 20.0  Eminem 79.0  Daft Punk 8.0  Tiesto 16.0  Other 0.0
1-Palabra- real
  Queen 3.0  Eminem 24.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
2-Palabra- life
  Queen 5.0  Eminem 6.0  Daft Punk 11.0  Tiesto 5.0  Other 2.0
3-Palabra- just
  Queen 16.0  Eminem 91.0  Daft Punk 14.0  Tiesto 17.0  Other 4.0
4-Palabra- fantasy
  Queen 1.0  Eminem 0.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
5-Palabra- caught
  Queen 1.0  Eminem 3.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
6-Palabra- landslide
  Queen 1.0  Eminem 0.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
7-Palabra- escape
  Queen 1.0  Eminem 1.0  Daft Punk 0.0  Tiesto 12.0  Other 0.0
8-Palabra- from
  Queen 7.0  Eminem 28.0  Daft Punk 1.0  Tiesto 5.0  Other 2.0
9-Palabra- reality
  Queen 1.0  Eminem 1.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
10-Palabra- open
  Queen 1.0  Eminem 3.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
11-Palabra- your
  Queen 5.0  Eminem 37.0  Daft Punk 11.0  Tiesto 18.0  Other 14.0
12-Palabra- eyes
  Queen 2.0  Eminem 0.0  Daft Punk 0.0  Tiesto 2.0  Other 0.0
13-Palabra- look
  Queen 4.0  Eminem 8.0  Daft Punk 1.0  Tiesto 0.0  Other 2.0
14-Palabra- skies
  Queen 1.0  Eminem 0.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
15-Palabra- poor
  Queen 4.0  Eminem 1.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
```

519-Palabra- asap	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
520-Palabra- mommas	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
521-Palabra- sister	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
522-Palabra- dont	Queen 0.0	Eminem 38.0	Daft Punk 1.0	Tiesto 3.0	Other 0.0
523-Palabra- understand	Queen 0.0	Eminem 2.0	Daft Punk 1.0	Tiesto 0.0	Other 0.0
524-Palabra- site	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
525-Palabra- front	Queen 0.0	Eminem 6.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
526-Palabra- buries	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
527-Palabra- nose	Queen 0.0	Eminem 3.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
528-Palabra- colors	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
529-Palabra- crayon	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
530-Palabra- dull	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
531-Palabra- hand	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 1.0	Other 3.0
532-Palabra- while	Queen 0.0	Eminem 8.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
533-Palabra- brother	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
534-Palabra- mother	Queen 0.0	Eminem 2.0	Daft Punk 2.0	Tiesto 0.0	Other 0.0
755-Palabra- wreck	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
756-Palabra- cameras	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
757-Palabra- directions	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
758-Palabra- press	Queen 0.0	Eminem 2.0	Daft Punk 7.0	Tiesto 0.0	Other 0.0
759-Palabra- apeshit	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
760-Palabra- bananas	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
761-Palabra- networks	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
762-Palabra- commando	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
763-Palabra- extra	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0

058-Palabra- reports	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
059-Palabra- number	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
060-Palabra- fatalities	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
061-Palabra- santa	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
062-Palabra- texas	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
063-Palabra- galveston	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
064-Palabra- mass	Queen 0.0	Eminem 3.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
065-Palabra- southern	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
066-Palabra- california	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
067-Palabra- schoolve	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
068-Palabra- following	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
069-Palabra- deadly	Queen 0.0	Eminem 3.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
070-Palabra- houston	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
071-Palabra- twenty-six	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
072-Palabra- twenty	Queen 0.0	Eminem 2.0	Daft Punk 1.0	Tiesto 0.0	Other 0.0
073-Palabra- other	Queen 0.0	Eminem 11.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
992-Palabra- roof	Queen 0.0	Eminem 11.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
993-Palabra- dogs	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
994-Palabra- caged	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
995-Palabra- playin	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
996-Palabra- beginnin	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
997-Palabra- mood	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
998-Palabra- changed	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
999-Palabra- chewed	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1000-Palabra- boood	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1001-Palabra- kept	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1002-Palabra- rhywin	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1003-Palabra- stepwritin	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0

1129-Palabra- squeezing	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1129-Palabra- neck	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1130-Palabra- popsicle	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1131-Palabra- stick	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1132-Palabra- possible	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1133-Palabra- house	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1134-Palabra- streamed	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1135-Palabra- cheeks	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1136-Palabra- temple	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1137-Palabra- wouldve	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1138-Palabra- anything	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 2.0	Other 0.0
1139-Palabra- adored	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1140-Palabra- save	Queen 0.0	Eminem 1.0	Daft Punk 7.0	Tiesto 0.0	Other 0.0
1141-Palabra- Just	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1142-Palabra- youll	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1143-Palabra- attention	Queen 0.0	Eminem 4.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1144-Palabra- Will	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1416-Palabra- spilling	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1417-Palabra- killing	Queen 0.0	Eminem 3.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1418-Palabra- pinning	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1419-Palabra- hiphop	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1420-Palabra- subliminal	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1421-Palabra- thoughts	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1422-Palabra- sending	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1423-Palabra- webs	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1424-Palabra- spin	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1425-Palabra- hock	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1426-Palabra- venom	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0

1563-Palabra- disaster	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1564-Palabra- catastrophe	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1565-Palabra- nana	Queen 0.0	Eminem 2.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1566-Palabra- kshh	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1567-Palabra- bent	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1568-Palabra- antenna	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1569-Palabra- tune	Queen 0.0	Eminem 1.0	Daft Punk 6.0	Tiesto 0.0	Other 0.0
1570-Palabra- enter	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1571-Palabra- endin	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1572-Palabra- splinter	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1573-Palabra- center	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1574-Palabra- winter	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1575-Palabra- interesting	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1576-Palabra- wrestling	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1577-Palabra- infesting	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1578-Palabra- ears	Queen 0.0	Eminem 1.0	Daft Punk 0.0	Tiesto 0.0	Other 0.0
1567-Palabra- celebrate	Queen 0.0	Eminem 0.0	Daft Punk 72.0	Tiesto 0.0	Other 0.
1558-Palabra- dancing	Queen 0.0	Eminem 0.0	Daft Punk 16.0	Tiesto 0.0	Other 1
1559-Palabra- uhmm	Queen 0.0	Eminem 0.0	Daft Punk 4.0	Tiesto 0.0	Other 0.1
1560-Palabra- celebration	Queen 0.0	Eminem 0.0	Daft Punk 4.0	Tiesto 0.0	Other 0.1
1561-Palabra- musics	Queen 0.0	Eminem 0.0	Daft Punk 30.0	Tiesto 0.0	Other 0
1562-Palabra- often	Queen 0.0	Eminem 0.0	Daft Punk 3.0	Tiesto 0.0	Other 0.1
1563-Palabra- speeding	Queen 0.0	Eminem 0.0	Daft Punk 3.0	Tiesto 0.0	Other 0.1
1564-Palabra- shirt	Queen 0.0	Eminem 0.0	Daft Punk 3.0	Tiesto 0.0	Other 0.1
1565-Palabra- ahead	Queen 0.0	Eminem 0.0	Daft Punk 3.0	Tiesto 0.0	Other 1.1
1566-Palabra- wipe	Queen 0.0	Eminem 0.0	Daft Punk 3.0	Tiesto 0.0	Other 0.1
1567-Palabra- dances	Queen 0.0	Eminem 0.0	Daft Punk 4.0	Tiesto 0.0	Other 0.1
1568-Palabra- everybodyon	Queen 0.0	Eminem 0.0	Daft Punk 1.0	Tiesto 0.0	Other 0.1

2109-Palabra- fields	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2110-Palabra- meals	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2111-Palabra- forgiven	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2112-Palabra- teenage	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 5.0
2113-Palabra- wasteland	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 5.0
2114-Palabra- sally	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2115-Palabra- south	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2116-Palabra- past	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 2.0
2117-Palabra- shoulder	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2118-Palabra- exodus	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2119-Palabra- ones	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2120-Palabra- older	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2121-Palabra- record	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2122-Palabra- seventeen	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2123-Palabra- wouldn	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2124-Palabra- lord	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 4.0
2125-Palabra- highway	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 2.0
2126-Palabra- adventure	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 2.0
2127-Palabra- darlin	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 3.0
2128-Palabra- embrace	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 3.0
2129-Palabra- lightning	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2130-Palabra- metal	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2131-Palabra- thunder	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 1.0
2132-Palabra- matures	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 2.0
2133-Palabra- child	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 2.0
2134-Palabra- climb	Queen 0.0	Eminem 0.0	Daft Punk 0.0	Tiesto 0.0	Other 2.0
.....					
total de palabras: 9238					

Fig. 6 Tabla de ocurrencia de las 9238 palabras del texto de entrenamiento

Para el conteo de textos de entrenamiento se convierte la categoría a mayúsculas; con un if y la función **equals** compara que la categoría no este vacia para poder obtener la cantidad de textos incrementando la variable **cantidadTextos**, tambien se lleva un conteo por categoría del texto de entrenamiento y se guarda e incrementa en el vector **categoria** para poder calcular las probabilidades a priori de cada clase.

```

90         switch (catego.toUpperCase()) { //con
91                                     //c
92             case "QUEEN":
93                 categoria[0]++;
94
95                 break;
96             case "EMINEM":
97                 categoria[1]++;
98                 break;
99             case "DAFT PUNK":
100                categoria[2]++;
101                break;
102             case "TIESTO":
103                 categoria[3]++;
104                 break;
105             case "OTHER":
106                 //revisar el toupp
107                 categoria[4]++;
108                 break;

```

Fig. 7 Conteo de textos por categoría.

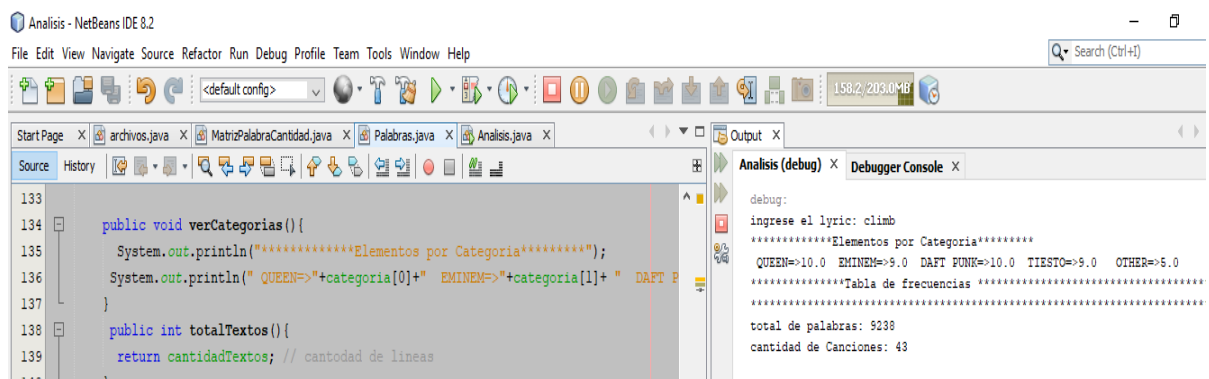


Fig. 8 En VerCategorias se imprime la cantidad de textos que tiene cada categoría, en totalTextos nos da el total de textos de entrenamiento.

d) Entrenamiento

✓ Crear Modelo

✓ Probabilidades a priori de cada clase/categoría

Una vez obtenido el total de textos de entrenamiento y categorizado se calcula la probabilidad a priori de cada clase.

$$P(y_i) = \frac{\text{textos de categoría}}{\text{total de textos}} \quad (8)$$

En el vector proCategoría(probabilidadCategoría) almacena el resultado de la división del valor categoría [i] entre cantidadtexto de acuerdo con la fórmula (8) y se realiza en el método **CalProbabilidadCategorías** teniendo un total de 1 entre el total de clases, los resultados se observan en la fig8.

