

PONTIFICIA UNIVERSIDAD CATÓLICA DEL ECUADOR

FACULTAD DE HÁBITAT, INFRAESTRUCTURA Y
CREATIVIDAD

SISTEMAS DE INFORMACIÓN



PROYECTO DE TITULACIÓN

APLICACIÓN DE TÉCNICAS DE MINERÍA DE DATOS Y APRENDIZAJE
SUPERVISADO PARA EL ANÁLISIS DE SENTIMIENTOS EN
ENTREVISTAS SOBRE TORMENTAS E INUNDACIONES

AUTORES:

JOSÉ DAVID ACUÑA OCHOA,
JAIR ALEJANDRO CURAY MAURIZACA

DIRECTOR:

MIGUEL ORTIZ NAVARRETE ING. MTR.

QUITO DM, 22 DE ENERO DE 2026

DEDICATORIA

A mi padre, que con su ejemplo de perseverancia y trabajo duro fue, es y seguirá siendo mi guía en esta aventura que llamamos vida.

A mi madre, que cada mañana me recuerda el valor del amor y nunca rendirse para seguir los sueños.

A mi hermano, que con alegría y firmeza es mi escudo para sortear los baches y calamidades que vivo a diario. Sé que no es sido fácil para él, pero sé que él está feliz por cumplir su tarea autoimpuesta de “hermano mayor”.

A mis amigos, que ellos saben quiénes son, por ser cada día gris una luz que ilumina cualquier abismo.

A todos los que se fueron antes de ver este logro materializarse, sé que desde el cielo me vigilan; espero que sigan cuidando de mí.

Y a mi Dios, que le estoy agradecido por todo lo que he logrado y por todo lo que he vivido; nunca me abandones.

En definitiva, a todos ustedes, con todo mi cariño,

JOSÉ DAVID ACUÑA OCHOA

Dedicado a las increíbles personas que Dios me ha puesto delante.

A mi mamá y a mi papá, que con su ayuda incondicional me han permitido llegar hasta este punto, me he esforzado por cada día ser un mejor ser humano y un mejor hijo para que puedan estar orgullosos de mí.

A mi hermano que me ha apoyado incluso cuando yo soy el que debería apoyarlo más, le deseo mucha suerte en lo que le queda de etapa universitaria y en su brillante futuro.

A mi cachorra Lola que ha sido más como una hermana que como una mascota, por reconformarme cuando no me siento bien y estar conmigo cuando no me sentía bien.

A mi familia cercana que siempre han estado presentes y atentas de mí.

Se lo dedico a todos los amigos que conocí y a los que perduraron incluso después de emprender cada uno su camino de vida.

Quiero hacer una mención especial a mi abuelito Fausto y a mi perrito Tomás, sé que desde algún lugar me han observado y pese a que su luz se apagó ellos viven en mis recuerdos, necesito pedirles disculpas por no poder terminar a tiempo mi carrera, quería con toda mi alma pudiesen ver que lo logré, no tienen idea la falta que me hacen.

Agradezco, como no, a mi grupo de música favorito, LINKIN PARK, sus canciones fueron un refugio cuando todo se ponía difícil y fueron la inspiración y el empujón para que pudiese levantarme nuevamente y seguir avanzando incluso en mis peores momentos. Felicidades por su retorno.

No sé lo que será de mi vida después de esta etapa, no hay una línea clara de lo que se debe hacer o seguir, infinidades de caminos e infinidades de posibilidades, millones de promesas, miles formas de ganar y miles formas de perder; he decidido afrontar este vacío con optimismo y esperanza porque como lo dice mi canción favorita "THE HARDEST PART OF ENDING IS STARTING AGAIN"

JAIR ALEJANDRO CURAY MAURIZACA

AGRADECIMIENTO

Queremos expresar nuestro sincero agradecimiento a todos los profesores que han sido guías en este camino educativo y formativo de la universidad. Hemos aprendido muchas cosas, dentro del área profesional, como seres humanos, haciéndonos mejores personas y crecido como individuos.

Ante este último escalón de un logro para nuestras vidas, queremos recalcar también la participación de nuestro director de trabajo de titulación el ingeniero Miguel Ortiz, que no solo fue un docente, sino amigo, con el que no solo resolvió cualquier duda ante algún imprevisto, si no también colaborando con una parte muy importante que es a veces es olvidada en trabajo de ciencia de datos, la misma recolección de la información a trabajar.

Sin querer alargar estas humildes palabras, no nos queda nada más que decir GRACIAS TOTALES.

José David Acuña Ochoa y Jair Alejandro Curay Maurizaca

RESUMEN

El presente proyecto de titulación tiene como objetivo analizar el impacto emocional de desastres naturales a partir de una metodología CRISP-DM, para lo cual se prestó atención a los archivos de audio que recogen las entrevistas en francés que las personas afectadas por tormentas e inundaciones en Francia vivieron. El trabajo correspondiente a la metodología CRISP-DM se desarrolló con una fase de preparación de datos exhaustiva y cuidada. Utilizamos el programa Audacity para realizar limpieza acústica de los archivos y Pydub para segmentar por silencios, y el modelo large-v3 de Faster Whisper para realizar la transcripción automática.

El núcleo del trabajo técnico se enfoca en un sistema híbrido de detección de emociones, basando su diseño en modelos de aprendizaje profundo (Transformers) y un diccionario manual de más de 150 términos en el idioma de la lengua francesa, con el que se clasifican los testimonios aportados por los individuos a través de un espectro de 15 emociones. Los algoritmos de aprendizaje automático cuentan con una matriz de características de 130 dimensiones que combina vectorización semántica (TF-IDF), score de intensidad emocional y frecuencias léxicas, lo que permitió transformar un conjunto de narrativas cualitativas con ambigüedad a datos estructurados a fin de cuantificar la subjetividad de las víctimas.

El estudio concluye que la utilización de técnicas de minería de datos permite automatizar el análisis de impacto social en contextos de crisis de forma objetiva y escalable.

ÍNDICE DE CONTENIDO

CAPÍTULO 1: INTRODUCCIÓN	1
1.1 Tema.....	1
1.2 Justificación.....	1
1.3 Planteamiento del problema.....	2
1.4 Objetivos	3
1.4.1 Objetivo General.....	3
1.4.2 Objetivos Específicos	3
1.5 Antecedentes	4
1.6 Alcance.....	5
CAPÍTULO 2: MARCO TEORICO CONCEPTUAL	7
2.1 Tormentas.....	7
2.2 Inundaciones.....	7
2.3 Damnificación	8
2.4 Inteligencia Artificial	8
2.5 Aprendizaje Automático	9
2.6 Metodología CRISP-DM.....	10
2.7 Minería de Texto	11
2.8 Procesamiento del Lenguaje Natural (PLN)	11
2.9 Algoritmos no supervisados	12
2.9.1 K-Means	12
2.9.2 DBSCAN (Density-Based Spatial Clustering of Applications with Noise).....	12
2.10 Algoritmos supervisados	12
2.10.1 Support Vector Machine (SVM)	13
2.10.2 Random Forest.....	13
2.10.3 Redes neuronales	13
2.11 Análisis de sentimientos.....	13
2.12 Métricas de evaluación de modelos	14
2.12.1 Accuracy.....	14
2.12.2 Recall	15

2.12.3 Especificidad	15
2.12.4 F1-score	15
2.12.5 Matriz de Confusión	16
CAPÍTULO 3: METODOLOGÍA	17
3.1 Materiales	17
3.1.1 Audacity.....	17
3.1.2 Hardware con Soporte CUDA	17
3.1.3 Visual Studio Code.....	17
3.1.4 Python.....	18
3.1.5 Bibliotecas de Procesamiento de Audio y Texto.....	18
3.1.6 Scikit-learn.....	18
3.1.7 Modelo de Análisis de Emociones (DistilRoBERTa)	19
3.1.8 Seaborn (sns)	19
3.1.9 WordCloud	19
3.1.10 Matplotlib Colors (mcolors).....	20
3.1.11 Collections (Counter)	20
3.1.12 Time.....	20
3.2 Metodología	20
3.2.1 Comprensión del Negocio	21
3.2.2 Comprensión de los Datos.....	21
3.2.3 Preparación de los Datos	21
3.2.4 Modelado	21
3.2.5 Evaluación	21
3.2.6 Despliegue	22
CAPÍTULO 4: DESARROLLO	23
4.1 Descripción del conjunto de datos	23
4.2 Aplicación de la metodología CRISP-DM.....	24
4.2.1 Comprensión del negocio	24
4.2.2 Comprensión de los datos.....	26
4.2.3 Preparación de los datos	28
a. Tratamiento Acústico y Estandarización Inicial	28

b. Segmentación Inteligente de Audio	29
c. Transcripción Automática y Corrección Gramatical	31
d. Consolidación y Unificación de Testimonios	33
4.2.4 Modelado	34
a. Procesamiento de Lenguaje Natural y Detección de Emociones	35
i. Arquitectura del Sistema Híbrido	35
ii. Catálogo de las 15 Emociones Analizadas	36
iii. Procesamiento y Resultados Iniciales.....	36
iv. Normalización y Codificación	41
b. Algoritmos no Supervisados (Clustering).....	41
i. K-MEANS	42
ii. DB SCAN	44
c. Algoritmos Supervisados (Clasificación)	45
i. Balanceo de Clases y Partición del Dataset	46
ii. Máquinas de vectores de soporte	46
iii. Bosques Aleatorios	50
iv. Redes Neuronales	53
v. Evaluación.....	56
d. Análisis de resultados	57
i. Análisis Del Panel De Perfiles Emocionales	58
ii. Análisis Diversidad y Estadísticas.....	62
e. Despliegue (No aplica).....	66
CAPÍTULO 5: CONCLUSIONES Y RECOMENDACIONES.....	67
5.1 Conclusiones	67
5.2 Recomendaciones.....	68
BIBLIOGRAFÍA	70

ÍNDICE DE ILUSTRACIONES

Ilustración 1. Ciclo de vida CrispDm	10
Ilustración 2. Visualización de la onda de audio original y visualización de la onda de audio resultante tras la limpieza en Audacity	29
Ilustración 3. Estructura de carpetas y archivos tras la segmentación inteligente.	31
Ilustración 4. Resultado de la transcripción masiva con archivos .txt generados.	33
Ilustración 5. Archivo de texto final consolidado que unifica todas las partes de una entrevista.	34
Ilustración 6. Distribución de las emociones predominantes detectadas en las entrevistas.....	37
Ilustración 7. Nube de palabras General.	38
Ilustración 8: Nube de palabras por emoción	39
Ilustración 9: Extracto de código para la vectorización.....	40
Ilustración 10: Aplicación de Standar Scaler.....	41
Ilustración 11. Obtención del K óptimo.....	42
Ilustración 12. Distribución del dataset por clusters.....	44
Ilustración 13: Agrupamiento DBSCAN y detección de testimonios atípicos (ruido).....	45
Ilustración 14: Extracto del código para aplicación de balanceo.....	46
Ilustración 15: Extracto del código de los hiperparámetros para SVM	47
Ilustración 16: Mejor resultado de SVM	48
Ilustración 17: Matriz de confusión del modelo SVM.....	49
Ilustración 18: Extracto del código de los hiperparámetros para Random Forest	50
Ilustración 19: Mejor resultado de Random Forest.....	51
Ilustración 20: Matriz de confusión del modelo Random Forest.....	52
Ilustración 21: Extracto del código de los hiperparámetros para Redes Neuronales.....	54
Ilustración 22: Mejor resultado de Redes Neuronales	55
Ilustración 23: Matriz de confusión del modelo de Red Neuronal (MLP)	56
Ilustración 24: Comparativa de métricas de desempeño entre modelos supervisados.	57
Ilustración 25: Total de emociones detectadas	59
Ilustración 26: Intensidad promedio de las emociones detectadas.	60
Ilustración 27: Distribución de emociones por Cluster.....	61
Ilustración 28: Distribución relativa de las emociones dominantes.....	62
Ilustración 29. Keywords de las 10 emociones detectadas	63
Ilustración 30. Emocional por Documento.	63
Ilustración 31. Matriz de correlación entre las emociones detectadas.....	64
Ilustración 32. Frecuencia de palabra clave por emoción detectada en el análisis	65

ÍNDICE DE TABLAS

Tabla 1. Tabla de emociones básicas y ampliadas a detectar	36
Tabla 2. Tabla de emociones encontradas en los .txt.....	39
Tabla 3. Cantidad de atributos obtenidos de cada Feature Set	41
Tabla 4. Métricas obtenidas de cada modelo supervisado.....	57

CAPÍTULO 1: INTRODUCCIÓN

1.1 Tema

APLICACIÓN DE TÉCNICAS DE MINERÍA DE DATOS Y APRENDIZAJE SUPERVISADO PARA EL ANÁLISIS DE SENTIMIENTOS EN ENTREVISTAS SOBRE TORMENTAS E INUNDACIONES

1.2 Justificación

La inteligencia artificial (IA) ha logrado afirmarse como un recurso extraordinario en el análisis de extensos volúmenes de datos, sobre todo en el ámbito social, el cual mediante su uso se puede lograr un mejor entendimiento de las percepciones, emociones y comportamientos humanos ante diversas situaciones. En este sentido, el análisis de sentimientos se ha afianzado como una de las técnicas más empleadas para extraer las opiniones plasmadas en textos (entrevistas, aportes de testimonios, etc.), y por tanto permite describir y analizar los patrones emocionales y las tendencias en la comunicación de personas expuestas a diversos sucesos.

En el caso de inundaciones y tormentas, los sentimientos y emociones que son logrados indagar a través de los testimonios de las víctimas, testigos o incluso de comunidades que han visto algún impacto sobre su existencia pudiera contribuir a crear un aporte valioso para la orientación sobre la toma de decisiones, en la planificación de las políticas públicas así como también en la elaboración de estrategias para la asistencia psicológica o social. Sin embargo, los datos textuales que provienen de entrevistas suelen estar contaminados por la ambigüedad, por ruido, por un léxico que puede ser errático, lo que hace que se necesiten técnicas de minería de datos, técnicas de aprendizaje supervisado para llegar a extraer información que sea útil pero que a su vez esté estructurada.

La presente propuesta de titulación buscará explorar el potencial de la inteligencia artificial en lo que concierne al análisis de sentimientos derivados de entrevistas sobre inundaciones y tormentas que sean muestras empíricas de cómo las tecnologías emergentes pueden servir para que los fenómenos sociales sobre el impacto humano sean comprendidos socialmente. Se plantea poder reforzar e intensificar la formación profesional a partir de la aplicación práctica de modelos de machine learning y de la minería de textos sobre datos reales, y poder mostrar la aplicabilidad de estos modelos de machine learning en los fenómenos sociales y humanitarios.

1.3 Planteamiento del problema

Actualmente, los medios de comunicación, las organizaciones dedicadas a la acción humanitaria y los profesionales de la investigación social elaboran grandes cantidades de información textual resultante de entrevistas, ciertas evidencias, relatos y reportes sobre desastres naturales, en este caso sobre tormentas e inundaciones. Estos datos, sin embargo, contienen interpretaciones, sentimientos y vivencias de las personas expuestas. Pero si no son debidamente procesados y analizados, se convierten en información poco utilizada para intentar profundizar e intentar comprender los efectos que el impacto social y emocional de estas experiencias pueden provocar en una comunidad.

Una de las principales dificultades asumida aquí consiste en cómo pasar de esos audios, generalmente poco formales, vagos e imprecisos, y con ciertas variantes lingüísticas a información estructurada que nos permita determinar la existencia de ciertos sentimientos o emociones predominantes, como el miedo, la tristeza, la esperanza o la nostalgia. Pero también está el desafío de saber seleccionar y saber implementar correctamente ciertos algoritmos que sean adecuados para el análisis de sentimientos, principalmente desde la perspectiva de

aprendizaje supervisado debido a que se hace necesario entrenar modelos que sean orientativos o que tengan la capacidad de clasificar las emociones a partir de datos previamente etiquetados.

En función de esta problemática, se plantea la siguiente pregunta principal:

- ¿De qué manera el uso de técnicas de minería de datos y aprendizaje supervisado puede contribuir al análisis de sentimientos en entrevistas sobre tormentas e inundaciones en territorio francés?

Y como preguntas secundarias:

- ¿Qué técnicas de preprocesamiento son más adecuadas para limpiar, normalizar y preparar los textos de entrevistas para su análisis?
- ¿Qué algoritmos de aprendizaje supervisado generan mejores resultados en la clasificación de sentimientos en información relacionada con tormentas e inundaciones?
- ¿Qué métricas pueden utilizarse para evaluar la precisión y efectividad de los modelos propuestos en el análisis de sentimientos?

1.4 Objetivos

1.4.1 Objetivo General

Analizar datos textuales provenientes de entrevistas sobre tormentas e inundaciones siguiendo la metodología CRISP-DM, con el fin de desarrollar un modelo de aprendizaje supervisado que permita realizar análisis de sentimientos y contribuir a la comprensión del impacto emocional de estos eventos.

1.4.2 Objetivos Específicos

1. Aplicar el modelo CRISP-DM con sus fases de comprensión y preparación de datos a un conjunto de entrevistas de audio con temática de tormentas e inundaciones, priorizando la calidad y ajustamiento de estos para su análisis.

2. Realizar un análisis exploratorio de los datos y aplicar técnicas de procesamiento del lenguaje natural con el fin de llevar a cabo un estudio de análisis de sentimientos, modelado de temas (Topic Modeling) y extracción de entidades nombradas (NER), para así obtener una caracterización más completa del contenido textual.
3. Implementación de algoritmos de clustering como K-Means y DBSCAN con el fin de identificar patrones ocultos en los datos y encontrar nuevas formas de agrupamiento de los mismos.
4. Realizar la implementación y ajuste de algoritmos de clasificación supervisada, véase support vector machine, random forest, redes neuronales, con el objetivo de conocer y clasificar sentimientos que predominan en los textos a analizar.
5. Evaluar el desempeño obtenido de cada modelo por sus métricas como precisión, recall, F1-score y matriz de confusión, y según estos determinar la técnica más apropiada para el análisis de sentimientos.
6. Comparar los resultados obtenidos para identificar el algoritmo con mejor equilibrio entre precisión y capacidad de detección de sentimientos relevantes en los testimonios sobre tormentas e inundaciones.

1.5 Antecedentes

Dentro de varias regiones del territorio francés, se han vuelto frecuentes la repetición de tormentas intensas que van acompañadas de inundaciones. En esta situación concreta juega un papel fundamental el cambio climático y su repercusión en términos de un desarrollo de las infraestructuras físicas y psicológicas del conjunto de la población. Una sucesión de fenómenos que producen inquietud, percepción de riesgo y otras reacciones de orden emocional. Sin embargo, una gran parte de los testimonios que generan estos fenómenos se encuentran dispersos a raíz de

entrevistas, de pautas de comunicación impresas y de recursos cualitativos que, por lo general, no son bien abordados.

La falta de un procedimiento documentado, que integre elementos de minería de datos y modelos que integren de forma supervisada y/o no, la evaluación de los testimonios genera: la repetición de formas manuales de proceder, la subjetividad de la valoración en sí de los testimonios, la dificultad de procesar grandes volúmenes de información y la falta de datos cuantificables que puedan ayudar a la toma de decisiones. En consecuencia, instituciones y responsables pueden no contar con información suficiente y comprensible del estado de la población ante tormentas e inundaciones, que entorpece la planificación de una respuesta consensuada y la planificación preventiva y de intervención ante la emergencia.

Mediante modelos de análisis de sentimientos estos datos se pueden transformar, permitiendo una conversión de entrevistas de tipo audio no estructurado en un conocimiento que permite identificar tendencias emocionales, percepciones mayoritarias y preocupaciones del conjunto de la población. Este planteamiento tecnológico está pensado para ser un apoyo de los investigadores y de los gestores de riesgos, mediante una herramienta capaz de procesar grandes volúmenes de datos textuales y poder presentar una cierta visualización de los fenómenos intensos y extremos y el impacto en el ámbito emocional y social de un conjunto de fenómenos meteorológicos extremos.

1.6 Alcance

Incluye:

- Aplicación completa de las cinco primeras fases del modelo CRISP-DM a un conjunto de datos textuales provenientes de entrevistas sobre tormentas e inundaciones.

- Construcción, ajuste y validación de varios modelos de clasificación supervisada en Python (Jupyter Notebook) para la identificación de sentimientos en los textos.
- Medición del desempeño de los modelos utilizando métricas como accuracy, recall y F1-score, garantizando una evaluación objetiva de los resultados.
- Análisis comparativo de los resultados obtenidos para recomendar la técnica más adecuada en el análisis de sentimientos.

Queda fuera de alcance:

- Implementación del modelo en una plataforma en tiempo real para el procesamiento continuo de entrevistas.
- Estudios de usabilidad o aceptación por parte de analistas sociales o instituciones humanitarias.
- Despliegue final del modelo en una aplicación web o móvil para uso público.

CAPÍTULO 2: MARCO TEORICO CONCEPTUAL

2.1 Tormentas

Desde una perspectiva meteorológica, una tormenta se define como "una perturbación violenta de la atmósfera, que incluye vientos fuertes y usualmente lluvia, truenos, relámpagos o nieve" (National Severe Storms Laboratory, s.f.).

Los fenómenos que se describe en el presente escrito representan diferentes manifestaciones de inestabilidad atmosférica y pueden oscilar drásticamente en la intensidad y duración, desde tormentas eléctricas de carácter local hasta sistemas ciclónicos de carácter general que originan efectos importantes en el espacio físico y social.

A diferencia del concepto de tormenta que queda abordado en la física, en el presente trabajo de titulación, el término tormenta se encuentra asociado al desencadenamiento de la data. Las entrevistas llevadas a cabo tienen como objetivo de indagar la experiencia humana frente a este fenómeno; por lo tanto, resulta altamente importante no perder de vista la naturaleza destructible e impredecible de las tormentas relativas para poder contextualizar los sentimientos detectados.

2.2 Inundaciones

La Oficina de las Naciones Unidas para la Reducción del Riesgo de Desastres (UNDRR) define técnicamente a la inundación como el "desbordamiento de agua fuera de los confines normales de un curso de agua o de otros cuerpos de agua, que resulta en la acumulación temporal de agua en tierras que normalmente están secas" (UNDRR, s.f). Este fenómeno del agua puede darse de muchas maneras, desde crecidas de ríos lentas o normales hasta inundaciones repentinas, más peligrosas que las anteriores porque pueden aparecer muy pronto tras lluvias muy intensas, dejando poco tiempo para reaccionar a las poblaciones que corren el riesgo.

En esta investigación, la inundación no se considera solamente como una variable física, sino como la situación de crisis que gatilla el relato de los entrevistados. Al igual que con las tormentas, el interés reside en cómo este fenómeno altera la vida cotidiana y genera la correspondiente emoción.

2.3 Damnificación

En el ámbito de la gestión de riesgos, la damnificación se entiende como la materialización del daño directo sobre la población y sus bienes. La Comisión Económica para América Latina y el Caribe (CEPAL) establece que este concepto abarca los efectos que ocurren simultáneamente con el desastre o inmediatamente después, incluyendo "la destrucción total o parcial de los activos físicos, como viviendas, enseres domésticos y capital productivo" (CEPAL, 2014). A diferencia de una simple afectación, la damnificación implica un grado de severidad tal que las personas sufren pérdidas cuantificables en su patrimonio o medios de vida, requiriendo asistencia externa para su recuperación.

Dentro del presente estudio, el concepto de damnificación es clave para contextualizar la intensidad de los sentimientos analizados. No todas las entrevistas reflejan el mismo nivel de impacto.

2.4 Inteligencia Artificial

La inteligencia artificial, específicamente la subdisciplina que se especializa en generación de código se define como "modelos entrenados sobre grandes repositorios de software capaces de producir o completar código automáticamente a partir de instrucciones en lenguaje natural o ejemplos previos" (Nair et al., 2022).

La herramienta se aplicará durante la fase de desarrollo. Su intervención se centra en la Fase 3 de CRISP-DM para la transcripción de audio y, crucialmente, en la Fase 4, facilitando la implementación tanto de algoritmos supervisados (para clasificación de emociones) como no supervisados (para detección de patrones) en el análisis de sentimientos. Se espera obtener una reducción significativa en los tiempos de implementación técnica asegurando el total anonimato de los datos.

2.5 Aprendizaje Automático

El Aprendizaje Automático, conocido como *Machine Learning*, es una disciplina de la Inteligencia Artificial que permite a los sistemas aprender y mejorar a partir de la experiencia. Brown (2021) lo define como "una subdisciplina de la inteligencia artificial que da a las computadoras la capacidad de aprender sin ser explícitamente programadas". En esta misma línea, IBM lo describe como "el subconjunto de la IA centrado en algoritmos que pueden 'aprender' los patrones de los datos de entrenamiento y, posteriormente, hacer inferencias precisas sobre nuevos datos" (Bergmann, s.f.).