```

145 public double getProCategorias(int pos){
146     return proCategorias[pos];
147 }
148
149 public void calProbabilidadCategorias(){
150
151     System.out.println("cantidad de Canciones: "+cantidadTextos);
152
153     System.out.println("*****Probabilidad de Categoría *****");
154     for(int j=0;j<4;j++){
155         proCategorias[j]=categoria[j]/cantidadTextos;
156         System.out.println("Categoría "+j+ " : " + proCategorias[j]);
157         // suma proCategorias[j] * prob;
158     }
159 }
160 System.out.println("*****");
161 //System.out.println("la Probabilidad de categoría \"+categoria[i]+"\": "+proCategorias[i]);
162 }
163
164 public void verPalabraCantidadCate(){
165     update_var_palabraCateCantidad();
166 }
167
168 //clasificando
169 public String clasificar(String texto){
170     double prob_0=0;
171     double prob_1=0;
172     double prob_2=0;
173     double prob_3=0;
174     double prob_4=0;
175 }
    
```

```

Ingreso el lyric: this is love
*****Elementos por Categoría*****
QUEEN=>10.0 EMINEM=>9.0 DAFT PUNK=>10.0 TIESTO=>9.0 OTHER=>5.0
*****Tabla de frecuencias *****
*****
total de palabras: 9238
cantidad de Canciones: 43

*****Probabilidad de Categoría *****
categoría0: 0.23255813953488372
categoría1: 0.20930232558139536
categoría2: 0.23255813953488372
categoría3: 0.20930232558139536
categoría4: 0.11627906976744186
*****
Fr 9238 ad de la categoría/clase --10.0/43=0.23255813953488372
0 palabras: this
Categoría 0: 10.0
Categoría 1: 9.0
Categoría 2: 8.0
Categoría 3: 14.0
Categoría 4: 0.0
probabilidad de la palabra= 0.0055459217977095
prob_total_0 = 0.001299207483431044
prob_total_1 = 0.001299207483431044
calcular la categoría mejor: 0: 0.001299207483431044
94 palabras: love
Categoría 0: 72.0
Categoría 1: 7.0
Categoría 2: 9.0
Categoría 3: 15.0
Categoría 4: 18.0
    
```

Fig. 9. Calcular probabilidad por categoría.

```

ingreso el lyric: this is love
*****Elementos por Categoría*****
QUEEN=>10.0 EMINEM=>9.0 DAFT PUNK=>10.0 TIESTO=>9.0 OTHER=>5.0
*****Tabla de frecuencias *****
*****
total de palabras: 9238
cantidad de Canciones: 43

*****Probabilidad de Categoría *****
categoría0: 0.23255813953488372
categoría1: 0.20930232558139536
categoría2: 0.23255813953488372
categoría3: 0.20930232558139536
categoría4: 0.11627906976744186
*****
    
```

Fig. 10 Probabilidad a priori de las clases

✓ **Suavizado de Laplace**

Para obtener la probabilidad de $P(X|y_j)$ que es la probabilidad de que la categoría y_j teniendo en cuenta la cantidad de ocurrencias de la palabra x_i entre todas las palabras de la categoría y_j , luego con esa frecuencia se divide entre el total de palabras del texto de entrenamiento.

$$P(x_i|y) = \frac{\text{count}(x_i|y_j)}{|V|} \quad (9)$$

Un problema de la clasificación de textos es cuando una palabra no está en todas las categorías del texto de entrenamiento, la probabilidad sería cero y nos daría un error en la clasificación, para solucionarlo se utiliza el suavizado de Laplace que consiste en incrementar en una unidad la ocurrencia de las palabras en las categorías como se muestra en la fórmula (9).

$$\text{count}(x_i|y_j) + 1 \quad | \quad (10)$$

Por ejemplo, con la palabra `climb` que solo está en la categoría `other` solo estará presente en esa categoría como se muestra en la Fig11, cuando es solo una palabra no hay un mayor problema, pero cuando es una oración la probabilidad cambiaría a cero como se observa en la Fig12.

```
2168-Palabra- climb
Queen 0.0 Eminem 0.0 Daft Punk 0.0 Tiesto 0.0 Other 2.
*****
```

```

*****Categoria 3*****
Probabilidad de la categoria/clase -->9.0/43=0.20930232558139536
2168 Palabra: climb
Categoria 0: 0.0
Categoria 1: 0.0
Categoria 2: 0.0
Categoria 3: 0.0
Categoria 4: 2.0
probabilidad de la palabra= 0.0
prob_total_c anterior0.20930232558139536
prob_total_c =0.0
Rango mayor : 0.0***0.0 valor de i: 3 categoria final: 0
*****Categoria 4*****
Probabilidad de la categoria/clase -->5.0/43=0.11627906976744186
2168 Palabra: climb
Categoria 0: 0.0
Categoria 1: 0.0
Categoria 2: 0.0
Categoria 3: 0.0
Categoria 4: 2.0
probabilidad de la palabra= 7.176175098672408E-4
prob_total_c anterior0.11627906976744186
prob_total_c =8.34438964961908E-5
calcular la categoria mayor 4: 8.34438964961908E-5
Rango mayor : 8.34438964961908E-5***8.34438964961908E-5 valor de i: 4
*****Resultado*****
palabraslocal [climb]
rango de probabilidad: 8.34438964961908E-5
categoria final: 4

```

Fig. 11. Probabilidad de la palabra climb sin suavizado

```

*****Categoria 3*****
Probabilidad de la categoria/clase -->9.0/43=0.20930232558139536
67 Palabra: never
Categoria 0: 9.0
Categoria 1: 17.0
Categoria 2: 30.0
Categoria 3: 3.0
Categoria 4: 5.0
probabilidad de la palabra= 9.372071227741331E-4
prob_total_c anterior0.20930232558139536
prob_total_c =1.961596303480744E-4
2168 Palabra: climb
Categoria 0: 0.0
Categoria 1: 0.0
Categoria 2: 0.0
Categoria 3: 0.0
Categoria 4: 2.0
probabilidad de la palabra= 0.0
prob_total_c anterior1.961596303480744E-4
prob_total_c =0.0
Rango mayor : 0.0***0.0 valor de i: 3 categoria final: 2

```

Fig. 12 Probabilidad de las palabras never climb sin suavizado

Para solucionar el error se incrementa una unidad en cada palabra en la Fig11 se observa que se incrementa en uno la ocurrencia del token en cada categoría y así no se pierde el valor de la probabilidad anterior.

```
*****Categoria 4*****
Probabilidad de la categoria/clase -->5.0/43=0.11627906976744186
2129 Palabra: blue
Categoria 0: 1.0
Categoria 1: 1.0
Categoria 2: 1.0
Categoria 3: 1.0
Categoria 4: 2.0
probabilidad de la palabra= 7.176175098672408E-4
prob_total_c anterior0.11627906976744186
prob_total_c =8.34438964961908E-5
calcular la categoria mayor 4: 8.34438964961908E-5
2168 Palabra: climb
Categoria 0: 1.0
Categoria 1: 1.0
Categoria 2: 1.0
Categoria 3: 1.0
Categoria 4: 3.0
probabilidad de la palabra= 0.001076426264800861
prob_total_c anterior8.34438964961908E-5
prob_total_c =8.982120182582431E-8
Rango mayor : 8.982120182582431E-8****8.982120182582431E-8 valor de i: 4 categ
*****Resultado*****
palabraslocal [blue, climb]
rango de probabilidad: 8.982120182582431E-8
categoria final: 4
```

Fig. 13. Incrementa en uno las palabras blue y climb con suavizado.

✓ **Normalizar los datos**

Entonces para conocer el valor $P(x_i|y_j)$, se normalizan los valores de acuerdo con a la formula (11), donde se realiza el suavizado de Laplace entre la suma del total de palabras de la categoría más el vocabulario total.

$$P(x_i|y) = \frac{\text{count}(x_i|y_j) + 1}{\text{count}(y_j) + |V|} \quad (11)$$

✓ **Palabras desconocidas**

Para las palabras que se ingresan en los datos, pero no está en el texto de entrenamiento la solución es ignorarlas para no incluir ninguna probabilidad.

Por ejemplo, al ingresar la oración:” this rain is real” la palabra que no está en el entrenamiento es “rain” y “is” es un stop words. En la Fig14. se observa el total de elementos por categoría y la tabla de frecuencias de las palabras this y real que ocurren en cada categoría.

```
*****Elementos por Categoría*****
QUEEN=>10.0  EMINEM=>9.0  DAFT PUNK=>10.0  TIESTO=>9.0  OTHER=>5.0
*****Tabla de frecuencias *****
0-Palabra- this
Queen 20.0  Eminem 79.0  Daft Punk 8.0  Tiesto 16.0  Other 0.0
1-Palabra- real
Queen 3.0  Eminem 24.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
```

Fig. 14. Ocurrencia de this y real por categoría.

En la Fig15. muestra la probabilidad a priori de las clases/categoría y se muestra el suavizado de las palabras this y rain se ignora con un valor de -1.

Al no estar en el entrenamiento la palabra “rain” no realiza operación y continua.

```
.....
total de palabras: 9238
cantidad de Canciones: 43

*****Probablilidad de Categoría *****
categoria0: 0.23255813953488372
categoria1: 0.20930232558139536
categoria2: 0.23255813953488372
categoria3: 0.20930232558139536
categoria4: 0.11627906976744186
*total*****
*****Categoría 0*****
Probabilidad de la categoría/clase -->10.0/43=0.23255813953488372
0 Palabra: this
Categoría 0: 21.0
Categoría 1: 80.0
Categoría 2: 9.0
Categoría 3: 17.0
Categoría 4: 1.0
probabilidad de la palabra= 0.005865921787709497
prob_total_c anterior0.23255813953488372
prob_total_c =0.0013641678576068598
calcular la categoría mayor 0: 0.0013641678576068598
-1 Palabra: rain
Categoría 0: -1.0
Categoría 1: -1.0
Categoría 2: -1.0
Categoría 3: -1.0
Categoría 4: -1.0
1 Palabra: real
Categoría 0: 4.0
Categoría 1: 25.0
Categoría 2: 1.0
Categoría 3: 1.0
Categoría 4: 1.0
probabilidad de la palabra= 0.0011173184357541898
prob_total_c anterior0.0013641678576068598
prob_total_c =1.524209896767441E-6
Rango mayor : 1.524209896767441E-6****1.524209896767441E-6 valor de i: 0 categoría final: 0
.....
```

Fig. 15 Realiza el suavizado de Laplace en los token´s this y real.

Al no estar la palabra rain en el texto de entrenamiento le agrega un valor de -1 en cada categoría y no realiza operaciones y pasa a la siguiente palabra.

✓ Clasificar

✓ Clasificar una oración nueva.

Para cada categoría disponible se determinan los valores de probabilidad de cada palabra/token del nuevo ejemplo.

Siguiendo con el ejemplo anterior de la Fig15; en la Fig16. se muestra la probabilidad de las palabras por categoría de acuerdo con la formula (11).