El ámbito de este trabajo de titulación presenta el aprendizaje automático como una de las herramientas que, mediante su uso en la fase de modelado de la metodología CRISP-DM, se utilizará para el entrenamiento de los algoritmos que permiten que las entrevistas sobre tormentas e inundaciones sean procesadas sin la intervención humana, de forma automática, de modo que el objetivo será llegar a que el sistema sea capaz de detectar patrones de sentimientos y clasificar las emociones en los testimonios, de tal forma que sea posible realizar un análisis mucho más ágil y objetivo que el tradicional, como el de procesar esas entrevistas de manera manual.

2.6 Metodología CRISP-DM

La metodología CRISP-DM (*Cross-Industry Standard Process for Data Mining*) es el estándar más utilizado para guiar proyectos de minería de datos. Se define según Shearer (2021) como "un enfoque sistemático, independiente de la industria y de la aplicación, para aplicar la minería de datos con el objetivo de obtener conocimiento relevante de los procesos de negocio". Básicamente, este modelo proporciona un ciclo de vida estructurado que permite transformar datos crudos en información útil, garantizando que el análisis sea riguroso y replicable sin importar el sector donde se aplique.

Dentro del marco del presente proyecto, la metodología CRISP-DM va a funcionar como la base que va a servir para organizar el flujo de trabajo. Esta metodología cuyas fases son visibles en la Ilustración 1, orientará las diferentes etapas del proceso, desde la comprensión y preparación de las transcripciones de las entrevistas, hasta la modelación y evaluación de los algoritmos de análisis de sentimientos, fases que se presentan detalladamente en el Capítulo 3 correspondiente a la Metodología que, junto a su implementación, van a permitir realizar el estudio de los testimonios sobre tormentas e inundaciones, de forma ordenada y secuencial, obteniendo así resultados concretos y que vayan en línea con los objetivos de investigación planteados.

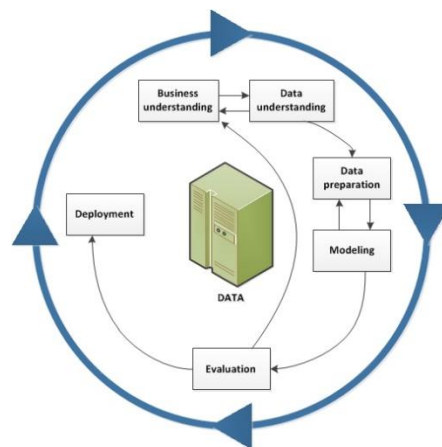


Ilustración 1. Ciclo de vida CrispDm
Fuentes: Google

2.7 Minería de Texto

La minería de texto es una técnica fundamental para gestionar y analizar información no estructurada. Talib, Hanif, Ayesha y Fatima (2016) la describen como "el proceso de extraer patrones interesantes y no triviales a partir de una gran cantidad de documentos de texto, con el objetivo de descubrir conocimiento útil que no es evidente mediante la simple lectura o análisis superficial".

La minería de texto será la técnica que usaremos en este trabajo de titulación para poder procesar adecuadamente el contenido de las entrevistas llevadas a cabo. Puesto que los relatos que encontramos son en lenguaje natural, la minería, junto con las técnicas pertinentes de texto, se hará cargo de limpiar y estructurar la información que contendrá antes del análisis que se llevará a cabo. Se trata, pues, de la técnica indispensable que ha de ser la que alimente los modelos de aprendizaje automático y que garantice que el análisis de sentimientos funcione con datos depurados de calidad.

2.8 Procesamiento del Lenguaje Natural (PLN)

El Procesamiento del Lenguaje Natural (PLN) es el campo que permite la interacción fluida entre el lenguaje humano y los sistemas informáticos. Zhang y Zhao (2023) lo definen como "una subdisciplina de la inteligencia artificial que se centra en desarrollar sistemas capaces de comprender, procesar y generar lenguaje humano natural mediante técnicas computacionales".

El PLN, en el desarrollo de este trabajo, actuará con una doble función, una esencial. En una primera parte, se aplicará a partir de técnicas de reconocimiento automático del habla para pasar de audios a texto procesable a tener en cuenta. Después, actuará como el motor tónico del análisis que, como es lógico, permitirá realizar la extracción de temas clave (*Topic Modeling*), tal y como se ha señalado, así como ejecutar el análisis de sentimientos. De esta forma, el PLN

facilitará pasar de grabaciones de voz en datos estructurados que puedan tener en cuenta los modelos de clasificación.

2.9 Algoritmos no supervisados

Los algoritmos no supervisados son métodos de aprendizaje automático que analizan conjuntos de datos sin etiquetas (u outputs definidos) para descubrir de manera autónoma patrones, estructuras o agrupaciones latentes en los datos (Sarker, 2021).

Para el estudio completo de este trabajo y encontrar nuevas formas de agrupamiento sin etiquetar se arrancará con el uso de dos conocidos algoritmos no supervisados, estos son:

2.9.1 K-Means

K-means es un método de agrupamiento que asigna cada muestra al clúster con el centro más cercano y actualiza iterativamente dichos centros para minimizar la variación dentro de los grupos (Bishop, 2006).

2.9.2 DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

DBSCAN forma clústeres conectando regiones densas del espacio de datos, permitiendo descubrir estructuras de forma irregular sin requerir especificar el número de clústeres (Bishop, 2006)

2.10 Algoritmos supervisados

Los algoritmos supervisados “son aquellos métodos de aprendizaje automático en los que se entrena un modelo a partir de un conjunto de datos etiquetado, es decir, en el que cada entrada (input) está asociada con una salida (output) correcta, con el objetivo de que el modelo aprenda una función que asigne nuevas entradas a salidas adecuadas” (Jiang, Gradus y Rosellini, 2020).

Durante el proyecto se los usara para el procesamiento de los datos con el fin de obtener los resultados deseados, los algoritmos que se usarán son los siguientes:

2.10.1 Support Vector Machine (SVM)

Al support vector machine se lo ha definido como “un algoritmo de aprendizaje supervisado que busca determinar un hiperplano óptimo que separa las clases de datos con el margen más amplio posible, y puede extenderse a tareas no lineales mediante la técnica de kernel” (Guido, Ferrisi, Lofaro y Conforti, 2023).

2.10.2 Random Forest

El Random Forest “es un método de aprendizaje automático basado en el ensamble de múltiples árboles de decisión, donde cada árbol se entrena con una muestra aleatoria de los datos y de las variables, y la predicción final se obtiene mediante el promedio (para regresión) o el voto mayoritario (para clasificación). Este enfoque mejora la precisión y reduce el riesgo de sobreajuste al combinar la diversidad de los modelos individuales” (Breiman, 2001).

2.10.3 Redes neuronales

Las redes neuronales están definidas como “un método de aprendizaje supervisado que emplea una arquitectura de múltiples capas de ‘neuronas’ interconectadas para modelar relaciones no lineales entre variables y extraer patrones complejos a partir de grandes volúmenes de datos” (Lee, 2024).

2.11 Análisis de sentimientos

El análisis de sentimientos es el campo de estudio que se ocupa de analizar emociones, actitudes y opiniones de las personas a partir del lenguaje escrito. Mao, Liu y Zhang (2024) lo caracterizan como "un proceso automático, rápido y eficiente de identificación de las opiniones y

sentimientos de los evaluadores, mediante el uso de técnicas de aprendizaje automático y procesamiento de lenguaje natural". Fundamentalmente, esta técnica busca determinar la polaridad emocional de un texto, ya sea positiva, negativa o neutra, para comprender la subjetividad detrás de las palabras.

La tarea primordial y última a la que se ve sometido todo el trabajo de transformación en el presente trabajo de titulación es el análisis de sentimientos. Su aplicación sobre las entrevistas permitirá la categorización y cuantificación de las emociones de las personas afectadas por tormentas e inundaciones. Más que una mera procedencia técnica de clasificar la información, esta tarea obtendrá los resultados críticos necesarios para alcanzar el propósito proyectado: traducir muchas verbalizaciones en índices claros que darán lugar a la posibilidad de finalizar el impacto emocional que se derive en la población objeto de estudio.

2.12 Métricas de evaluación de modelos

Estas métricas que se describirán a continuación se usarán para medir y comparar los resultados de cada uno de los algoritmos supervisados y no supervisados, según corresponda, y con esta comparación obtener los mejores modelos y poder obtener conclusiones de cada uno; entonces, se explicarán las diferentes métricas con las que se pone a prueba cada modelo que se usará:

2.12.1 Accuracy

La exactitud (accuracy) es la proporción de todas las predicciones que fueron correctas (positivas y negativas) respecto al total de predicciones realizadas. Matemáticamente:

$$\text{accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

donde TP = verdaderos positivos, TN = verdaderos negativos, FP = falsos positivos y FN = falsos negativos (Google Developers, 2025).

2.12.2 Recall

La sensibilidad (recall) es la proporción de los casos positivos reales que han sido correctamente identificados por el modelo. Matemáticamente:

$$\text{recall} = \frac{TP}{TP + FN}$$

donde TP = verdaderos positivos y FN = falsos negativos (Google Developers, 2025).

2.12.3 Especificidad

La especificidad es “la proporción de los casos negativos reales que han sido correctamente identificados como negativos por el modelo”. (StatisticsByJim, s. f.)

En formato matemático:

$$\text{Especificidad} = \frac{TN}{TN + FP}$$

donde TN = verdaderos negativos y FP = falsos positivos.

2.12.4 F1-score

La puntuación F1 (F1-score) es «la media armónica de la precisión y el recall (sensibilidad)», lo cual implica que penaliza valores extremos de cualquiera de los dos.

Matemáticamente:

$$F_1 = 2 \times \frac{\text{precisión} \times \text{recall}}{\text{precisión} + \text{recall}}$$

También puede expresarse como:

$$F_1 = \frac{2 \times TP}{2 \times TP + FP + FN}$$

donde TP = verdaderos positivos, FP = falsos positivos, FN = falsos negativos.

Esta métrica varía entre 0 y 1, siendo 1 el mejor valor (precisión = recall = 1) y 0 si la precisión o el recall son 0 (Mao, Liu y Zhang, 2024).

2.12.5 Matriz de Confusión

Una matriz de confusión es una tabla bidimensional que resume el rendimiento de un modelo de clasificación en relación con un conjunto de datos de prueba, donde una dimensión está indexada por la clase real de un objeto y la otra por la clase predicha por el clasificador (Ting, 2011). Para nuestro trabajo de titulación nos será de vital importancia para contrastar el funcionamiento de los modelos de clasificación.

CAPÍTULO 3: METODOLOGÍA

3.1 Materiales

Para el desarrollo del presente Proyecto de integración curricular, se han elegido herramientas tanto de hardware como de software que permiten el procesamiento de volúmenes elevados de audio y texto. Los materiales empleados son los que se describen a continuación:

3.1.1 Audacity

Audacity es un software editor de audio de código abierto: esto significa que es gratis. En este Proyecto, se usará como paso previo a la transcripción para realizar la limpieza de los audios originales, lo que implica eliminar distorsión, reducir la reverberación e igualar niveles de decibelios, asegurando que la calidad del sonido sea óptima para el reconocimiento automático del habla.

3.1.2 Hardware con Soporte CUDA

Con el fin de agilizar la transcripción y el entrenamiento de modelos, se utilizará un par de laptops con tarjetas gráficas NVIDIA RTX. La arquitectura CUDA permite hacer uso de la capacidad de procesamiento paralelo de las GPU y, por lo tanto, logra reducir significativamente los tiempos de ejecución en comparación con la opción de trabajar solo con la CPU.

3.1.3 Visual Studio Code

Se optará por este programa como el IDE principal; dada su multifuncionalidad y que soporta extensiones de Python, la escritura, depuración y organización de los scripts de processing y modeling se hace más fácil.

3.1.4 Python

El lenguaje de programación base del proyecto, es elegido por su sintaxis clara y su robusto ecosistema de bibliotecas especializadas en ciencia de datos e inteligencia artificial.

3.1.5 Bibliotecas de Procesamiento de Audio y Texto

Se implementarán diversas librerías especializadas para cubrir el flujo de trabajo:

- **Faster Whisper:** Implementada para la transcripción masiva de audio a texto, aprovechando modelos optimizados de Whisper para una mayor velocidad.
- **Pydub:** Empleada para la segmentación de archivos de audio largos en fragmentos más manejables para el sistema.
- **LanguageTool Python:** Utilizada para la corrección gramatical y normalización de los textos transcritos.
- **Transformers (Hugging Face):** Biblioteca fundamental para la implementación de modelos pre-entrenados de análisis de sentimientos y clasificación de secuencias.

3.1.6 Scikit-learn

Esta librería es fundamental para implementar algoritmos de aprendizaje automático supervisado (Random Forest, SVM, MLP), así como no supervisados (K-Means, DBSCAN). Igualmente, es también utilizada para la partición de conjuntos de datos y el cálculo de métricas de evaluación.

3.1.7 Modelo de Análisis de Emociones (DistilRoBERTa)

Es un modelo de Deep Learning que utiliza la arquitectura de los transformadores, en este caso específico, una versión optimizada y de menor peso del modelo RoBERTa. Este modelo es el responsable del procesamiento del texto de las entrevistas para poder clasificar automáticamente los estados emocionales de los participantes, determinando así si se trata de expresiones de alegría, tristeza, miedo, enojo, disgusto y sorpresa. Existen diferentes maneras de hacer esta implementación, la nuestra se basa en la librería *transformers* de Hugging Face, asignando al sistema el trabajo de cargar el modelo de un repositorio local en cache. Esto no solo mejora la rapidez de respuesta (al no tener que volver a realizar descargas), sino que hace posible que los testimonios puedan ser tratados sin suscitar problemas de privacidad, al poder hacerlo *offline*.

3.1.8 Seaborn (sns)

Este módulo, creado sobre Matplotlib, permite la creación de visualizaciones estadísticas de alto nivel. Por eso, es imprescindible para construir mapas de calor (heatmaps) de las matrices de correlación emocional, y de confusión porque proporciona una interpretación clara de cómo se entrelazan los sentimientos y cómo se distribuyen los errores de clasificación de los modelos.

3.1.9 WordCloud

Se utiliza para la creación de "Nubes de Palabras", un método de visualización que permite la identificación rápida, a simple vista, de los vocablos más recurrentes en los testimonios. Esta herramienta sirve para la validación semántica del proyecto, ya que da lugar a la representación gráfica del léxico predominante (por ejemplo, "*eau*", "*peur*" o "*merci*") relacionado a cada una de las categorías emocionales que se recogieron en los datos.

3.1.10 Matplotlib Colors (mcolors)

Este módulo de Matplotlib se ocupa de la gestión y personalización avanzada de las paletas de color del proyecto. Permite la asignación de colores explícitos y concretos para cada emoción (por ejemplo, lila para el miedo o verde para la esperanza) asegurando que todas las gráficas tengan identidad visual y que hagan posible el análisis en comparación.

3.1.11 Collections (Counter)

Esta utilidad se aplica para el análisis de frecuencias léxicas. Permite contabilizar de manera eficiente la aparición de cada palabra clave definida en el diccionario manual de francés. Gracias a Counter, se pudieron generar los gráficos de barras que muestran la densidad de términos por emoción, fundamentando estadísticamente la presencia de cada sentimiento en el corpus.

3.1.12 Time

Esta librería se utiliza para el monitoreo del rendimiento computacional. Permite medir el tiempo de ejecución exacto de cada algoritmo de aprendizaje automático y del proceso de extracción de emociones. Esta métrica es vital para evaluar la eficiencia técnica del sistema, especialmente considerando la carga que supone procesar modelos de lenguaje basados en transformadores como DistilRoBERTa.

3.2 Metodología

La metodología seleccionada para llevar a cabo este trabajo es la metodología CRISP-DM (*Cross-Industry Standard Process for Data Mining*), y el motivo es que es la más extendida para proyectos de minería de datos, ya que proporciona un ciclo de vida con un marco muy bien estructurado, pero al mismo tiempo flexible asegurando así que el análisis realizado tenga las características de ser riguroso y replicable.

CRISP-DM a diferencia de otras metodologías, permite una fuerte alineación de los objetivos sociales de la investigación con el resultado técnico de la misma. El proceso es dividido en las fases que a continuación se aplican a nuestro caso:

3.2.1 Comprensión del Negocio

En esta etapa inicial, se definirá la necesidad de entender el impacto emocional de las tormentas e inundaciones en territorio francés, traduciendo esta problemática social en un objetivo de clasificación de sentimientos y emociones.

3.2.2 Comprensión de los Datos

Se analizarán los archivos de audio recolectados, identificando factores como el ruido de fondo y las variaciones lingüísticas que podrían afectar la calidad de la información.

3.2.3 Preparación de los Datos

Esta será la fase más intensiva. Incluirá la limpieza de audios en Audacity, la segmentación con Pydub y la transcripción mediante Faster Whisper. Posteriormente, el texto se someterá a procesos de normalización y limpieza para eliminar ruido textual.

3.2.4 Modelado

Se implementarán y ajustarán diversos algoritmos. Se utilizarán técnicas supervisadas como Random Forest y SVM para clasificar emociones específicas, y técnicas no supervisadas como DBSCAN para identificar agrupamientos naturales de temas en los testimonios.

3.2.5 Evaluación

Los modelos se pondrán a prueba utilizando métricas de Accuracy, Recall y F1-score, comparando los resultados para determinar qué algoritmo logrará capturar mejor la subjetividad de los afectados.

3.2.6 Despliegue

Aunque el alcance de este trabajo no incluye una plataforma en tiempo real, se documentarán los hallazgos y recomendaciones técnicas para futuras implementaciones en organismos de gestión de riesgos.

CAPÍTULO 4: DESARROLLO

4.1 Descripción del conjunto de datos

El conjunto de datos (dataset) que se utiliza en este trabajo de titulación tiene un carácter cualitativo está constituido por distintas grabaciones en audio, extraídas de entrevistas y testimonios. Estos datos fueron obtenidos con el propósito de recoger las experiencias relatadas por las personas que fueron afectadas por fenómenos climáticos extremos como inundaciones y tormentas en el territorio francés.

Enumeramos ahora las características estructurales y técnicas de la data:

- Tipo de archivos: Los datos son representaciones en soportes digitales de audio, como los que utilizan la mayoría de grabadores de voz, en concreto formato .m4a y .wav.
- Nomenclatura: El dataset está compuesto por una serie de grabaciones interpeladas a partir de una estructuración que garantiza el anonimato a los entrevistados y se realiza siguiendo la forma "ENTRETIEN AVEC CAS [Número]".
- Volumen y Formato de Datos: Se procesaron 88 archivos audio originales, la mayoría en formato .m4a, cuyo tamaño total asciende a 34.8 GB, es decir, un volumen considerable que requiere buena gestión tanto de memoria como de almacenamiento en las fases de transcripción y análisis.
- Participantes y Dinámica de Entrevista: Se considera que el número de personas entrevistadas es no menor a 84 personas, cantidad que presenta oscilaciones, ya que en una misma entrevista pueden aparecer el testimonio de 2 o máximo 3 personas.

- Duración de los registros: Hay una gran cantidad de variabilidad en la longitud de las entrevistas. Hay registros breves de unos 3 minutos, entrevistas con un tiempo estándar de 30 minutos, y tenemos entrevistas más largas de hasta 3 horas de duración.
- Idioma y contexto: Todas las entrevistas fueron realizadas en idioma francés, capturando no solo el contenido textual del relato, sino también los matices emocionales, tonos de voz y expresiones propias del lenguaje natural de los afectados por las tormentas e inundaciones.
- Calidad inicial: Al realizarse las grabaciones en contextos diversos, los audios originales presentan distintos niveles de ruido de fondo, reverberaciones y variaciones en la intensidad del volumen (decibelios) lo cual motivó la necesidad de una fase posterior de limpieza y pre-procesado técnico.

4.2 Aplicación de la metodología CRISP-DM

En cuanto a la ejecución práctica de esta investigación, se ha seguido la secuencia de trabajo característica de la metodología CRISP-DM, organizando cada una de sus fases hacia la transformación de los testimonios en audio hacia indicadores de análisis de sentimientos. A continuación, se detallan las etapas que se han implementado, desde la definición de los objetivos hasta la validación de los modelos, haciendo que el procesado técnico trate exclusivamente la necesidad de entender cómo la población se ve afectada por las tormentas y por las inundaciones.

4.2.1 Comprensión del negocio

Esta fase inicial es esencial para conseguir que los objetivos técnicos del proyecto coincidan con las ya mencionadas necesidades del sector humanitario y del sector de gestión de

riesgos. El enfoque se separa en dos pilares estratégicos. Por una parte, la comprensión del impacto social derivado de los desastres; por otra parte, el derribo de las trabas tecnológicas en el tratamiento de información cualitativa.

Enfoque Social y Humanitario: Las tormentas e inundaciones en el francés han pasado de ser reivindicaciones puntuales a ser unas catástrofes hábiles, las cuales provocan no solo daños importantes, sino también alteraciones en la salud emocional de la población. A día de hoy, las organizaciones humanitarias y los investigadores producen documentos y grandes series de testimonios que incluyen las percepciones sobre temas como el miedo, la pérdida y la resiliencia. Sin embargo, los testimonios suelen subutilizarse debido a que permanecen en un estado disperso y sin un análisis que permita transformar unas narrativas en indicadores de utilidad determinadas para la planificación de políticas públicas o estrategias de atención psicológica.

Innovación y Necesidad Técnica: Desde el punto de vista técnico, el análisis que habitualmente se hace de estas entrevistas es manual, lo que representa dificultades significativas:

- **Subjetividad:** La interpretación de los sentimientos depende directamente del criterio del analista y puede introducir sesgos en el resultado.
- **Escalabilidad:** Aunque haya que procesar cientos de horas de grabación, su procesamiento manual es costoso e ineficaz para cualquiera de las instituciones.
- **Complejidad Lingüística:** Las narrativas suelen ser o informales o ambiguas, o simplemente tienden a cargar variaciones que impiden la estructuración automática.

Por lo tanto, en el presente proyecto se intenta definir una solución tecnológica mediante la inteligencia artificial que vuelve estos "audios no estructurados" en datos procesables. Se intenta que el sistema aprenda patrones emocionales automáticamente y entregue una clara y objetiva

representación del sentir colectivo. Automáticamente y a través del sistema de audios se detectan emociones como el miedo o la esperanza y se las ofrece a los gestores de riesgos como una acotada herramienta para ayudar directamente a la toma de decisiones y fundamentar su propia decisión en estadística científica socialmente orientada y así ayudar a mejorar la respuesta de las instituciones ante emergencias futuras.

4.2.2 Comprensión de los datos

Esta etapa no se limitó únicamente a una catalogación técnica, sino que se concentró en un diagnóstico profundo en lo que respecta a la utilidad y calidad de los datos recolectados. Se trató del análisis exploratorio de los registros de audio con el objetivo de identificar patrones, inconsistencias y problemáticas lingüísticas que pudieran desencadenar problemas en las etapas de modelado.

Los hallazgos principales tras la exploración de los datos se detallan a continuación:

- Evaluación de la Calidad y Ruido Ambiental: Se detectó que la calidad de la señal no es homogénea, dado que, al tratarse de entrevistas de campo, muchos registros muestran un exceso de ruido ambiental, distorsión y saturación, lo cual supone un riesgo para la exactitud de la transcripción automática. La mencionada observación fue decisiva para inferir que el preprocesado dejaría de ser opcional para convertirse en un requerimiento obligatorio que salve la integridad del análisis del sentimiento.
- Análisis de la Complejidad Lingüística: Al trabajar con audios en lengua francesa, se pudo observar una diversidad de registros lingüísticos que iban desde el lenguaje formal hasta las expresiones coloquiales llenas de regionalismos. Dicha diversidad es un activo para el análisis de sentimientos, pero lo que puede complicar un

proceso técnico para los modelos de Procesado de Lenguaje Natural (PLN) debiendo tener en cuenta localidades, pausas dubitativas variaciones en la entonación emocional.

- **Desafíos de la Gran Escala (Duración):** La existencia de audios de hasta 3 horas de duración planteaba una dificultad infraestructural. Se determinó que si se procesan archivos de tal peso de forma lineal la memoria de los modelos de transcripción quedaría saturada, haciéndose necesario abordar una determinada forma de segmentación técnica, a partir del silencio, para no perder el contexto del testimonio.

- **Representatividad y Relaciones:** A partir de la nomenclatura de los archivos (ejemplo "CAS 0124", "CAS 0202") se pudo establecer una estructura de seguimiento. Se observó que determinados casos pueden tener varias "partes" o entrevistas de seguimiento en diferentes fechas (ejemplo: los casos de abril y de julio de 2021), lo que permite concluir que el dataset no solo ofrece una foto estática del sentimiento, sino que permite observar la evolución emocional del testimonio a lo largo del tiempo.