```
*****Categoría 0*****
Probabilidad de la categoría/clase -->10.0/43=0.23255813953488372
0 Palabra: this
Categoría 0: 21.0
Categoría 1: 80.0
Categoría 2: 9.0
Categoría 3: 17.0
Categoría 4: 1.0
probabilidad de la palabra= 0.005865921787709497

*****Categoría 1*****
Probabilidad de la categoría/clase -->9.0/43=0.20930232558139536
0 Palabra: this
Categoría 0: 21.0
Categoría 1: 80.0
Categoría 2: 9.0
Categoría 3: 17.0
Categoría 4: 1.0
probabilidad de la palabra= 0.012447487163528862

*****Categoría 2*****
Probabilidad de la categoría/clase -->10.0/43=0.23255813953488372
0 Palabra: this
Categoría 0: 21.0
Categoría 1: 80.0
Categoría 2: 9.0
Categoría 3: 17.0
Categoría 4: 1.0
probabilidad de la palabra= 0.002204261572373255

*****Categoría 3*****
Probabilidad de la categoría/clase -->9.0/43=0.20930232558139536
0 Palabra: this
Categoría 0: 21.0
Categoría 1: 80.0
Categoría 2: 9.0
Categoría 3: 17.0
Categoría 4: 1.0
probabilidad de la palabra= 0.005310840362386754

*****Categoría 4*****
Probabilidad de la categoría/clase -->5.0/43=0.11627906976744186
0 Palabra: this
Categoría 0: 21.0
Categoría 1: 80.0
Categoría 2: 9.0
Categoría 3: 17.0
Categoría 4: 1.0
probabilidad de la palabra= 3.588087549336204E-4
```

Fig. 16 Probabilidades del token this en las cinco categorías.

✓ Naive Bayes

Aplicar la formula (7) del capítulo 3 Independencia Naive Bayes Multinomial para obtener la probabilidad $P(Y_j|X)$, que es la probabilidad a posteriori del token para determinar la clasificación de la categoría Y_j a la que pertenece.

$$P(x_1, x_2, \dots, x_n | y)$$

$$P(x_1, x_2, \dots, x_n | y) = P(x_1 | y) * P(x_2 | y) * P(x_3 | y) * \dots * P(x_n | y)$$

$$Y_{MAP} = \operatorname{argmax}_{y \in Y} P(x_1, x_2, \dots, x_n | y) P(y)$$

$$Y_{NB} = \operatorname{argmax}_{y \in Y} P(y) \prod_i^n P(x_i | y_j) \quad (7)$$

Siguiendo con el ejemplo, ahora después de conocer la probabilidad del primer token se aplicará la formula (7) que nos dice: multiplica la probabilidad a priori de la clase por la probabilidad del token i por categoría, hasta n -token para estimar la hipótesis más probable o MAP (máximum a posteriori hypothesis).

Esto se realiza dentro de la clase **Palabras** en el método **clasificar**, realiza las operaciones dentro de dos estructuras **for**; dentro del primer **for** se obtiene las palabras locales que es la frase ingresada al inicio, en cada categoría, el ArrayList **palabraslocal**, llama al método **separar palabras** y convierte en minúsculas las palabras de la frase ingresada.

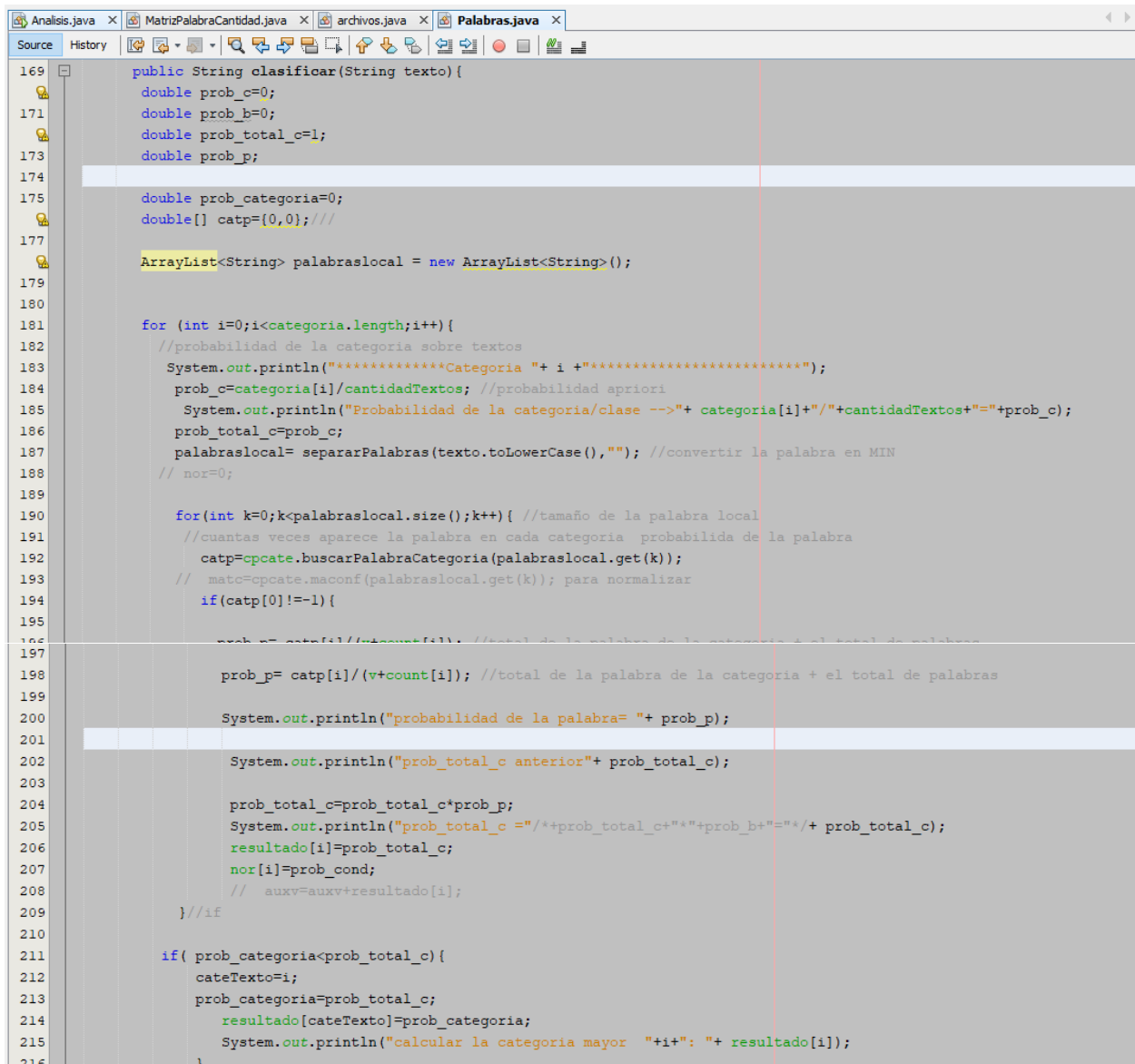
Dentro del segundo **for** se realizan las operaciones de la formula (7)

Primero se debe de conocer la cantidad de ocurrencia de cada token de palabraslocal,

y se realiza dentro el método **buscarPalabraCategoria** de la clase

MatrizPalabraCantidad en base a la matriz **cpcate** que es donde se guarda la

frecuencia del token y lo almacena dentro de arreglo **catp** (cantidad de palabras).



```
169 public String clasificar(String texto){
170     double prob_c=0;
171     double prob_b=0;
172     double prob_total_c=1;
173     double prob_p;
174
175     double prob_categoria=0;
176     double[] catp={0,0};
177
178     ArrayList<String> palabraslocal = new ArrayList<String>();
179
180
181     for (int i=0;i<categoria.length;i++){
182         //probabilidad de la categoria sobre textos
183         System.out.println("*****Categoria " + i + "*****");
184         prob_c=categoria[i]/cantidadTextos; //probabilidad apriori
185         System.out.println("Probabilidad de la categoria/clase -->" + categoria[i]+"/"+cantidadTextos+"="+prob_c);
186         prob_total_c=prob_c;
187         palabraslocal= separarPalabras(texto.toLowerCase(),""); //convertir la palabra en MIN
188         // nor=0;
189
190         for(int k=0;k<palabraslocal.size();k++){ //tamaño de la palabra local
191             //cuantas veces aparece la palabra en cada categoria probabilidad de la palabra
192             catp=cpcate.buscarPalabraCategoria(palabraslocal.get(k));
193             // matc=cpcate.maconf(palabraslocal.get(k)); para normalizar
194             if(catp[0]!=-1){
195
196                 prob_p=catp[i]/(v+count[i]); //total de la palabra de la categoria + el total de palabras
197
198                 prob_p= catp[i]/(v+count[i]); //total de la palabra de la categoria + el total de palabras
199
200                 System.out.println("probabilidad de la palabra= " + prob_p);
201
202                 System.out.println("prob_total_c anterior"+ prob_total_c);
203
204                 prob_total_c=prob_total_c*prob_p;
205                 System.out.println("prob_total_c ="+prob_total_c+"*"+prob_b+"="+prob_total_c);
206                 resultado[i]=prob_total_c;
207                 nor[i]=prob_cond;
208                 // auxv=auxv+resultado[i];
209             }//if
210
211         if( prob_categoria<prob_total_c){
212             cateTexto=i;
213             prob_categoria=prob_total_c;
214             resultado[cateTexto]=prob_categoria;
215             System.out.println("calcular la categoria mayor "+i+" : "+ resultado[i]);
216         }
```

Fig. 17 Método clasificar

Siguiendo con los pasos de la formula se obtiene la probabilidad a priori de la categoría `prob_c`, en seguida la probabilidad de la palabra `prob_p` para cada palabra de la frase, se multiplican ambas guardando el resultado en `prob_total_c`; para poder comparar los resultados y obtener la categoría que maximiza la probabilidad (argmax) se guarda el resultado de cada categoría en el arreglo `resultado` para finalizar dentro de un if compara la probabilidad total anterior con la actual para obtener la mayor. Como resultados se imprime las palabras que se consideraron separadas por comas, el resultado de la probabilidad mayor y la categoría.

```

242 System.out.println("*****Resultado*****");
243 System.out.println("palabraslocal "+palabraslocal);
244 System.out.println("rango de probabilidad: "+resultado[cateTexto)+"\n categoria final: "+cateTexto);
245 // System.out.println("normalizar: "+nor[cateTexto)+"\n categoria final: "+cateTexto);
246 System.out.println("*****");

```

Fig. 18: Código para la impresión de los resultados

✓ **Continuando con la clasificar de la oración anterior aplicando NBM.**

Probabilidad de las palabras this y real por categoría de acuerdo con la formula (11) para el token rain al no estar en el entrenamiento no se toma en cuenta.

$$P(x_i|y) = \frac{\text{count}(x_i|y_j) + 1}{\text{count}(y_j) + |V|} \quad (11)$$

$$P(\text{this} | \text{categoria0}) = \frac{20 + 1}{1412 + 2168} = 0.00586592178$$

$$P(\text{real} | \text{categoria0}) = \frac{3 + 1}{1412 + 2168} = 0.00111731843$$

$$P(\text{this real} | \text{categoria0}) = P(\text{this} | \text{categoria0}) * P(\text{real} | \text{categoria0})$$

$$P(\text{categoria0} | \text{this real}) = P(\text{this} | \text{categoria0}) * P(\text{real} | \text{categoria0}) * P(\text{categoria0})$$

$$P(\text{categoria0} | \text{this real}) = \frac{21}{3580} * \frac{4}{3580} * \frac{10}{43} = 0.0000015242$$

Ahora hay que calcular la probabilidad de las demás categorías usando el mismo procedimiento:

$$P(\text{categoria1} | \text{this real}) = \frac{80}{6427} * \frac{25}{6427} * \frac{9}{43} = 0.00001013415$$

$$P(\text{categoria2} | \text{this real}) = \frac{9}{4083} * \frac{1}{4083} * \frac{10}{43} = 0.00000012554$$

$$P(\text{categoria3} | \text{this real}) = \frac{17}{3201} * \frac{1}{3201} * \frac{9}{43} = 0.00000034725$$

$$P(\text{categoria4} | \text{this real}) = \frac{1}{2787} * \frac{1}{2787} * \frac{5}{43} = 0.00000001497$$

Conociendo la probabilidad a posteriori de cada categoría, la probabilidad más alta es el resultado de la clasificación. En este ejemplo la categoría mayor es **categoria1**, lo verificaremos en las Fig18 a Fig22.

```
*****Categoria 0*****
Probabilidad de la categoria/clase -->10.0/43=0.23255813953488372
0 Palabra: this
Categoria 0: 21.0
Categoria 1: 80.0
Categoria 2: 9.0
Categoria 3: 17.0
Categoria 4: 1.0
probabilidad de la palabra= 0.005865921787709497
prob_total_c anterior0.23255813953488372
prob_total_c =0.0013641678576068598
calcular la categoria mayor 0: 0.0013641678576068598
-1 Palabra: rain
Categoria 0: -1.0
Categoria 1: -1.0
Categoria 2: -1.0
Categoria 3: -1.0
Categoria 4: -1.0
1 Palabra: real
Categoria 0: 4.0
Categoria 1: 25.0
Categoria 2: 1.0
Categoria 3: 1.0
Categoria 4: 1.0
probabilidad de la palabra= 0.0011173184357541898
prob_total_c anterior0.0013641678576068598
prob_total_c =1.524209896767441E-6
-----
```

Fig. 19. Probabilidad de la oración ingresada para la categoría 0.