- **Validación de Metadatos:** Se verificó que, aunque los formatos varían entre .mp3 y .wav, estos ofrecen la información suficiente para poder mantener la trazabilidad de cada testimonio, de manera tal que cada "sentimiento" detectado puede ser enlazado con el evento meteorológico y el sujeto al que le pertenece sin romper el anonimato.

Este diagnóstico nos permitió concluir que, a pesar de que el conjunto de datos es rico en matices emocionales, su carácter "sucio" y heterogéneo necesita un proceso de curación técnica con mucho rigor para poder transformar el ruido acústico en conocimiento social estructurado, en ocasiones se necesitan de varios intentos.

4.2.3 Preparación de los datos

La preparación de los datos se convirtió en una etapa extremadamente importante para garantizar el éxito del modelado posterior dada la naturaleza no estructurada y la variabilidad de calidad observada en la fase anterior. Esta etapa imponía un alto volumen de procesamiento para transformar registros de sonido crudo en entradas limpias y estandarizadas y adecuadas para la transcripción automática.

Este proceso se dividió en varias subetapas secuenciales, comenzando con el tratamiento de la señal acústica.

a. Tratamiento Acústico y Estandarización Inicial

Este primer paso, la mejora de la calidad de la señal acústica, se llevó a cabo con la ayuda de software especializado, en este caso, Audacity. El principal objetivo fue la maximización de la relación señal-ruido (SNR); por tanto, se trató de que la voz de los entrevistados fuese lo más clara posible en el contexto de ruidos de fondo.

Se aplica un flujo de trabajo de limpieza que incluía las siguientes intervenciones técnicas en cada uno de los archivos:

- Nivelación de decibeles (Normalización): Se escuchó la amplitud de la señal para que el volumen fuese el mismo en todas las grabaciones, impidiendo la existencia de picos que podrían causar distorsiones digitales (*clipping*) y elevando así aquellos segmentos con un volumen muy por debajo de la media.
- Reducción de ruido y saturación: Se aplican en el flujo de trabajo filtros específicos para poder extraer y atenuar los ruidos de fondo invariables (siseo, zumbido eléctrico, etc.) y

controlar la saturación en aquellos puntos donde la grabación original sobrepasaba el límite de captura.

- Eliminación de reverberación: En aquellas entrevistas que se llevaron a cabo en espacios cerrados y con eco, se aplicaron técnicas de des-reverberación para obtener una señal en la cual las palabras fueran más legibles.

Para poder garantizar la consistencia en el procesamiento posterior de las grabaciones, todos los archivos de resultado se exportan y se estandarizan al formato .wav (*Waveform Audio File Format*), manteniendo estrictamente la nomenclatura original de los archivos para poder garantizar la trazabilidad de cada caso concreto.

En la ilustración 2 se puede encontrar el reflejo visual de este procesamiento en la forma de onda de uno de los audios de muestra:

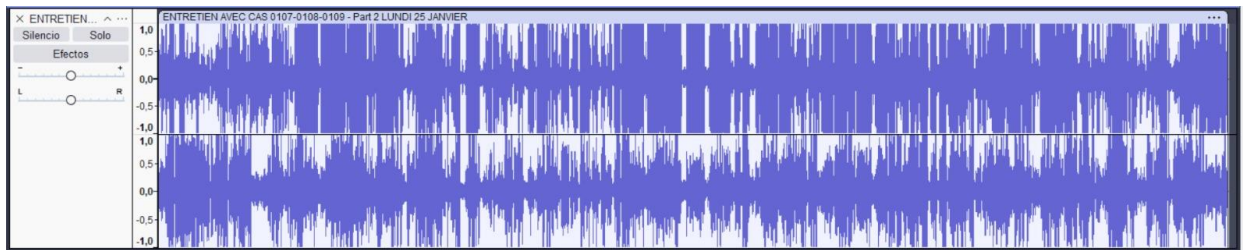


Ilustración 2. Visualización de la onda de audio original y visualización de la onda de audio resultante tras la limpieza en Audacity
Fuente: Autores del documento

b. Segmentación Inteligente de Audio

Concluida la estandarización de la calidad acústica, se procedió a la segmentación de los archivos. Teniendo en cuenta que el dataset contiene entrevistas de larga duración (como de 3 horas), su procesamiento continuo provocaría errores de memoria y cuellos de botella durante la transcripción. Para resolver esta problemática, se desarrolló un script en Python, aprovechando la biblioteca Pydub, que permitiera realizar una segmentación inteligente basada en los silencios.

Esta segmentación se ejecutó bajo los siguientes criterios técnicos:

- **Detección de Umbral Dinámico:** El algoritmo calculó automáticamente el nivel medio de los decibeles (dBFS) de cada entrevista para definir adecuadamente un umbral de silencio. Esto permitió que el sistema fuera capaz de diferenciar correctamente los silencios que interrumpen el habla (siendo estos últimos silencios naturales), sin importar el volumen de grabación que tuviera cada caso.
- **Puntos de Corte Estratégicos:** En lugar de realizar cortes arbitrarios por tiempo, la segmentación identificó el punto medio de los silencios que duraban más de 1.2 segundos. Esto garantizaba que las oraciones no fueran abruptamente cortadas, manteniendo la coherencia semántica necesaria para el análisis de sentimientos posterior.
- **Estandarización de Fragmentos:** Los cortes se agruparon en fragmentos de entre 5 y 10 minutos. Este intervalo de tiempo se determinó como el rango sensato de tiempo en el que se maximiza el rendimiento de la tarjeta gráfica (GPU) durante la transcripción sin comprometer la estabilidad del entorno de desarrollo.
- **Gestión de Salida:** Como podemos ver en la Ilustración 3, los segmentos obtenidos mediante los cortes se exportan en formato .wav, con la misma nomenclatura secuencial (ej. parte_000.wav, parte_001.wav). Los grupos de fragmentos se agruparon automáticamente en carpetas, nombrándolas a partir del audio original, facilitando la trazabilidad y el orden del dataset. Este proceso nos da como resultado final una estructura de archivos manejable y lista para alimentar los modelos de lenguaje natural.

Name	#	Title	Contributing artists	Album
parte_000				
parte_001				
parte_002				
parte_003				
parte_004				
parte_005				
parte_006				
parte_007				

*Ilustración 3. Estructura de carpetas y archivos tras la segmentación inteligente.
Fuente: Autores del documento*

c. Transcripción Automática y Corrección Gramatical

Se realizó una transcripción de los segmentos de audio previamente procesados en archivos de texto estructurados. Para esto, nos pusimos en marcha con un protocolo de transcripción masiva a partir de la librería Faster Whisper, y utilizamos su modelo large-v3, el que ofrece un rendimiento aceptable y rápido para el idioma francés.

El flujo de trabajo realizado estaba pensando para poder procesar automáticamente los fragmentos de audio segmentados que topamos en la fase anterior, y está predeterminado a ciertos criterios técnicos:

- **Procesamiento Masivo y Trazabilidad:** El script que se preparó recorrió de manera iterativa cada carpeta audio, identificando los archivos .wav y generando un archivo .txt correspondiente para cada archivo .wav. La numeración y la nomenclatura se mantuvieron de manera estricta y siguiendo así el resto de las operaciones (ej. “parte_000.txt”), lo que asegura que el resultado final pueda ser vinculado en todo momento con su audio original.

- Corrección Gramatical Avanzada (*Offline*): Para subsanar los errores de base en el reconocimiento de la voz, se integró la herramienta LanguageTool de forma offline. Permitiendo con ello realizar la corrección automática de concordancias de gramática y errores de escritura de la transcripción del testimonio a partir de textos sin requerir conexión a servidores remotos, preservando de este modo la privacidad de los testimonios.

- Normalización del Contexto Regional: Siendo las entrevistas de regiones como Tende, Breil-sur-Roya o Sospel, se creó además un diccionario de nombres geográficos correctos. Con el uso de algoritmos de comparación de cadenas (*SequenceMatcher*), con esto el sistema logra identificar y corregir la fonética errónea de ciudades o parajes para que el contexto geográfico de las inundaciones no se perdiera.

- Refinamiento de Patrones Específicos: Se establecieron manualmente "gold standard rules" para corregir ambigüedades críticas del habla. Un ejemplo fue la distinción de términos con pronunciaciones similares y sentidos opuestos fonéticamente (la confusión entre "Plus" y "Police" en francés), buscando así que lo que se intente narrar siga siendo comprensible para el análisis posterior.

De este modo, como se observa en la Ilustración 4 se logra obtener un repositorio de texto limpio y normalizado, siendo cada testimonio con fragmentos corregidos.

Name	Date modified	Type	Size
parte_000	12/18/2025 9:26 AM	Documento de tex...	8 KB
parte_001	12/18/2025 9:27 AM	Documento de tex...	9 KB
parte_002	12/18/2025 9:28 AM	Documento de tex...	9 KB
parte_003	12/18/2025 9:29 AM	Documento de tex...	8 KB
parte_004	12/18/2025 9:30 AM	Documento de tex...	8 KB
parte_005	12/18/2025 9:31 AM	Documento de tex...	8 KB
parte_006	12/18/2025 9:32 AM	Documento de tex...	8 KB
parte_007	12/18/2025 9:33 AM	Documento de tex...	8 KB

*Ilustración 4. Resultado de la transcripción masiva con archivos .txt generados.
Fuente: Autores del documento*

d. Consolidación y Unificación de Testimonios

Como etapa final del preprocesamiento, se procedió a la reconstrucción de las narrativas completas. Este paso se fue necesario dado que todo el trabajo de transcripción era realizado a partir de fragmentos de audio de corta duración (en el rango de franja temporal de 5 a 10 minutos) desde un punto de vista informático (con el fin de no sobrecargar este último). Esta práctica tenía entonces como resultado un repertorio de archivos de texto separados para una misma entrevista.

Para poder consolidar debidamente esta información, se creó un script de consolidación en Python que realizaba estas operaciones para así tener “armada” toda la entrevista en un mismo documento .txt:

- Reagrupación de los archivos: El código recorre cada carpeta de caso jerárquicamente y, dado el recorrido de archivos de texto, va reagrupando (ej. de la parte_000 a la parte_n). De esta manera, se garantizaba que el testimonio unificado conservaría el orden del discurso y la cronología de la evolución narrada por las personas afectadas.

- **Preservación de Trazabilidad:** Un atributo esencial era que cada archivo siguiese la herencia del nombre de la carpeta de origen; es decir, el cual corresponde con el nombre de audio original (ej. "ENTRETIEN AVEC CAS 0214 - VENDREDI 09 JUILLET 2021"). Esta práctica garantizaba en las fases de modelado y de evaluación, que cada indicador de sentimiento puede ser relacionado sin error con el testimonio sonoro y el contexto temporal correspondiente.

El resultado de esta fase es un repositorio de documentos de texto completos, depurados y normalizados que representan exactamente las realidades de los entrevistados frente a tormentas e inundaciones y una base fuerte y estructurada para el análisis. Como se ve en la Ilustración 5 se tiene un ejemplo de cómo termina la secuencia de unificación para una entrevista.

Name	Date modified	Type	Size
ENTRETIEN AVEC CAS 0214 - VENDREDI 09 JUILLET 2021	12/18/2025 8:47 PM	Documento de tex...	95 KB

Ilustración 5. Archivo de texto final consolidado que unifica todas las partes de una entrevista.

Fuente: Autores del documento

4.2.4 Modelado

Una vez normalizados, unificados y transcritos los testimonios, se pasó a la fase de modelado. La idea central de esta fase es aplicar técnicas de inteligencia artificial para traducir el lenguaje natural utilizado por los entrevistados en datos estructurados que permitan encontrar patrones emocionales y niveles de sentimiento. En este proyecto, el modelado no consistió en un único método. Se ejecutó una estrategia híbrida que combina técnicas de Deep Learning para la detección inicial de emociones y algoritmos de Machine Learning (supervisados y no supervisados) para la correspondiente clasificación y la identificación de agrupamientos naturales en los relatos.

Este paso permitía que a partir de una interpretación subjetiva de las entrevistas se pudiera realizar una medida técnica del impacto emocional de las inundaciones, tomando como vehículo

la capacidad de cálculo de las GPU para procesar la complejidad semántica de cada uno de los testimonios.

a. Procesamiento de Lenguaje Natural y Detección de Emociones

Esta fase era el núcleo analítico del proyecto, donde las transcripciones textuales se transformaron en datos cuantitativos emocionales. Para poder alcanzar una alta precisión en el contexto de las entrevistas en francés se implementó un sistema híbrido de detección que es capaz de resolver las limitaciones de los modelos de lenguaje genéricos.

i. Arquitectura del Sistema Híbrido

Visto que el lenguaje natural de los afectados por las tormentas era complejo, no se utilizó una sola técnica, sino que se optó por implementar una función de análisis de doble vía que combina:

- **Modelo de Aprendizaje Profundo (Transformers):** Se implementó el modelo michellejeili/emotion_text_classifier (o en su defecto distilroberta-base) cargado desde un repositorio local para la privacidad de los datos. Este modelo realiza un análisis de la estructura semántica y del contexto de los enunciados.
- **Análisis por Palabras Clave (Keywords):** Se desarrolló un diccionario exhaustivo en francés con más de 300 términos asociados a 15 emociones específicas. Esta vía tiene un peso mayor en la decisión final (80%) debido a su fiabilidad para detectar matices locales y términos técnicos del desastre.

El cálculo del puntaje emocional final de cada testimonio se rige por una suma ponderada:

$$S_{final} = (S_{keywords} \times 0.8) + (S_{transformer} \times 0.2)$$

ii. Catálogo de las 15 Emociones Analizadas

A diferencia de los análisis de polaridad básica (que termina siendo positivo, negativo o neutro), este modelo clasifica los testimonios en un espectro emocional grande, esencial para entender el complejo impacto humano que provocan las inundaciones o tormentas:

Categoría	Emociones Específicas Detectadas
Básicas	Alegría, Tristeza, Enojo, Miedo, Sorpresa, Disgusto
Ampliadas	Amor, Gratitud, Frustración, Ansiedad, Esperanza, Decepción, Orgullo, Vergüenza, Nostalgia

*Tabla 1. Tabla de emociones básicas y ampliadas a detectar
Fuente: Autores del documento*

iii. Procesamiento y Resultados Iniciales

Un total de 88 archivos de texto unificados, que recogen la total de las entrevistas recogidas, fueron procesados. El sistema recorrió el documento, lo dividió en diferentes bloques para evitar sobrecargar la GPU con información y, después, fue capaz de extraer las características emocionales más dominantes.

La ejecución del procedimiento genera la Ilustración 6, la cual muestra la emoción principal asignada a cada uno de los 88 archivos de audio estudiados: cada una de las barras representa cuántas veces una emoción determinada fue considerada como dominante en un testimonio dado; la suma de todas las barras da el total del corpus, lo que confirma que a un archivo le ha sido asignada sólo una categoría, cualitativamente determinada por aquel puntaje de confianza más alto. Esta visualización también es importante para identificar rápidamente los sentimientos que predominan en el relato general de los afectados por las inundaciones.

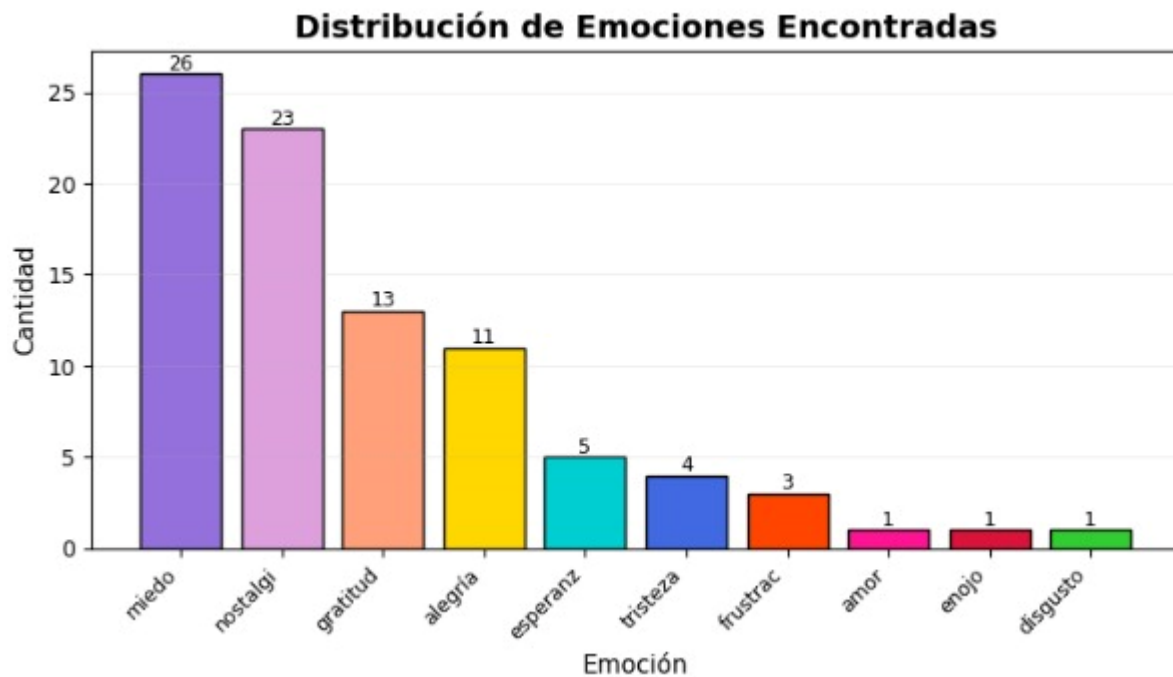


Ilustración 6. Distribución de las emociones predominantes detectadas en las entrevistas
Fuente: Autores del documento

A continuación, en la Ilustración 7 podemos observar una nube de palabras que nos da a entender las expresiones más recurrentes en los 88 archivos que han sido analizados, a partir de los cuales quedan evidenciados algunos de los contenidos que se repiten con mayor frecuencia en el corpus. Tal y como se puede observar, el tamaño de cada palabra es determinante con relación a la frecuencia de su repetición. A partir de esta nube podemos evidenciar que el discurso de los sujetos entrevistados se articula en torno a conceptos clave como el agua (eau), la casa (maison) y el acontecimiento mismo de la inundación. Pero también la nube puede ser considerada en sí misma una forma de justificación semántica global que da cuenta de cuáles han sido los elementos que preocupan más explícitamente a los entrevistados en las tormentas e inundaciones sobre el cual hemos tratado de obtener información.

De esta manera, para que los algoritmos de aprendizaje supervisado y no supervisado puedan procesar la información resulta necesario convertir tanto la nomenclatura emocional como los resultados previos en una representación numérica multifactorial. En este paso se realizó la construcción de la matriz de atributos (feature set) definitiva con la que se alimentan todos los modelos del proyecto.

El proceso de ingeniería de atributos que se describe se estructuró mediante la combinación de tres tipos de datos diferentes resultando en un total de 130 dimensiones o variables por cada entrevista:

- **Vectorización Semántica (TF-IDF):** Se aplicó la técnica Term Frequency-Inverse Document Frequency limitando, como se puede observar en la Ilustración 9, el análisis a las 100 palabras y combinaciones de palabras (n-gramas de 1 y 2 términos) cuyo análisis había mostrado mayor relevancia, y esto permitió representar no sólo palabras sueltas sino también pequeñas frases que permiten la representación de contexto climático o emocional.

```
vectorizer = TfidfVectorizer(max_features=100, ngram_range=(1, 2))
```

Ilustración 9: Extracto de código para la vectorización
Fuente: Autores del documento

- **Distribución Probabilística de Emociones:** Se añadieron las 15 puntuaciones numéricas de acuerdo al modelo de detección de emociones. Cada columna indica la posibilidad de (0-1) de una emoción concreta en el testimonio.

- **Frecuencia de Palabras Clave (Keywords):** Se añadieron 15 variables de recuento bruto de términos clave detectados con el diccionario manual de francés. Lo que refuerza el peso de términos que los modelos automáticos pueden pasar por alto.

iv. Normalización y Codificación

Dado que las variables tienen escalas muy diferentes (ej. el TF-IDF tiene valores decimales pequeños mientras que el recuento de palabras clave son números enteros), se aplicó una normalización (*Standard Scaler*), tal como puede observarse en la Ilustración 10. Lo cual centra los datos en una media de cero y una varianza de uno, de modo que se asegure que ninguna de las variables domine injustificadamente sobre el resto durante el entrenamiento de los modelos.

```
scaler = StandardScaler()  
X_scaled = scaler.fit_transform(X_combined)
```

Ilustración 10: Aplicación de *Standard Scaler*
Fuente: Autores del documento

Finalmente, la variable objetivo ("emoción principal") fue procesada mediante un *LabelEncoder*, transformando las 10 categorías emocionales detectadas en valores numéricos interpretables por las máquinas de vectores de soporte y redes neuronales. Todo esto lo podemos ver referenciado en la Tabla 3.

Componente del Feature Set	Cantidad de Atributos
Características de Texto (TF-IDF)	100
Puntuación de Emociones Ampliadas	15
Conteo de Palabras Clave (Keywords)	15
Total de Dimensiones (Input para ML)	130

Tabla 3. Cantidad de atributos obtenidos de cada Feature Set
Fuente: Autores del documento

b. Algoritmos no Supervisados (Clustering)

Una vez preparadas las 130 características técnicas se implementaron algoritmos de agrupamiento (clustering) para llevar a cabo un descubrimiento de patrones ciego, esto es, sin depender de las etiquetas previas. El uso de K-Means y DBSCAN se utilizó para conocer de qué

manera se agrupan los testimonios de forma natural partiendo de las similitudes léxicas y de los perfiles emocionales. K-Means permitió identificar grupos a partir de centroides y DBSCAN aportó una perspectiva de densidad para intentar detectar casos atípicos o "ruido" en las entrevistas.

i. K-MEANS

El primer método de agrupamiento implementado fue K-Means, el cual particiona los datos en k grupos basándose en la distancia hacia un centroide. Para este estudio, se realizó una búsqueda automatizada del número óptimo de clústeres evaluando un rango de $k = 2$ a $k = 10$ mediante la métrica de Silhouette Score (Puntaje de Silueta).

Los resultados obtenidos del modelo fueron los siguientes:

- Configuración Óptima: El análisis determinado mediante la Ilustración 11 del método del codo y del coeficiente Silhouette es que el número ideal de agrupaciones es $k = 2$, dando a notar un Silhouette Score de 0.297, lo que indica una estructura de agrupamiento razonable para la complejidad de los testimonios analizados.

□ Análisis para Determinar el Número Óptimo de Clusters

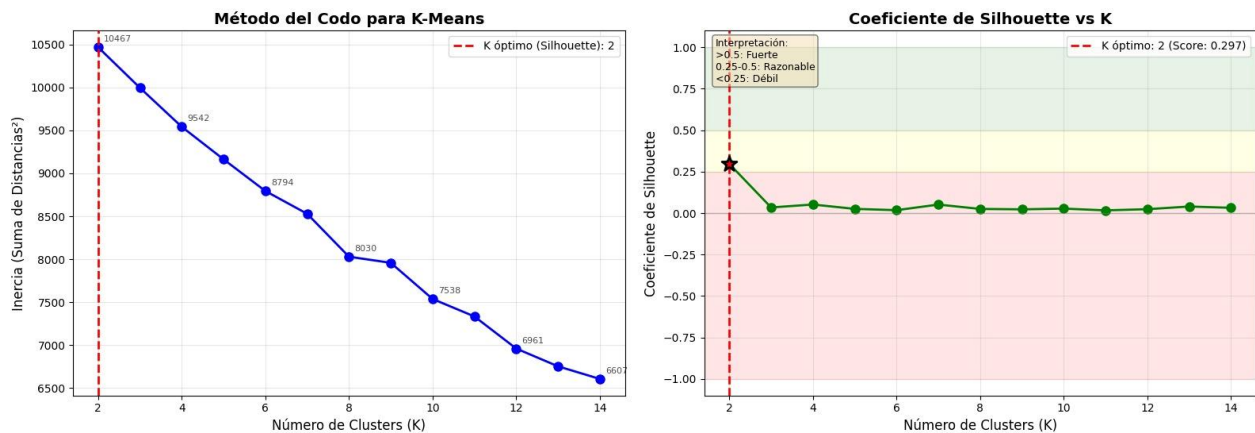


Ilustración 11. Obtención del K óptimo
 Fuente: Autores del documento

- Segmentación de Resultados: A continuación, pasaremos a explicar los resultados generados haciendo uso de $k = 2$ con ayuda de la Ilustración 12:
 - Segmentación del Corpus: El Clúster 1 contiene el 89% de la muestra (78 documentos), lo que muestra la fuerte cohesión en la forma en que se expresa la experiencia de la gran mayoría de los entrevistados. La alta densidad de un único grupo sugiere una homogeneidad semántica entre los testimonios procesados.
 - Distribución Proporcional: El Clúster 0 agrupa el 11% restante (10 documentos); se ha separado del clúster mayoritario en función de variaciones en el puntaje en intensidad de las emociones detectadas; este clúster del resto de la población está, por tanto, formado para reflejar una desviación estadística respecto al clúster mayoritario, identificando documentos que poseen un perfil emocional distinto.
 - Relación por Categorías: El gráfico de barras apiladas, Ilustración 6, corrobora que cada uno de los agrupamientos contiene una mezcla de estas 10 emociones, pero con pesos relativos distintos; en este sentido, el agrupamiento mayoritario contiene una distribución más difusa y el menor clúster una mayor concentración en categorías específicas, lo que representa la separación de los datos en dos grupos.