```
*****Categoria 1*****
Probabilidad de la categoria/clase -->9.0/43=0.20930232558139536
0 Palabra: this
Categoria 0: 21.0
Categoria 1: 80.0
Categoria 2: 9.0
Categoria 3: 17.0
Categoria 4: 1.0
probabilidad de la palabra= 0.012447487163528862
prob_total_c anterior0.20930232558139536
prob_total_c =0.0026052880109711575
calcular la categoria mayor 1: 0.0026052880109711575
-1 Palabra: rain
Categoria 0: -1.0
Categoria 1: -1.0
Categoria 2: -1.0
Categoria 3: -1.0
Categoria 4: -1.0
1 Palabra: real
Categoria 0: 4.0
Categoria 1: 25.0
Categoria 2: 1.0
Categoria 3: 1.0
Categoria 4: 1.0
probabilidad de la palabra= 0.0038898397386027696
prob_total_c anterior0.0026052880109711575
prob_total_c =1.0134152835580976E-5
```

Fig. 20 Probabilidad de la oración ingresada para la categoría 1.

```

*****Categoria 2*****
Probabilidad de la categoria/clase -->10.0/43=0.23255813953488372
0 Palabra: this
Categoria 0: 21.0
Categoria 1: 80.0
Categoria 2: 9.0
Categoria 3: 17.0
Categoria 4: 1.0
probabilidad de la palabra= 0.002204261572373255
prob_total_c anterior0.23255813953488372
prob_total_c =5.126189703193616E-4
-1 Palabra: rain
Categoria 0: -1.0
Categoria 1: -1.0
Categoria 2: -1.0
Categoria 3: -1.0
Categoria 4: -1.0
1 Palabra: real
Categoria 0: 4.0
Categoria 1: 25.0
Categoria 2: 1.0
Categoria 3: 1.0
Categoria 4: 1.0
probabilidad de la palabra= 2.449179524859172E-4
prob_total_c anterior5.126189703193616E-4
prob_total_c =1.2554958861605722E-7

```

Fig. 21 Probabilidad de la oración ingresada para la categoría 2.

```

*****Categoria 3*****
Probabilidad de la categoria/clase -->9.0/43=0.20930232558139536
0 Palabra: this
Categoria 0: 21.0
Categoria 1: 80.0
Categoria 2: 9.0
Categoria 3: 17.0
Categoria 4: 1.0
probabilidad de la palabra= 0.005310840362386754
prob_total_c anterior0.20930232558139536
prob_total_c =0.001111571238639088
-1 Palabra: rain
Categoria 0: -1.0
Categoria 1: -1.0
Categoria 2: -1.0
Categoria 3: -1.0
Categoria 4: -1.0
1 Palabra: real
Categoria 0: 4.0
Categoria 1: 25.0
Categoria 2: 1.0
Categoria 3: 1.0
Categoria 4: 1.0
probabilidad de la palabra= 3.1240237425804435E-4
prob_total_c anterior0.001111571238639088
prob_total_c =3.4725749410780635E-7

```

Fig. 22 Probabilidad de la oración ingresada para la categoría 3.

```

*****Categoria 4*****
Probabilidad de la categoria/clase -->5.0/43=0.11627906976744186
0 Palabra: this
Categoria 0: 21.0
Categoria 1: 80.0
Categoria 2: 9.0
Categoria 3: 17.0
Categoria 4: 1.0
probabilidad de la palabra= 3.588087549336204E-4
prob_total_c anterior0.11627906976744186
prob_total_c =4.17219482480954E-5
-1 Palabra: rain
Categoria 0: -1.0
Categoria 1: -1.0
Categoria 2: -1.0
Categoria 3: -1.0
Categoria 4: -1.0
1 Palabra: real
Categoria 0: 4.0
Categoria 1: 25.0
Categoria 2: 1.0
Categoria 3: 1.0
Categoria 4: 1.0
probabilidad de la palabra= 3.588087549336204E-4
prob_total_c anterior4.17219482480954E-5
prob_total_c =1.4970200304304056E-8

```

Fig. 23 Probabilidad de la oración ingresada para la categoría 4

En la Fig24. muestra las probabilidades de las cinco categorías y se elige la categoría con la probabilidad a posteriori mayor como resultado final.

```

*****Resultado*****
palabraslocal [this, rain, real]
rango de probabilidad: 1.0134152835580976E-5
categoria final: 1
*****
Probabilidad de la categoria 0: 1.524209896767441E-6
total de palabras de categoria: 1412
Probabilidad de la categoria 1: 1.0134152835580976E-5
total de palabras de categoria: 4259
Probabilidad de la categoria 2: 1.2554958861605722E-7
total de palabras de categoria: 1915
Probabilidad de la categoria 3: 3.4725749410780635E-7
total de palabras de categoria: 1033
Probabilidad de la categoria 4: 1.4970200304304056E-8
total de palabras de categoria: 619
*****
Frase ingresada: this rain is real
Categoria::1
Genero: Hip Hop
Interprete: EMINEM

```

Fig. 24. Resultado de la clasificación.

Capítulo 5: Datos y Aprendizaje

Se asigna la etiqueta de la clase que maximiza la probabilidad de la categoría y se compara con la etiqueta de la clase para evaluar el desempeño.

Ahora que se tienen los criterios para la clasificación, se aplicarán para una muestra de 144 datos (coros / frase de las canciones).

Dentro de tabla:4: Pruebas de clasificador, se observan los resultados para ambos métodos Naive Bayes Multinomial NMB y Naive Bayes NB, con ello se obtendrán las matrices de confusión.

En la siguiente tabla, **C** es igual a correcto, **I** es igual a incorrecto, **C0** categoría 0, **C1** categoría 1, **C2** categoría 2, **C3** categoría 3, **C** categoría 4, **X** es igual a incorrecto más la categoría que arroja el clasificador y **✓** corresponde a la categoría a la que pertenece.

Pruebas para el texto de entrenamiento con los metodos Naive Bayes Multinomial y Naive Bayes Simple

Textofrase	Categoria a la que pertenece	Categoria que arroja el clasificador	NMB		NB	
			C	I	C	I
Is this the real life?	Categoria 0	<pre>*****Resultado***** palabraslocal [this, real, life] rango de probabilidad: 1.1037664516736711E-8 categoria final: 1</pre>		X C1	✓	
Open your eyes	Categoria 0	<pre>*****Resultado***** palabraslocal [open, your, eyes] rango de probabilidad: 1.1983750046742712E-10 categoria final: 1</pre>		X C1		X C1
Mama! Ooh!/Galileo	Categoria 0	<pre>rango mayor : 3.248018708877615E-4 *****Resultado***** palabraslocal [mama, /galileo] rango de probabilidad: 3.248018708877615E-4 categoria final: 0</pre>	✓		✓	
Mama, Just killed a man	Categoria 0	<pre>*****Resultado***** palabraslocal [mama, just, killed] rango de probabilidad: 7.253322396712694E-10 categoria final: 1</pre>		X C1	✓	
Nothing really matters	Categoria 0	<pre>rango mayor : 7.402802757219878E-10 *****Resultado***** palabraslocal [nothing, really, matters] rango de probabilidad: 7.402802757219878E-10 categoria final: 0</pre>	✓			X C1
Any way the wind blows	Categoria 0	<pre>*****Resultado***** palabraslocal [wind, blows] rango de probabilidad: 1.6330820322508258E-7 categoria final: 0</pre>	✓		✓	

I want to break free	Categoría 0	<pre> rango mayor : 9.051004888888888E-9 *****Resultado***** palabraslocal [want, break, free] rango de probabilidad: 9.347801046505361E-9 categoría final: 2 </pre>		X C2		X C3
I've fallen in love for the first time	Categoría 0	<pre> rango mayor : 2.100000000000000E-11 *****Resultado***** palabraslocal [fallen, love, first, time] rango de probabilidad: 2.583820303896885E-11 categoría final: 0 </pre>	✓		✓	
Living without you by my side	Categoría 0	<pre> *****Resultado***** palabraslocal [living, without, side] rango de probabilidad: 3.041121102887951E-10 categoría final: 0 </pre>	✓		✓	
The show must go on	Categoría 0	<pre> *****Resultado***** palabraslocal [show, must] rango de probabilidad: 2.6129312516013277E-6 categoría final: 0 </pre>	✓			X C1
Does anybody want to take it anymore?	Categoría 0	<pre> *****Resultado***** palabraslocal [does, anybody, want, take, anymore] rango de probabilidad: 3.55525345890135E-15 categoría final: 0 </pre>	✓			X C2
Inside my heart is breaking	Categoría 0	<pre> rango mayor : 2.100000000000000E-11 *****Resultado***** palabraslocal [inside, heart, breaking] rango de probabilidad: 1.892171060012007E-10 categoría final: 1 </pre>		X C1	✓	
What are we living for?	Categoría 0	<pre> rango mayor : 0.705770030672032E-7 *****Resultado***** palabraslocal [what, living] rango de probabilidad: 8.709770888671052E-7 categoría final: 0 </pre>	✓			X C3

I have to find the will to carry on	Categoría 0	<pre> Rango mayor : 4.923558076072225E-11 *****Resultado***** palabraslocal [have, find, will, carry] rango de probabilidad: 4.923558076072225E-11 categoría final: 0 </pre>	✓		✓	
Pressure on people, people on streets	Categoría 0	<pre> Rango mayor : 6.19267397746048E-12 *****Resultado***** palabraslocal [pressure, people, people, streets] rango de probabilidad: 6.19267397746048E-12 categoría final: 0 </pre>	✓		✓	
These are the days it never rains but it pours	Categoría 0	<pre> *****Resultado***** palabraslocal [these, days, never, rains, pours] rango de probabilidad: 5.840553154736432E-10 categoría final: 1 </pre>		X C1		X C1
Why can't we give love that one more chance	Categoría 0	<pre> *****Resultado***** palabraslocal [give, love, that, more, chance] rango de probabilidad: 1.37066369354733E-14 categoría final: 2 </pre>		X C2		X C2
Turned away from it all like a blind man	Categoría 0	<pre> *****Resultado***** palabraslocal [turned, away, from, like, blind] rango de probabilidad: 2.3247691960852047E-16 categoría final: 1 </pre>		X C1		X C2
'Cause love's such an old fashioned word	Categoría 0	<pre> *****Resultado***** palabraslocal [cause, love, such, fashioned, word] rango de probabilidad: 3.2333040674463017E-15 categoría final: 0 </pre>		X C1		X C1
This is our last dance	Categoría 0	<pre> Rango mayor : 3.0749348179293955E-9****4.2971611 *****Resultado***** palabraslocal [this, last, dance] rango de probabilidad: 3.0749348179293955E-9 </pre>		X C2	✓	

I don't wanna stop at all	Categoría 0	<pre> *****Resultado***** palabraslocal [wanna, stop] rango de probabilidad: 3.5927804709618247E-6 categoria final: 0 </pre>	✓			X C2
'Cause I'm having a good time, having a good time	Categoría 0	<pre> *****Resultado***** palabraslocal [cause, having, good, time, having, good, time] rango de probabilidad: 1.63841816369657E-18 categoria final: 0 </pre>	✓			X C1
I'm a shooting star leaping through the sky	Categoría 0	<pre> *****Resultado***** palabraslocal [shooting, star, leaping, through] rango de probabilidad: 1.104036327220325E-14 categoria final: 1 </pre>		X C1	✓	
That's why they call me Mister Fahrenheit	Categoría 0	<pre> *****Resultado***** palabraslocal [that, they, call, mister, fahrenheit] rango de probabilidad: 1.02450280006614E-8 categoria final: 1 </pre>		✓	✓	
I wanna make a supersonic man out of you	Categoría 0	<pre> *****Resultado***** palabraslocal [wanna, make, supersonic] rango de probabilidad: 5.850368442701889E-10 categoria final: 1 </pre>		X C1		X C2
Are you ready, hey, are you ready for this?	Categoría 0	<pre> *****Resultado***** palabraslocal [ready, ready, this] rango de probabilidad: 1.5327250358555273E-8 categoria final: 0 </pre>	✓		✓	

How do you think I'm going to get along Without you when you're gone	Categoría 0	<pre> *****Resultado***** palabraslocal [think, going, along] rango de probabilidad: 1.1824049125075042E-10 categoria final: 1 </pre>		X C1		X C1
And another one gone, and another one gone	Categoría 0	<pre> *****Resultado***** palabraslocal [another, gone, another, gone] rango de probabilidad: 8.36010984985245E-11 categoria final: 0 </pre>	✓		✓	
Can anybody find me somebody to love?	Categoría 0	<pre> *****Resultado***** palabraslocal [anybody, find, somebody, love] rango de probabilidad: 2.3440417794952553E-9 categoria final: 0 </pre>	✓		✓	
I have spent all my years in believing you	Categoría 0	<pre> *****Resultado***** palabraslocal [have, spent, years, believing] rango de probabilidad: 2.439845391793411E-13 categoria final: 0 </pre>	✓		✓	
I work hard everyday of my life	Categoría 0	<pre> rango mayor : 1.5262963256601131E-12****b.791bb *****Resultado***** palabraslocal [work, hard, everyday, life] rango de probabilidad: 1.5262963256601131E-12 categoria final: 2 </pre>		X C2		X C2
They say I'm goin' crazy	Categoría 0	<pre> *****Resultado***** palabraslocal [they, goin, crazy] rango de probabilidad: 2.4440542050408427E-10 categoria final: 1 </pre>		X C1		X C1
Can anybody find me somebody to love?	Categoría 0	<pre> *****Resultado***** palabraslocal [anybody, find, somebody, love] rango de probabilidad: 2.3440417794952553E-9 categoria final: 0 </pre>	✓		✓	

We are the champions, my friends	Categoría 0	*****Resultado***** palabraslocal [champions, friends] rango de probabilidad: 1.651227388164728E-6 categoria final: 0	✓		✓	
and we'll keep on fighting till the end	Categoría 0	*****Resultado***** palabraslocal [keep, fighting] rango de probabilidad: 5.08069965589147E-7 categoria final: 0	✓		✓	
No time for losers	Categoría 0	*****Resultado***** palabraslocal [time, losers] rango de probabilidad: 1.8145355913898105E-6 categoria final: 0	✓		✓	
Crazy little thing called love, yeah, yeah	Categoría 0	*****Resultado***** palabraslocal [crazy, little, thing, called, love, yeah, yeah] rango de probabilidad: 3.5624595032970E-16 categoria final: 0	✓		✓	
This thing called love	Categoría 0	*****Resultado***** palabraslocal [this, thing, called, love] rango de probabilidad: 8.681636221093537E-10 categoria final: 0	✓		✓	
I gotta be cool, relax, get hip!	Categoría 0	*****Resultado***** palabraslocal [gotta, cool, relax] rango de probabilidad: 5.321961930063915E-10 categoria final: 0	✓		✓	
This world has only one sweet moment	Categoría 0	*****Resultado***** palabraslocal [this, world, only, sweet, moment] rango de probabilidad: 9.894584168024071E-16 categoria final: 1		x C1	✓	
Oh, love when must die?	Categoría 0	*****Resultado***** palabraslocal [when, love, must] rango de probabilidad: 2.2200184061082042E-8 categoria final: 0	✓			X C3

But my touch with tears lips your	Categoría 0	<pre>Resultado..... palabraslocal [touch, tears, with, your, lips] rango de probabilidad: 6.237558975058385E-16 categoría final: 1 </pre>		X C1		X C3
Who wants to live forever?	Categoría 0	<pre>Resultado..... palabraslocal [wants, live, forever] rango de probabilidad: 8.363083032941866E-9 categoría final: 0 </pre>	✓		✓	
Time for me to just stand up, and travel new land	Categoría 1	<pre>Resultado..... palabraslocal [time, just, stand, travel, land] rango de probabilidad: 9.735171423054794E-15 categoría final: 1 </pre>	✓			X C0
Only way that I know how to escape from this 8 Mile Road	Categoría 1	<pre>Resultado..... palabraslocal [only, that, know, escape, from, this, mile, road] rango de probabilidad: 6.8213373361182E-12 categoría final: 1 </pre>	✓		✓	
Sorry momma I'm grown, I must travel alone	Categoría 1	<pre>Resultado..... palabraslocal [sorry, momma, grown, must, travel, alone] rango de probabilidad: 2.645897795815956E-10 categoría final: 1 </pre>	✓			X C2
I don't wanna be alone in the darkness	Categoría 1	<pre>Resultado..... palabraslocal [wanna, alone, darkness] rango de probabilidad: 8.325552664043833E-9 categoría final: 1 </pre>	✓			X C2

Hello darkness, my old friend	Categoria 1	<pre>*****Resultado***** palabraslocal (hello, darkness, friend) rango de probabilidad: 3.443751329221854E-9 categoria final: 1</pre>	✓			X C3
Yeah, we just wanted to share that with you	Categoria 1	<pre>*****Resultado***** palabraslocal (yeah, just, wanted, share, that, with) rango de probabilidad: 3.94747876454493E-17 categoria final: 1 normalizer: 0.31663947327462029</pre>	✓			X C0
you better lose yourself in the music, the moment	Categoria 1	<pre>*****Resultado***** palabraslocal (better, lose, yourself, music, moment) rango de probabilidad: 3.532820088123396E-14 categoria final: 2</pre>		X C2		X C2
the soul's escaping, through this hole that it's gaping	Categoria 1	<pre>*****Resultado***** palabraslocal (soul's, escaping, through, this, hole, that, it's, gaping) rango de probabilidad: 1.330461307130063E-18 categoria final: 1</pre>	✓			X C0
you only get one shot, do not miss your chance to blow	Categoria 1	<pre>*****Resultado***** palabraslocal (only, shot, miss, your, chance, blow) rango de probabilidad: 4.8751905812057825E-19 categoria final: 1</pre>	✓		✓	
I'm a space bound rocket ship and your heart's the moon	Categoria 1	<pre>*****Resultado***** palabraslocal (space, bound, rocket, ship, your, heart, moon) rango de probabilidad: 1.730140407164912E-12 categoria final: 1</pre>	✓			X C4

And I'm aiming right at you	Categoría 1	*****Resultado***** palabraslocal [aiming, right] rango de probabilidad: 5.168417946146298E-7 categoria final: 1	✓			X C2
Two hundred fifty thousand miles on a clear night in june	Categoría 1	*****Resultado***** palabraslocal [hundred, fifty, thousand, miles, clear, night, june] rango de probabilidad: 7.870497500852261E-14 categoria final: 1	✓			X C0
I'm Slim Shady, yes I'm the real Shady	Categoría 1	*****Resultado***** palabraslocal [slim, shady, real, shady] rango de probabilidad: 1.0363220991508114E-10 categoria final: 1	✓		✓	
Please stand up, please stand up?	Categoría 1	*****Resultado***** palabraslocal [please, stand, please, stand] rango de probabilidad: 1.4547776280149492E-10 categoria final: 1	✓		✓	
And if I'm lucky, you might just give it a little kiss	Categoría 1	*****Resultado***** palabraslocal [lucky, might, just, give, little, kiss] rango de probabilidad: 3.0294605330931367E-10 categoria final: 1	✓			X C2
I'm so sick and tired of bein' admired	Categoría 1	*****Resultado***** palabraslocal [sick, tired, bein, admired] rango de probabilidad: 1.1776387490380131E-14 categoria final: 1	✓		✓	
That I	Categoría 1	*****Resultado***** palabraslocal [that, wished, that] rango de probabilidad: 7.726365161715695E-9 categoria final: 1	✓			

wished that I would just die or get fired						X C0
In the paper, the news, everyday	Categoría 1	<pre> *****Resultado***** palabraslocal [paper, news, everyday] rango de probabilidad: 6.064947392038422E-10 categoria final: 1 </pre>	✓		✓	
'Til I collapse I'm spilling these raps long as you feel 'em '	Categoría 1	<pre> *****Resultado***** palabraslocal [collapse, spilling, these, raps, long, feel] rango de probabilidad: 4.549763151343109E-19 categoria final: 1 </pre>	✓			X C2
Music is like magic there's a certain feeling you get	Categoría 1	<pre> *****Resultado***** palabraslocal [music, like, magic, there, certain, feeling] rango de probabilidad: 2.234109662676756E-19 categoria final: 1 </pre>	✓		✓	
Adrenaline shots of penicillin could not get the illin' to stop	Categoría 1	<pre> *****Resultado***** palabraslocal [adrenaline, shots, penicillin, could, illin, stop] rango de probabilidad: 1.131779551633004E-21 categoria final: 1 </pre>	✓			X C2
Guess who's back, guess	Categoría 1	<pre> *****Resultado***** palabraslocal [guess, back, guess, back] rango de probabilidad: 2.032027928630216E-11 categoria final: 1 </pre>	✓		✓	

who's back?						