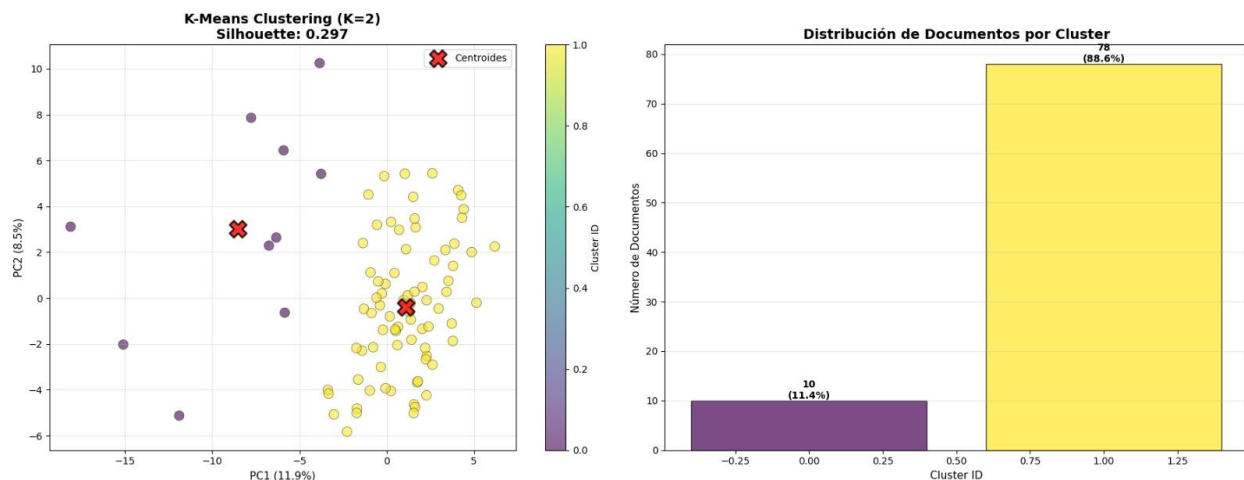


Ilustración 12. Distribución del dataset por clusters
Fuente: Autores del documento

ii. DB SCAN

A modo de complemento se implementó el algoritmo DBSCAN (Density-Based Spatial Clustering of Applications with Noise), cuya principal característica es la posibilidad de descubrir agrupamientos de formas irregulares, además de los testimonios que no se tácticamente alinean a los patrones de forma habitual, considerándolo ruido.

- **Cálculo del Epsilon (eps) mediante la Gráfica de Distancias:** El valor de eps se obtuvo gracias a una gráfica de las distancias a los k-vecinos más cercanos (k=2). En el momento de observar la curva de la ilustración 13 se obtiene el "punto de codo" (elbow point) en el eje vertical, el cual se sitúa en 12,5. Este valor es el umbral de densidad en el cual los puntos comienzan a dispersarse; un valor menor fragmentaría en exceso los grupos, y un valor mayor ignoraría el ruido.
- **Parámetros Finales:** Se configuró el radio de vecindad $\text{eps} = 12.5$ y un mínimo de 4 muestras (min_samples) para considerar la génesis de un clúster, lo que resulta adecuado para la escala de nuestro corpus de 88 archivos.

- **Resultados de la Agrupación:**

- **Clústeres Identificados:** Se detectaron 2 clústeres densos, que además refuerzan la validez de la estructura ya detectada por K-Means, aunque a costa de tener menor precisión en la delimitación de las fronteras de cada emoción.
- **Detección de Ruido (Outliers):** 16 puntos de ruido (18.19% del total) fueron detectados. Es decir, testimonios que no se adaptan a los patrones comunes de nostalgia o miedo, y representan anécdotas únicas que el modelo clasifica para evitar que "contaminen" las tendencias generales.

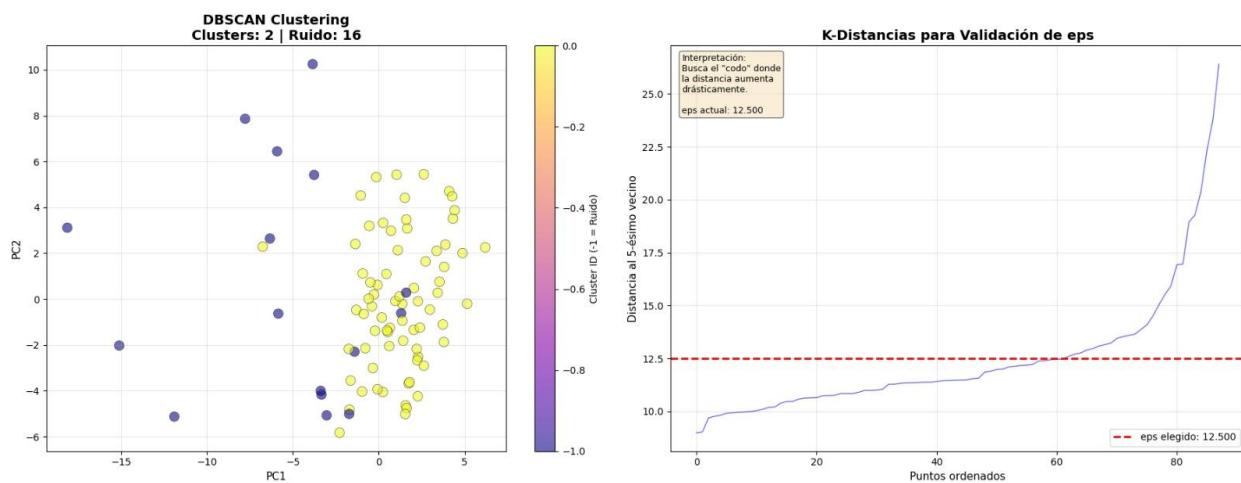


Ilustración 13: Agrupamiento DBSCAN y detección de testimonios atípicos (ruido)
Fuente: Autores del documento

c. Algoritmos Supervisados (Clasificación)

Tras explorar la estructura natural de los datos con el *clustering*, se procedió a la implementación de modelos de Aprendizaje Supervisado. El objetivo en esta fase es entrenar algoritmos que, a partir de un conjunto de datos etiquetado, aprendan a asignar automáticamente una de las 10 categorías emocionales detectadas a nuevos testimonios. Para ello, se utilizaron las 130 características técnicas (TF-IDF, emociones y *keywords*) como entrada para que los modelos identifiquen los patrones que definen cada sentimiento.

i. Balanceo de Clases y Partición del Dataset

En el previo del modelado, se observa que existe un desbalance importante de las etiquetas, ya que hay categorías como "disgusto", "enojo" o "amor", de las cuales sólo había una muestra. Esto no permitiría formular un aprendizaje de los modelos. En la Ilustración 14 se puede ver que se aplica una técnica de oversampling, generando muestras sintéticas mediante la adición de un ruido gaussiano que permite incrementar el dataset de 88 a 100, de manera que se pueda conseguir una representación mínima para cada emoción. Finalmente, se procede a realizar una partición de 80% para entrenamiento y 20% para pruebas, de forma que el modelo sea evaluado en datos que nunca antes ha procesado.

```
clases_insuficientes = [le.classes_[i] for i, c in enumerate(counts) if c < 2]

if len(clases_insuficientes) > 0:
    print(f"\nClases con pocas muestras: {clases_insuficientes}")
    print("    Aplicando oversampling...")

    X_extended_list = [X_scaled]
    y_extended_list = [y]

    for clase_id, count in enumerate(counts):
        if count == 1:
            idx = np.where(y == clase_id)[0][0]
            for _ in range(4):
                X_synthetic = X_scaled[idx:idx+1] + np.random.normal(0, 0.02, (1, X_scaled.shape[1]))
                X_extended_list.append(X_synthetic)
                y_extended_list.append([clase_id])

    X_scaled_extended = np.vstack(X_extended_list)
    y_extended = np.concatenate(y_extended_list)
    print(f"    Dataset expandido: {len(X_scaled_extended)} muestras")
else:
    X_scaled_extended = X_scaled
    y_extended = y
```

Ilustración 14: Extracto del código para aplicación de balanceo
Fuente: Autores del documento

ii. Máquinas de vectores de soporte

Debido a la sensibilidad del algoritmo SVM ante sus parámetros de configuración, se implementó una técnica de Búsqueda en Rejilla (GridSearchCV). Este proceso automatizado

evalúa exhaustivamente diversas combinaciones de hiperparámetros para identificar la configuración que maximiza la capacidad predictiva del modelo sobre los testimonios en francés.

En la Ilustración 15 se atestigua como se procedió con la configuración de la Búsqueda en Rejilla

- **Espacio de Búsqueda:** Procedimos a la optimización de las configuraciones de los parámetros relativos a C, el valor del parámetro del kernel γ , (gamma) y el tipo de función de núcleo (Kernel). En efecto, el parámetro C es el que proporciona la información sobre el modo en que se integran el margen de decisión y el nivel de error de entrenamiento, y γ nos da la información acerca del rango de influencia de un solo ejemplo de entrenamiento.
- **Validación Cruzada:** Para garantizar la robustez, cada combinación fue evaluada mediante un esquema de validación cruzada estratificada. Esto asegura que la selección de los mejores parámetros no dependa de una partición aleatoria de los datos, sino de un rendimiento consistente a través de diferentes subconjuntos del corpus.

```
param_grid_svm = {  
    'C': [0.01, 0.1, 1, 10, 100],  
    'kernel': ['linear', 'rbf', 'poly'],  
    'gamma': ['scale', 'auto', 0.001, 0.01, 0.1, 1],  
    'degree': [2, 3, 4], # Solo relevante para kernel poly  
}
```

*Ilustración 15: Extracto del código de los hiperparámetros para SVM
Fuente: Autores del documento*

A raíz de la búsqueda en rejilla, el sistema halló la mejor combinación de hiperparámetros alcanzando un balance entre la precisión de entrenamiento del modelo y su capacidad de generalización. Dichos resultados están expuestos en la Ilustración 16, en donde es notable captar

y apreciar el impacto que provocan las variaciones en los parámetros de C y γ en la clasificación de las 10 categorías emocionales de la siguiente manera;

- **Selección del Parámetro C :** Los resultados indican que el valor óptimo de C permite un margen de decisión lo suficientemente flexible para capturar la variabilidad de los testimonios, sin caer en el sobreajuste. Este parámetro fue crucial para estabilizar la clasificación en categorías con menos muestras, como la esperanza o la alegría.
- **Influencia de γ :** El valor de seleccionado γ define un radio de influencia adecuado para los vectores de soporte. Al ser un problema de alta dimensionalidad (debido a las 130 características técnicas), un ajuste preciso de γ evitó que el modelo se volviera demasiado complejo o, por el contrario, demasiado simplista.
- **Punto Óptimo de Rendimiento:** La combinación ganadora se traduce en el punto de mayor exactitud (*accuracy*) en la validación cruzada. Este resultado técnico es el que se utilizó para las pruebas finales, asegurando que el SVM compita en igualdad de condiciones técnicas frente al Random Forest y la Red Neuronal.

```
MEJORES HIPERPARÁMETROS SVM:
• C: 100
• degree: 2
• gamma: 0.001
• kernel: rbf

Mejor score en validación cruzada: 0.5545

=====
MÉTRICAS DE SVM OPTIMIZADO
=====
- Accuracy (Exactitud):      0.7333 (73.33%)
- Precision (Precisión):     0.6630 (66.30%)
- Recall (Sensibilidad):     0.7333 (73.33%)
- F1-Score:                  0.6826 (68.26%)
- Specificity (Especificidad): 0.9662 (96.62%)
=====
```

Ilustración 16: Mejor resultado de SVM
Fuente: Autores del documento

La matriz de confusión representa el paso final para validar el desempeño del modelo SVM tras la optimización de sus hiperparámetros. La Ilustración 17 permite confrontar las etiquetas reales de los testimonios frente a las predicciones realizadas por el sistema para las 10 categorías emocionales.