Some vodka that will jump start my heart quicker	Categoría 1	<pre>*****Resultado***** palabraslocal [vodka, that, will, jump, start, heart, quicker] rango de probabilidad: 8.116325499420002E-05 categoria final: 1</pre>	✓			X C3
Cause we need a little controversy	Categoría 1	<pre>*****Resultado***** palabraslocal [cause, need, little, controversy] rango de probabilidad: 1.8501440771558208E-11 categoria final: 1</pre>	✓		✓	
I'm beginning to feel like a Rap God, Rap God	Categoría 1	<pre>*****Resultado***** palabraslocal [beginning, feel, like] rango de probabilidad: 5.5188922589348354E-9 categoria final: 1</pre>	✓			X C2
But for me to rap like a computer must be in my genes	Categoría 1	<pre>*****Resultado***** palabraslocal [like, computer, must, genes] rango de probabilidad: 3.09130171621169098E-13 categoria final: 1</pre>	✓			X C0
Rappers are hungry looking at me like it's lunchtime	Categoría 1	<pre>*****Resultado***** palabraslocal [rappers, hungry, looking, like, lunchtime] rango de probabilidad: 4.151927495070768E-14 categoria final: 1</pre>	✓		✓	
I listened to your	Categoría 2	<pre>*****Resultado***** palabraslocal [listened, your, problems, listen, mine] rango de probabilidad: 1.9820630952070072E-14 categoria final: 2</pre>	✓		✓	

problems, now listen to mine						
Kinda counted on you being a friend	Categoría 2	<p>*****Resultado*****</p> <p>palabraslocal [kinda, counted, being, friend]</p> <p>rango de probabilidad: 4.264267716866193E-12</p> <p>categoría final: 2</p>	✓		✓	
The summer memory that just never dies	Categoría 2	<p>*****Resultado*****</p> <p>palabraslocal [summer, memory, that, just, never, dies]</p> <p>rango de probabilidad: 3.402648776068274E-15</p> <p>categoría final: 1</p>		X C1	✓	
We're up all night to get lucky	Categoría 2	<p>*****Resultado*****</p> <p>palabraslocal [night, lucky]</p> <p>rango de probabilidad: 2.4649569231619232E-5</p> <p>categoría final: 2</p>	✓			X C0
What keeps the planet spinning	Categoría 2	<p>*****Resultado*****</p> <p>palabraslocal [what, keeps, planet, spinning]</p> <p>rango de probabilidad: 1.2049707834150785E-13</p> <p>categoría final: 2</p>	✓			X C3
She's up all night for good fun, I'm up all night to get lucky	Categoría 2	<p>*****Resultado*****</p> <p>palabraslocal [night, good, night, lucky]</p> <p>rango de probabilidad: 5.899612246278105E-10</p> <p>categoría final: 2</p>	✓			X C0
She got a kung fu star as a body guard,	Categoría 2	<p>*****Resultado*****</p> <p>palabraslocal [kung, star, body, guard]</p> <p>rango de probabilidad: 1.339856426017643E-14</p> <p>categoría final: 2</p>	✓		✓	

an

Going coat to coast on a campaign trail,	Categoría 2	<pre>*****Resultado***** palabraslocal [going, coat, coast, campaign, trail] rango de probabilidad: 9.83729923598562E-18 categoria final: 2</pre>	✓		✓	
Baby' got an atom bomb,	Categoría 2	<pre>*****Resultado***** palabraslocal [baby, atom, bomb] rango de probabilidad: 1.0933101574860074E-9 categoria final: 2</pre>	✓			X C0
Do it faster, makes us stronger	Categoría 2	<pre>*****Resultado***** palabraslocal [faster, makes, stronger] rango de probabilidad: 1.6785727511682384E-8 categoria final: 2</pre>	✓		✓	
Harder, better, faster, stronger	Categoría 2	<pre>*****Resultado***** palabraslocal [harder, better, faster, stronger] rango de probabilidad: 6.988914222338722E-11 categoria final: 2</pre>	✓		✓	
Work harder, make it better	Categoría 2	<pre>*****Resultado***** palabraslocal [work, harder, make, better] rango de probabilidad: 1.83791515881572E-10 categoria final: 2</pre>	✓		✓	
One more time	Categoría 2	<pre>*****Resultado***** palabraslocal [more, time] rango de probabilidad: 4.467198809156014E-5 categoria final: 2</pre>	✓		✓	
Music's got me feeling so free	Categoría 2	<pre>*****Resultado***** palabraslocal [music, feeling, free] rango de probabilidad: 1.0516277077318533E-7 categoria final: 2</pre>	✓			X C1

You know we're gonna do it right	Categoria 2	<p>*****Resultado*****</p> <p>palabraslocal [know, gonna, right]</p> <p>rango de probabilidad: 5.301135931376441E-8</p> <p>categoria final: 2</p>	✓		✓	
Celebrate and dance so free	Categoria 2	<p>*****Resultado*****</p> <p>palabraslocal [celebrate, dance, free]</p> <p>rango de probabilidad: 7.108224320780118E-7</p> <p>categoria final: 2</p>	✓		✓	
Lose yourself to dance	Categoria 2	<p>*****Resultado*****</p> <p>palabraslocal [lose, yourself, dance]</p> <p>rango de probabilidad: 1.1549088538445949E-7</p> <p>categoria final: 2</p>	✓		✓	
I know you don't get chance to take a break this often	Categoria 2	<p>*****Resultado*****</p> <p>palabraslocal [know, chance, take, break, this, often]</p> <p>rango de probabilidad: 9.424727799983189E-10</p> <p>categoria final: 1</p>		X C1	✓	
You take my shirt and just go ahead and wipe up all the sweat	Categoria 2	<p>*****Resultado*****</p> <p>palabraslocal [take, shirt, just, ahead, wipe, sweat]</p> <p>rango de probabilidad: 4.818462373737763E-18</p> <p>categoria final: 2</p>	✓			X C0
technologic. technologic. technologic. technologic.	Categoria 2	<p>*****Resultado*****</p> <p>palabraslocal [technologic, technologic, technologic, technologic]</p> <p>rango de probabilidad: 1.094473708742304E-12</p> <p>categoria final: 1</p>	✓		✓	

View it, jam unlock it,	Categoría 2	*****Resultado***** palabraslocal [view, code, unlock] rango de probabilidad: 5.1248913632156659E-10 categoria final: 2	✓		✓	
Touch bring it, pay watch it	Categoría 2	*****Resultado***** palabraslocal [touch, bring, watch] rango de probabilidad: 1.1718918250553141E-9 categoria final: 2	✓			X C3
Just to leave this unresolved	Categoría 2	*****Resultado***** palabraslocal [just, leave, this, unresolved] rango de probabilidad: 3.611425497040707E-12 categoria final: 1		X C1	✓	
When you finally get involved face to face	Categoría 2	*****Resultado***** palabraslocal [when, finally, involved, face, face] rango de probabilidad: 7.081355449309648E-16 categoria final: 2	✓			X C3
All because I hoped that you'd be someone different	Categoría 2	*****Resultado***** palabraslocal [because, hoped, that, someone, different] rango de probabilidad: 2.1977188017335558E-17 categoria final: 1		X C1	✓	
It's not hard to go the distance	Categoría 2	*****Resultado***** palabraslocal [hard, distance] rango de probabilidad: 1.6739946148807626E-7 categoria final: 2	✓		✓	
Let the music in	Categoría 2	*****Resultado***** palabraslocal [music, tonight, just, turn, music] rango de probabilidad: 5.533480920241912E-14 categoria final: 2	✓			X

tonight, just turn on the music						C1
Music	Categoría 2	<pre>*****Resultado***** palabraslocal [music] rango de probabilidad: 8.543649605322694E-4 categoria final: 2</pre>	✓			X C1
These fragments of time	Categoría 2	<pre>*****Resultado***** palabraslocal [these, fragments, time] rango de probabilidad: 1.8089462011512043E-6 categoria final: 1</pre>		X C1		X C1
Turning our days into melodies	Categoría 2	<pre>*****Resultado***** palabraslocal [turning, days, into, melodies] rango de probabilidad: 7.657698760930179E-11 categoria final: 3</pre>		X C3		X C0
And it's crystal clear	Categoría 2	<pre>*****Resultado***** palabraslocal [crystal, clear] rango de probabilidad: 5.5799817162692094E-8 categoria final: 2</pre>	✓		✓	
Now put your hands in the air for the century!	Categoría 3	<pre>*****Resultado***** palabraslocal [your, hands, century] rango de probabilidad: 5.941097947596666E-9 categoria final: 3</pre>	✓		✓	
Hello..if you hear me	Categoría 3	<pre>*****Resultado***** palabraslocal [hello, hear] rango de probabilidad: 2.767633041444344E-6 categoria final: 3</pre>	✓		✓	

can you hear me?	Categoría 3	<pre>*****Resultado***** palabraslocal [hear] rango de probabilidad: 5.884788910442231E-4 categoria final: 3</pre>	✓		✓	
Escape me	Categoría 3	<pre>*****Resultado***** palabraslocal [escape] rango de probabilidad: 8.500250448416557E-4 categoria final: 3</pre>	✓			X C0
Forget shout friends	Categoría 3	<pre>*****Resultado***** palabraslocal [forget, shout, friends] rango de probabilidad: 4.085382283621251E-7 categoria final: 3</pre>	✓		✓	
don't shred me down to strips	Categoría 3	<pre>*****Resultado***** palabraslocal [shred, down, strips] rango de probabilidad: 1.5315397501860359E-10 categoria final: 3</pre>	✓		✓	
you're way to good at it	Categoría 3	<pre>*****Resultado***** palabraslocal [good] rango de probabilidad: 0.0010393459867480836 categoria final: 0</pre>		X C0		X C0
I feel you in my bones	Categoría 3	<pre>*****Resultado***** palabraslocal [feel, bones] rango de probabilidad: 1.0724128454505784E-5 categoria final: 3</pre>	✓			X C1
Take hand my up again	Categoría 3	<pre>*****Resultado***** palabraslocal [take, hand, again] rango de probabilidad: 1.9867796130126077E-10 categoria final: 1</pre>		X C1		X C0
You're	Categoría 3	<pre>*****Resultado***** palabraslocal [knocking, windows] rango de probabilidad: 6.128073425431878E-7 categoria final: 3</pre>	✓		✓	

knocking at my windows						
For as long as you are, here on this earth	Categoría 3	<pre>*****Resultado***** palabraslocal {long, here, this, earth} rango de probabilidad: 3.600193404049477E-11 categoria final: 3</pre>	✓			X C0
Feel alive, feel alive, feel alive, feel alive	Categoría 3	<pre>*****Resultado***** palabraslocal {feel, alive, feel, alive, feel, alive, feel, alive} rango de probabilidad: 1.894662119322609E-13 categoria final: 3</pre>	✓		✓	
I feel alive	Categoría 3	<pre>*****Resultado***** palabraslocal {feel, alive} rango de probabilidad: 1.0009186594872066E-5 categoria final: 3</pre>	✓			X C1
Make it so, it shakes your heart	Categoría 3	<pre>*****Resultado***** palabraslocal {make, shakes, your, heart} rango de probabilidad: 1.1933405990221465E-12 categoria final: 1</pre>		X C1		X C0
It's a not the things you say	Categoría 3	<pre>*****Resultado***** palabraslocal {things} rango de probabilidad: 0.0013731174124365206 categoria final: 3</pre>	✓			X C0
In my mind I can feel you coming	Categoría 3	<pre>*****Resultado***** palabraslocal {mind, feel, coming} rango de probabilidad: 1.2060875532715034E-8 categoria final: 3</pre>	✓		✓	
I'm still waiting	Categoría 3	<pre>*****Resultado***** palabraslocal {still, waiting} rango de probabilidad: 1.317439848625527E-7 categoria final: 1</pre>		X C1		X C0

I will be here	Categoría 3	*****Resultado***** palabraslocal [will, here] rango de probabilidad: 2.022264230352519E-5 categoria final: 3	✓		✓	
When the big world falls apart	Categoría 3	*****Resultado***** palabraslocal [when, world, falls, apart] rango de probabilidad: 6.977492874168896E-12 categoria final: 3	✓		✓	
You don't mind if life's not that pretty	Categoría 3	*****Resultado***** palabraslocal [mind, life, that, pretty] rango de probabilidad: 6.455164832087318E-12 categoria final: 3	✓		✓	
We could just run them, red lights	Categoría 3	*****Resultado***** palabraslocal [could, just, them, lights] rango de probabilidad: 9.143306662310266E-11 categoria final: 3	✓			X C1
Don't ever turn around	Categoría 3	*****Resultado***** palabraslocal [don't, ever, turn, around] rango de probabilidad: 1.2299739271717582E-9 categoria final: 2	✓			X C0
There ain't no reason to stay	Categoría 3	*****Resultado***** palabraslocal [there, ain't, reason, stay] rango de probabilidad: 1.8378477002232431E-9 categoria final: 3	✓			X C0
Oh, you'll always be my ritual	Categoría 3	*****Resultado***** palabraslocal [always, ritual] rango de probabilidad: 6.577465476630215E-6 categoria final: 3	✓		✓	
You know what I	Categoría 3	*****Resultado***** palabraslocal [know, what, want, meet, dare] rango de probabilidad: 8.657761147021062E-16 categoria final: 1		X		X

want, now meet me if you dare				C1		C2
Be there when the sun is rising	Categoría 3	<pre>*****Resultado***** palabraslocal {there, when, rising} rango de probabilidad: 1.2762831261650302E-8 categoria final: 3</pre>	✓			X C0
You are my diamond when I'm with you I will shine	Categoría 3	<pre>*****Resultado***** palabraslocal {diamond, when, with, will, shine} rango de probabilidad: 5.549110680487382E-13 categoria final: 3</pre>	✓		✓	
All my life I been looking for love	Categoría 3	<pre>*****Resultado***** palabraslocal {life, been, looking, love} rango de probabilidad: 9.823469898200403E-13 categoria final: 1</pre>		X C1		X C0
Teenage wasteland	Categoría 4	<pre>*****Resultado***** palabraslocal {teenage, wasteland} rango de probabilidad: 5.389272109549458E-7 categoria final: 4</pre>	✓		✓	
The happy ones are near	Categoría 4	<pre>*****Resultado***** palabraslocal {happy, ones, near} rango de probabilidad: 1.9144246877325448E-11 categoria final: 3</pre>		X C3		X C0
I get my back into my living	Categoría 4	<pre>*****Resultado***** palabraslocal {back, into, living} rango de probabilidad: 4.083935871192583E-10 categoria final: 1</pre>		X C1		X C1
I love rock and roll	Categoría 4	<pre>*****Resultado***** palabraslocal {love, rock, roll} rango de probabilidad: 1.4800122700721359E-9 categoria final: 0</pre>		X C0		X C0

Sing for the year	Categoría 4	<pre>*****Resultado***** palabraslocal [sing, year] rango de probabilidad: 1.5710710319619255E-6 categoria final: 4</pre>	✓			X C2
Born to be wild	Categoría 4	<pre>*****Resultado***** palabraslocal [born, wild] rango de probabilidad: 9.431226191711554E-7 categoria final: 4</pre>	✓			X C0
Yeah, darlin', gonna make it happen	Categoría 4	<pre>*****Resultado***** palabraslocal [yeah, darlin, gonna, make, happen] rango de probabilidad: 2.1949473920292917E-14 categoria final: 2</pre>	✓			X C0
We were born, born to be wild	Categoría 4	<pre>*****Resultado***** palabraslocal [were, born, born, wild] rango de probabilidad: 5.463950259438169E-12 categoria final: 4</pre>	✓			X C2

Tabla 4 Pruebas del clasificador

Se puede observar que las categorías en las cuales tuvo mejor desempeño son las categorías que mayor cantidad de palabras aportan al vocabulario, asignando la mayoría de las oraciones de prueba a la categoría1(Eminen) que cuenta con un total de 4259 palabras, como el clasificador selecciona la probabilidad de pertenencia mayor, asigna las oraciones a esta categoría.

5.2 Aprendizaje Supervisado

Al realizar la tabla 4, se obtuvieron datos para guiar el proceso de aprendizaje con algunas oraciones, la mayoría son coros de las canciones, con base a las pruebas realizadas, clasifica algunas oraciones/frases ingresadas con las canciones y categoría que pertenecen.

Por lo tanto, la Arquitectura del sistema que se había mencionado en el capítulo 4.4 cambia, quedando como se muestra en la Fig.25.

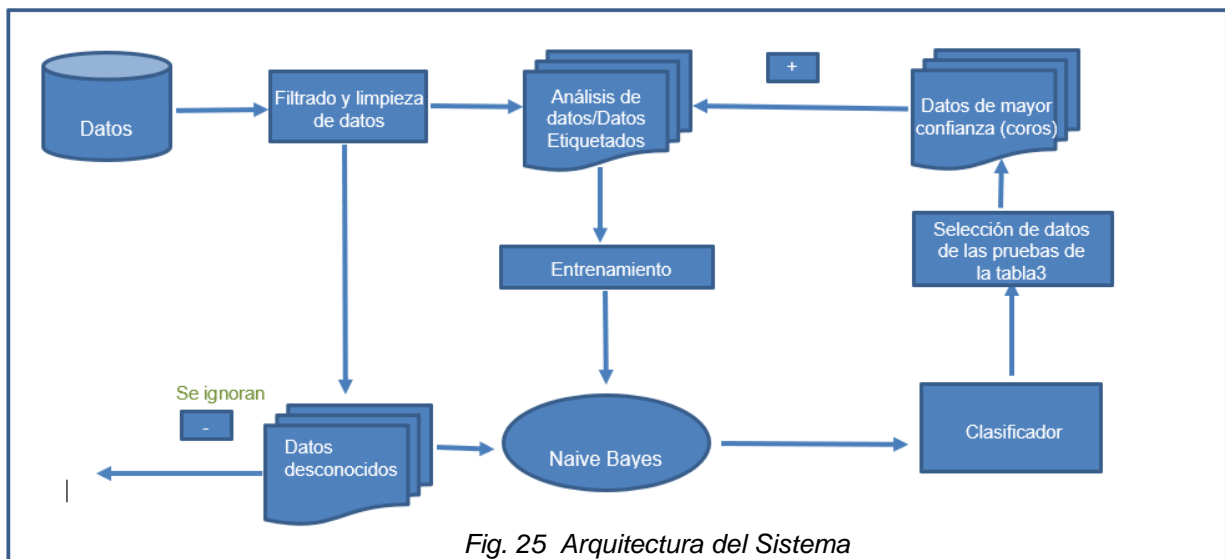


Fig. 25 Arquitectura del Sistema

De acuerdo con la tabla4 se agregan más datos para aprendizaje supervisado, pero ahora clasifica algunas canciones ya no solo el género al que pertenecen. Esto se realiza dentro de la clase Palabras en el método clasificar a partir de la línea 262 con los datos obtenidos en la tabla4, como se muestra en la Fig26.

Dentro de la clasificación se llama a el método leer, de la clase Archivo donde se leen las líneas del archivo. De acuerdo con la clasificación, se asigna el archivo a imprimir. Se puede ver un ejemplo en la línea 274 de la Fig26

```

260 //espacio muestral dataset
261 String clase=" ";
262 switch (cateTexto) { ///agente de la inteligencia artificial un agente inteligente deben basarse en el
263 /*razonamiento y en las conclusiones obtenidas a partir de la información que se
264 posee. Estos agentes tomarán la decisión más conveniente a la vista de esos datos y
265 del tiempo del que disponen:*/
266 case 0:
267     clase="QUEEN"; gene="ROCK";
268     if(resultado[cateTexto]==0.23255813953488372)
269     {
270         System.out.println("Error");
271     }else if((resultado[cateTexto]== 1.664101177351583E-4) || (resultado[cateTexto]==6.532328129003319E-7) || (resultado[cate
272 {
273     System.out.println("Canción: Bohemian Rhapsody ");
274     System.out.println(a.leer("C:\\Users\\makaf\\Desktop\\tratamiento de la informacion\\Letra\\Queen\\BohemianRhapsody
275
276
277     } else if((resultado[cateTexto]== 6.401969153457704E-13) || (resultado[cateTexto]== 1.9171611197051015E-10) || (resultado[c
278 {
279     System.out.println("Canción: I Want To Break Free ");
280     System.out.println(a.leer("C:\\Users\\makaf\\Desktop\\tratamiento de la informacion\\Letra\\Queen\\IWantToBreakFree
281
282     //////
283     }else if((resultado[cateTexto]== 1.9171611197051015E-10) || (resultado[cateTexto]==1.2781074131367343E-10) || (resultado[ca
284 {
285     System.out.println("Canción: The Show Must Go On ");
286     System.out.println(a.leer("C:\\Users\\makaf\\Desktop\\tratamiento de la informacion\\Letra\\Queen\\TheShowMustGoOn

```

Fig. 26. Dentro de los if anidados se ingresaron los valores verdaderos positivos de la tabla 3.

```

4  * and open the template in the editor.
5  */
6  package analisis;
7  import java.io.*;
8  /**
9   *
10  * @author makaf
11  */
12  public class archivos {
13      public String leer(String dir){
14          String cancion="";
15          try{
16              BufferedReader bf = new BufferedReader(new FileReader(dir));
17              //String temp="";
18              String linea;
19              while((linea=bf.readLine())!=null){
20                  System.out.println(linea);
21              }
22              //cancion=temp;
23          }catch(Exception e ){
24              System.err.println("No se encontro el archivo");
25          }
26          return cancion;
27      }
28  }
29

```

Fig. 27 Clase Archivo, método leer imprime las líneas del dataset.

Capítulo 6: Evaluación del desempeño del clasificador

Para observar el funcionamiento de la arquitectura de sistema Fig25, que se obtiene de la tabla4: Pruebas del clasificador, se ingresara el coro:” You don't mind if life's not that pretty” de la categoría:3 que corresponde a la canción I will be here, by Tiesto, se observa el resultado de la clasificación dentro del clasificador.

1. Filtrado y limpieza de datos: Las palabras menores de 4 caracteres son ignoradas y solo se consideran: “mind life that pretty” de la oración ingresada.
2. Análisis de Datos: aquí se encuentra BoW/ Tabla de frecuencias

```
debug:
ingrese el lyrics: You don't mind if life's not that pretty
*****Elementos por Categoría*****
  QUEEN=>10.0  EMINEM=>9.0  DAFT PUNK=>10.0  TIESTO=>9.0  OTHER=>5.0
*****Tabla de frecuencias *****
0-Palabra- this
  Queen 20.0  Eminem 79.0  Daft Punk 8.0  Tiesto 16.0  Other 0.0
1-Palabra- real
  Queen 3.0  Eminem 24.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
2-Palabra- life
  Queen 5.0  Eminem 6.0  Daft Punk 11.0  Tiesto 5.0  Other 2.0
117-Palabra- walk
  Queen 1.0  Eminem 6.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
118-Palabra- that
  Queen 10.0  Eminem 69.0  Daft Punk 9.0  Tiesto 11.0  Other 6.0
1288-Palabra- mind
  Queen 0.0  Eminem 2.0  Daft Punk 0.0  Tiesto 8.0  Other 0.0
1289-Palabra- song
  Queen 0.0  Eminem 0.0  Daft Punk 0.0  Tiesto 0.0  Other 0.0
2042-Palabra- pretty
  Queen 0.0  Eminem 0.0  Daft Punk 0.0  Tiesto 4.0  Other 0.0
2043-Palabra- ahead
```

Fig. 28 Frecuencia del token's de la oración

- Entrenamiento: se calcula la probabilidad para los token's para las cinco categorías

```

*****Categoria 3*****
Probabilidad de la categoria/clase -->9.0/43=0.20930232558139536
1288 Palabra: mind
Categoria 0: 1.0
Categoria 1: 3.0
Categoria 2: 1.0
Categoria 3: 9.0
Categoria 4: 1.0
probabilidad de la palabra= 0.0028116213683223993
prob_total_c anterior0.20930232558139536
prob_total_c =5.884788910442231E-4
calcular la categoria mayor 3: 5.884788910442231E-4
2 Palabra: life
Categoria 0: 6.0
Categoria 1: 7.0
Categoria 2: 12.0
Categoria 3: 6.0
Categoria 4: 3.0
probabilidad de la palabra= 0.0018744142455482662
prob_total_c anterior5.884788910442231E-4
prob_total_c =1.1030532165777379E-6
118 Palabra: that
Categoria 0: 11.0
Categoria 1: 70.0
Categoria 2: 10.0
Categoria 3: 12.0
Categoria 4: 7.0
probabilidad de la palabra= 0.0037488284910965324
prob_total_c anterior1.1030532165777379E-6
prob_total_c =4.1351573255022976E-9
2042 Palabra: pretty
Categoria 0: 1.0
Categoria 1: 1.0
Categoria 2: 1.0

```

Fig. 29. Probabilidad de los token's dentro de la categoría 3.

- Naive Bayes se calcula la probabilidad de la categoría y de los token's y se van multiplicando como lo dice la fórmula para cada categoría.

```

Categoria 3: 5.0
Categoria 4: 1.0
probabilidad de la palabra= 0.001562011871290222
prob_total_c anterior4.1351573255022976E-9
prob_total_c =6.459164832087313E-12
Rango mayor : 6.459164832087313E-12****6.459164832087313E-12 valor de i: 3 categoria final: 3

```

Fig. 30 Probabilidad de la categoría 3

5. Clasificación: se comparan los resultados obtenidos y se obtiene la categoría que maximiza la probabilidad.

```

*****Resultado*****
palabraslocal [mind, life, that, pretty]
rango de probabilidad: 6.459164832087313E-12
categoria final: 3
*****
Probabilidad de la categoria 0: 9.344226852449011E-14
total de palabras de categoria: 1412
Probabilidad de la categoria 1: 1.8032593344598638E-13
total de palabras de categoria: 4259
Probabilidad de la categoria 2: 1.0041423195132321E-13
total de palabras de categoria: 1915
Probabilidad de la categoria 3: 6.459164832087313E-12
total de palabras de categoria: 1033
Probabilidad de la categoria 4: 4.047370562546792E-14
total de palabras de categoria: 619
*****

```

Fig. 31 Resultado categoría final 3

6. Selección de datos: las probabilidades se incluyen en un switch donde se clasifican las canciones para imprimir la letra.

```

}else if((resultado[cateTexto]== 1)|| (resultado[cateTexto]==2.022264230392519E-5 )|| (resultado[cateText
{
    System.out.println("Canción:I Will Be Here");
    System.out.println(a.leer("C:\\Users\\makaf\\Desktop\\tratamiento de la informacion\\Letra\\Tiesto

```

Fig. 32 Probabilidades para la canción

7. Datos de mayor confianza: con estas probabilidades se imprimen: la letra de la canción, nombre, categoría e interprete.

En este ejemplo se muestra la clasificación de la oración: You don't mind if life's not that pretty de la tabla4: Pruebas del clasificador.

```

Frase ingresada: You don't mind if life's not that pretty
Categoria::3
Genero: Electropop
Interprete: TIESTO

```

Canción: I Will Be Here

I don't know, what went wrong
If I did, would it matter cause
It just wasn't enough
You know when the moment comes
To be strong, to resistance
And that is what
We're lead to believe

When the big world falls apart
And you think that the feeling will linger
You need somewhere to start
I will be here

And when it all seems to fall apart
You can't breathe
You don't know what you're thinking
You need somewhere to start
I will be here

You don't mind if life's not that pretty
It will soon disappear
It will be miles away
Away from here

Guess that things didn't work out
It will soon disappear and will be miles away
Away from here
You don't mind if life's not that pretty
It will soon disappear and will be miles away
Away from here

When the big world falls apart
And you think that the feeling will linger
You need somewhere to start
I will be here

And when it all seems to fall apart
You can't breathe
You don't know what you're thinking

Fig. 33 .Letra de la canción: I Will Be Hero by Tiesto

6.1 Matriz de confusión

En la matriz de confusión se analizan los resultados, mostrando los aciertos y errores para cada categoría con base a cuatro conceptos básicos: verdadero y falso positivo, verdaderos y falsos negativos.

VP	FP
FN	VN

Tabla 5 Matriz de confusión

Donde:

1. **Verdadero Positivo [VP]**: Los datos son positivos y en la prueba dio positivo
2. **Verdadero Negativo [VN]**: Los datos son negativos y en la prueba dio negativo
3. **Falso Negativo [FN]**: Los datos son positivos y en la prueba dio negativo
4. **Falso Positivo [FP]**: los datos son negativos y la prueba dio positiva

6.2 Métrica para la matriz de confusión

Accuracy: La exactitud mide el rendimiento de un sistema de IA

$$accuracy = \frac{\text{numero de predicciones correctas}}{\text{total de predicciones}} \quad (12)$$

Para la exactitud se considera como excelente los rangos mayores o igual a 0.9, como buena los rangos mayores o igual a 0.7 y menor a 0.7 tienen una mala exactitud.

Dos métricas mejores para evaluar clases desequilibradas son: precisión y recall.

Precision (Precisión): es la probabilidad de que la canción etiquetada en la categoría i corresponda realmente a esa categoría

$$\mathbf{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}} \quad (13)$$

Recall (Sensibilidad): es la probabilidad de que la canción que pertenece a la clase i es etiquetada dentro de esa clase

$$\mathbf{Recall} = \frac{\text{true positive}}{\text{true positive} + \text{Falso negativo}} \quad (14)$$

Puntaje F1: Resume la precisión y la sensibilidad en una sola métrica.

$$\mathbf{PuntajeF1} = \frac{2 * \text{Precision} * \text{Sensibilidad}}{\text{Precision} + \text{Sensibilidad}} \quad (15)$$

Teniendo en cuenta que el valor bajo es menor o igual a 0.5 y los valores altos es mayor o igual a 0.7 para precisión y recall se realiza la evaluación como lo marca la tabla 5.

Valores para Precision y Recall	Evaluación
Alta precisión y alto recall	El clasificador detecta correctamente la categoría.
Alta precisión y bajo recall	El clasificador no detecta la categoría muy bien, pero cuando lo hace es realmente confiable.
Bajo precisión y alto recall	El clasificador detecta bien la categoría, pero también incluye muestras de otras categorías.
Bajo precisión y bajo recall	El clasificador no detecta las categorías correctamente.

Tabla 6 Casos posibles

6.3 Matriz de confusión del clasificador

Para evaluar al clasificador se realizan las matrices de confusión de ambos métodos, donde se almacenan el conjunto de posibilidades entre la clase/categoría correcta, evaluando mediante oraciones de las canciones.

Los valores de las columnas son los aciertos de las categorías mientras que los valores de las filas son los errores de la clasificación.

En las siguientes tablas 6 y 7 se muestran las matrices de confusión para ambos métodos de clasificación NB Y NBM, la evaluación de la matriz de acuerdo con la métrica establecida anteriormente y el nivel de confianza de ambos métodos.

C A T E G O R I A R E A L	CATEGORIA						ASIGNADA	POR	NB	EVALUACIÓN
	C0	C1	C2	C3	C4	PRECISION	RECALL	PUNTAJE F1		
	25	8	6	4	0	0.462962963	0.5813953488	0.5154639175	El clasificador no detecta las categorías correctamente.	
	6	10	8	2	1	0.3846153846	0.3703703704	0.3773584906	El clasificador no detecta las categorías correctamente.	
	5	4	19	3	0	0.5	0.6129032258	0.5501246377	El clasificador no detecta las categorías correctamente.	
	11	3	1	13	0	0.5909090909	0.4642857143	0.52	El clasificador no detecta las categorías correctamente.	
	7	1	4	0	3	0.75	0.2	0.3157894737	El clasificador no detecta la categoría muy bien, pero cuando lo hace es realmente confiable.	
	MACROPOROMEDIO:					0.5376974877	0.4457909319	0.4874498964	Teniendo en cuenta la métrica el nivel de confianza es del 48% entra en el rango de media	
	ACCURACY					70/144	0.48%		Tiene una baja exactitud	

Tabla 7 Matriz de confusión NB.

El método NB presenta un desbalance debido a que solo contempla que el token se presente dentro de la categoría, con la finalidad de mejorar el nivel de confianza se realiza la clasificación de Naive Bayes Multinomial que toma en cuenta la frecuencia de las palabras en cada categoría.

Para medir el nivel de confianza del clasificador la métrica es la siguiente:

- ✓ Porcentaje del 70% su nivel de confianza es optima_
- ✓ Porcentaje del 80% su nivel de confianza es excelente.
- ✓ Porcentaje del 100% su nivel de confianza es perfecto

C A T E G O R I A R E A L	CATEGORIA						ASIGNADA	POR	NBM	EVALUACIÓN
	C0	C1	C2	C3	C4	PRECISION	RECALL	PUNTAJE F1		
C0	26	13	4	0	0	0.9285714286	0.6046511628	0.73239436	El clasificador no detecta la categoría muy bien, pero cuando lo hace es realmente confiable.	
C1	0	26	1	0	0	0.52	0.962962963	0.6753246753	El clasificador detecta bien la categoría, pero también incluye muestras de otras categorías.	
C2	0	4	26	1	0	0.7878787879	0.8387096774	0.8125	El clasificador detecta correctamente la categoría.	
C3	1	5	0	22	0	0.9166666667	0.7857142857	0.84615384	El clasificador detecta correctamente la categoría	
C4	1	2	2	1	9	1	0.6	0.75	El clasificador no detecta la categoría muy bien, pero cuando lo hace es realmente confiable.	
MACROPOROMEDIO:						0.8306233767	0.7584076158	0.7952026415	Teniendo en cuenta la métrica su nivel de confianza es del 79% entra en el rango de excelente	
ACCURACY						109/144	0.75%		Tiene una buena exactitud	

Tabla 8 Matriz de confusión NBM y la evaluación.

La clasificación con el método NBM mejoro el desempeño del clasificador, ya que considera la frecuencia del token's, más el suavizado de Laplace y el normalizar los datos permite que se ajuste mejor la categoría a la que pertenece.

6.4 Macro promedio

Como el clasificador tiene más de dos categorías entonces se utiliza el macro promedio, que calcula el rendimiento de cada categoría/clase y luego promedio sobre el total categorías.

categoría	0	1	2	3	4
Precisión	0.93	0.52	0.78	0.91	1
Recall	0.60	0.96	0.84	0.78	0.6

Tabla 9 Valores de precisión y recall de las cinco categorías.

$$Precision\ macro\ promedio = \frac{0.93 + 0.52 + 0.78 + 0.91 + 1}{5} = \frac{4.153}{5} = 0.83$$

$$Recall\ macro\ promedio = \frac{0.60 + 0.96 + 0.84 + 0.78 + .6}{5} = \frac{3.791}{5} = 0.75$$

Tiene alta precisión, alto recall y el valor del puntaje F1 nos da un resultado de 78% de confiabilidad para el clasificador.

$$PuntajeF1macro\ promedio = \frac{2 * (precision * recall)}{precision + recall} = \frac{2 * (0.6225)}{1.58} = 0.78$$

Teniendo en cuenta la métrica su nivel de confianza es excelente.

Capítulo 7: Conclusión

Con la realización de este trabajo se presenta una solución de ciencias de datos para el análisis de datos cumpliendo con los objetivos que se plantearon al inicio.

Dado que el algoritmo Naive Bayes Multinomial se utiliza para el análisis de datos de texto de múltiples clases. El clasificador identifica las categorías/intérprete de las canciones, se clasifica por género y en algunos casos el nombre de la canción e imprime la letra.

La clasificación está basada en tres aspectos principales:

1. Preprocesamiento: Al tener el conjunto de datos se realiza la limpieza de los datos, la cual se realizó mediante la tokenización, remoción de stop words y el método Bag of words.
2. Análisis de datos: se lleva a cabo mediante la tabla de frecuencias y normalizando los datos
3. Clasificación: se realiza mediante el método Naive Bayes Multinomial obteniendo como resultados que el clasificador tiene mayor efectividad en las categorías con mayor frecuencia de palabras.

Una de las desventajas del método es que, al no estar equilibradas las categorías respecto a la cantidad de palabras, se asigna la probabilidad a posteriori a la que tiene mayor vocabulario.

7.1 Trabajo a futuro realizar la clasificación de sentimientos.

En el presente todo lo que vivimos es digital, los datos fluyen a través de las aplicaciones que utilizamos esto es una oportunidad para comprender lo que nos rodea dentro de las ciencias sociales y de comunicación.

La ciencia de datos nos permite que la comunicación se adapte logrando aumentar la efectividad de los mensajes considerando una gran cantidad de perspectivas sobre diferentes temas fomentando que la toma de decisiones sea basada en datos.

Teniendo en cuenta la desventaja que presento método, las mejoras que se tomarían en cuenta para el trabajo a futuro.

Evaluar el desempeño del clasificador en:

- ✓ Aumentar el nivel de categorías
- ✓ Agregar más idiomas
- ✓ Equilibrar el vocabulario de las categorías

Los datos no son solo un recurso, son una ventaja para la toma de decisiones.

Bibliografía

- Abdulwahab Alazeb, M. A. (2021). Review on Data Science and Prediction on Computing . *International Conference on Computing and Data Science*, 8.
- Araujo, B. S. (2006). *Aprendizaje automatico:Conceptos basicos y avanzados*. Madrid : Pearson Prentice Hall.
- Arce, J. I. (20 de 11 de 2022). Obtenido de La matriz de confusión y sus metricas - Inteligencia artificial La matriz de confusión y sus métricas Inteligencia Artificial : <https://www.juanbarrios.com/la-matriz-de-confusion-y-sus-metricas/>
- datos.gob.es. (05 de 04 de 2018). *Ciencia de datos, machine learning y deep learning*. Obtenido de <https://datos.gob.es/es/blog/ciencia-de-datos-machine-learning-y-deep-learning>
- Deshpande, V. k. (2019). *Data Science Concepts and Practice*. United States: Morgan Kaufman.
- Fernando Molina Graja, L. B. (2022). Demand and employability study of the data science engineering career in Ecuador. *Iberian Conference on Informattion Systems and Technologies* , 5.

Gonzalo Pajares Martinsanz, M. S. (2006). *Inteligencia Artificial e ingeniería del conocimiento*. Mexico: Alfaomega.

Hernández, L. d. (s.f.). *Programarfacil*. Obtenido de <https://programarfacil.com/tutoriales/fragmentos/servomotor-con-arduino/>

Iñigo Martinez, E. V. (2021). A survey study of success factors in data science projects. *International Conference on Big Data*, 6.

Juan J. Cuadrado-Gallego, Y. D. (2021). Classification and Analysis of Techniques and Tools for Data Visualizations Teaching. *IEEE Global Engineering Education Conference*, 7.

Jurafsky, D. &. (1999). *Speech and Language Processing*. New Jersey: Prentice-Hall.

Kotu, V. (219). Introduction. En B. Deshpande, *Data Science Concepts and Practice* (pág. 549). Cambridge, MA: Morgan Kaufmann Publishers , and imprint of Elsevier.

Martín, J. A. (Febrero de 2023). *Blog de Tecnología - IMF Smart Education*. Obtenido de Glosario de Procesamiento de Lenguaje Natural (NLP): <https://blogs.imf-formacion.com/blog/tecnologia/glosario-de-procesamiento-de-lenguaje-natural-nlp-202302/>

Scherz, A. (2018). Clasificación automática de papers de Ciencias de la Computación.

64.