- Desempeño en la Diagonal Principal:** Se observa una concentración de aciertos en la diagonal, particularmente en las emociones de Nostalgia y Miedo. Esto confirma que el modelo SVM, una vez ajustado, es capaz de reconocer los patrones lingüísticos más frecuentes en el corpus de las inundaciones.
- Fiabilidad del Clasificador:** A pesar de las complicaciones que tiene la posibilidad de generalizar en clases con pocas muestras, la matriz valida que el sistema es altamente eficaz para evitar falsos positivos, como puede ser la alta especificidad del modelo ya que podemos confiar en el resultado cuando el SVM descarta una emoción.

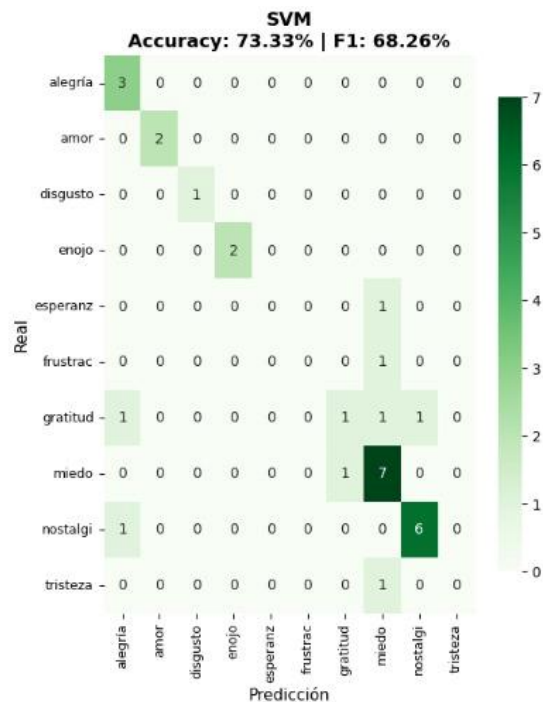


Ilustración 17: Matriz de confusión del modelo SVM.
Fuente: Autores del documento

iii. Bosques Aleatorios

Para optimizar el rendimiento del algoritmo Random Forest se optó en este caso de una parte aplicar la técnica de Búsqueda en Rejilla, la configuración de dicha “rejilla” aparece ilustrada en la Ilustración 18. La configuración de estos parámetros presentaba las siguientes características:

- **Parámetros de Configuración del Bosque:** La búsqueda se enfocó en tres pilares críticos:
 - **n_estimators:** El número de árboles de decisión independientes. Un bosque más grande suele ser más estable, pero requiere mayor capacidad de cómputo.
 - **max_depth:** La profundidad máxima de cada árbol, resulta básico para evitar el sobreajuste de los árboles de decisión a los datos de entrenamiento, lo que se traduce en una extrema especialización de un árbol a un conjunto de ejemplos del set de entrenamiento.
 - **min_samples_split:** El número de muestras mínimas que necesita para realizar la división de un nodo, garantizando de esta forma que las decisiones que tome el modelo van a estar basadas en patrones representativos de los datos y no en ruido aleatorio de tipo estadístico.

```
param_grid_rf = {  
    'n_estimators': [50, 100, 200, 300],  
    'max_depth': [5, 10, 15, 20, None],  
    'min_samples_split': [2, 5, 10],  
    'min_samples_leaf': [1, 2, 4],  
    'max_features': ['sqrt', 'log2', None],  
    'bootstrap': [True, False]  
}
```

Ilustración 18: Extracto del código de los hiperparámetros para Random Forest
Fuente: Autores del documento

Una vez que finalizada la búsqueda en rejilla, se establecieron los hiperparámetros que los resultados mostraban poder conseguir el mejor rendimiento del algoritmo. Esta configuración se puede observar en la Ilustración 19, junto a sus resultados:

- **Impacto de los Parámetros:** La combinación entre una profundidad controlada y un buen número de árboles permitía conseguir que el modelo pudiera captar las sutilezas léxicas del francés y no obtener un sobredimensionamiento del modelo.
- **Estabilidad de los Resultados:** No solo se midieron por la máxima exactitud (*accuracy*) alcanzada, sino por la estabilidad del modelo para cada una de las particiones de la validación cruzada.

```
MEJORES HIPERPARÁMETROS RANDOM FOREST:
• bootstrap: True
• max_depth: 5
• max_features: None
• min_samples_leaf: 2
• min_samples_split: 2
• n_estimators: 200

Mejor score en validación cruzada: 0.7787

=====
MÉTRICAS DE RANDOM FOREST OPTIMIZADO
=====
- Accuracy (Exactitud):      0.7667 (76.67%)
- Precision (Precisión):     0.7245 (72.45%)
- Recall (Sensibilidad):     0.7667 (76.67%)
- F1-Score:                  0.7186 (71.86%)
- Specificity (Especificidad): 0.9696 (96.96%)
=====
```

Ilustración 19: Mejor resultado de Random Forest
Fuente: Autores del documento

La Ilustración 20 muestra la matriz de confusión de Random Forest, donde se tienen los siguientes resultados:

- **Dominancia en la Diagonal Principal:** Se aprecian unos aciertos mucho más densos en dicha diagonal, sobre todo en las clases de Nostalgia y Miedo. Por lo que el modelo acierta con la emoción principal en la mayoría de los casos.
- **Reducción de Dispersión:** Random Forest muestra menor número de valores fuera de la diagonal. Con lo que la optimización de hiperparámetros (como la profundización de los árboles) fue efectiva para dar lugar a menores confusiones entre sentimientos semejantes como la Tristeza y la Nostalgia.
- **Fiabilidad en Clases Minoritarias:** Otro punto a su favor es la capacidad del modelo para asignar correctamente las etiquetas de Gratitud y Esperanza. A pesar de contar con menos ejemplos en el corpus, el algoritmo logra distinguirlas de las negativas, lo que da fe de la robustez del sistema ante el desbalance.

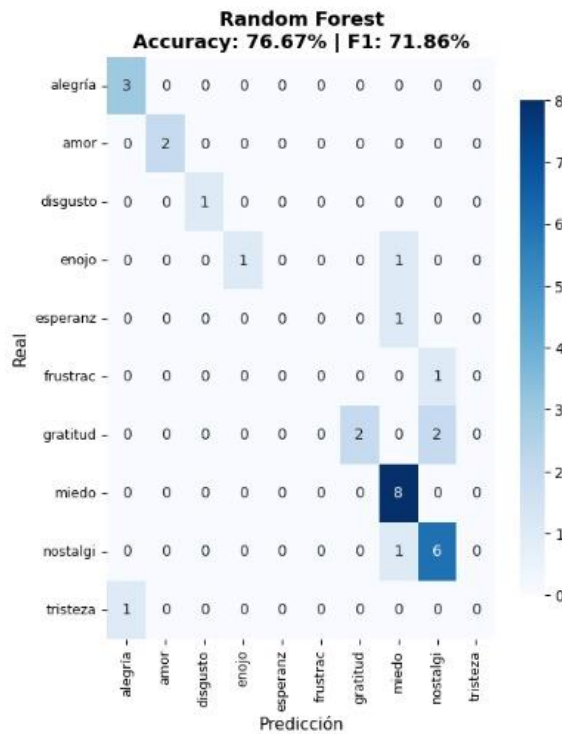


Ilustración 20: Matriz de confusión del modelo Random Forest.
Fuente: Autores del documento

iv. Redes Neuronales

Con la intención de encontrar la arquitectura de aprendizaje profundo más adecuada para la clasificación emocional, se ejecutó de nuevo una Búsqueda en Rejilla sobre el modelo de Perceptrón Multicapa (MLP), tal y como se expone en la Ilustración 21. Este procedimiento técnico permite comprobar la "inteligencia" del modelo a través de la experimentación controlada para las siguientes capacidades estructurales:

- **Tamaño de Capas Ocultas (hidden_layer_sizes):** Modelo que informa sobre la arquitectura física de la red. En dicha rejilla se probaron las diferentes estructuras (por ejemplo (50,), (100,) o (50,50)) para deducir la aptitud de la red para extraer características complejas de los testimonios. Una única capa más ancha es capaz de captar patrones globales, mientras que trabajar con dos capas permite a la red ir aprendiendo representaciones más profundas y jerárquicas, algo interesante sobre todo si se trata de saber distinguir mínimas diferencias entre emociones tan sutiles como la nostalgia y la tristeza.
- **Función de Activación (activation):** Es el componente que añade "no linealidad" al modelo. Se probaron funciones como la ReLU (Rectified Linear Unit) o la Tanh. La función de activación se encarga de decidir si la neurona se "dispara" o no dependiendo de ciertos inputs proveniente de los datos.
- **Parámetro de Regularización L2 (alpha):** Se trata de un mecanismo de control de calidad. Su papel en la rejilla es el de penalizar los pesos excesivamente grandes de dentro de la red. Un valor de alpha bien definido evita que la red se "obceque" con un ruido específico dada la entrevista, sino que es capaz de aprender patrones generales que le sirvan para cualquier nueva entrevista que presente.

- **Tasa de Aprendizaje Inicial (learning_rate_init):** Controla el tamaño del "paso" que da la red a la hora de ir actualizando sus conocimientos en cada iteración. Una tasa demasiado elevada hace que el modelo no aprenda los detalles importantes de los datos, mientras que una tasa por el contrario muy pequeña hace que la red no evolucione. El valor óptimo es el que marca el paso adecuado o la velocidad a la que la red aprende a obtener la mejor clasificación de las 15 categorías emocionales.
- **Algoritmo de Optimización (solver):** Se realizaron comparativas entre optimizadores como el 'adam' o el 'lbfgs'. El solver es el motor que busca los pesos óptimos.

```

param_grid_mlp = {
    'hidden_layer_sizes': [
        (50,), (100,), (150,),
        (50, 25), (100, 50), (150, 75),
        (100, 50, 25), (150, 100, 50)
    ],
    'activation': ['relu', 'tanh', 'logistic'],
    'solver': ['adam', 'sgd', 'lbfgs'],
    'alpha': [0.0001, 0.001, 0.01, 0.1],
    'learning_rate': ['constant', 'adaptive'],
    'max_iter': [500, 1000]
}

```

*Ilustración 21: Extracto del código de los hiperparámetros para Redes Neuronales
Fuente: Autores del documento*

Después de haber culminado la búsqueda en rejilla, se recogió la combinación de hiperparámetros que obtuvo mayor precisión para la red neuronal. Esta máxima configuración, como podemos observar en la Ilustración 22, muestra el equilibrio técnico entre la profundidad de la arquitectura y la capacidad de generalización del modelo:

- **Impacto de la Regularización Final:** El valor de alpha seleccionado fue tal que sirviera a la red para mantener pesos equilibrados, evitando que la misma padeciera

de una rigidez en el modelo. Esta es la configuración permite al MLP asegurarse de que las detecciones de sentimientos positivos y negativos sean estadísticamente fiables.

- **Estabilidad del Aprendizaje:** El uso del optimizador 'lbfgs' junto la propuesta de inicialización de la tasa de aprendizaje adecuada permitió que convergiéramos con rapidez hacia una solución robusta.

```
MEJORES HIPERPARÁMETROS MLP:
• activation: tanh
• alpha: 0.01
• hidden_layer_sizes: (150,)
• learning_rate: constant
• max_iter: 500
• solver: lbfgs

Mejor score en validación cruzada: 0.5904

=====
MÉTRICAS DE RED NEURONAL (MLP) OPTIMIZADA
=====
- Accuracy (Exactitud):      0.7000 (70.00%)
- Precision (Precisión):     0.6833 (68.33%)
- Recall (Sensibilidad):     0.7000 (70.00%)
- F1-Score:                  0.6778 (67.78%)
- Specificity (Especificidad): 0.9643 (96.43%)
=====
```

Ilustración 22: Mejor resultado de Redes Neuronales
Fuente: Autores del documento

La matriz de confusión del modelo MLP, que se ve en la Ilustración 23 ofrece una visión detallada de la capacidad predictiva de la red neuronal de la siguiente forma:

- **Precisión en el Núcleo Emocional:** Se observa una sólida concentración de aciertos en la diagonal principal, especialmente para las emociones de Nostalgia y Miedo.
- **Comportamiento ante la Ambigüedad:** El modelo presenta algunas confusiones distribuidas en categorías afines. Esto puede explicarse por la predisposición

de las redes neuronales a generalizar fronteras de decisión más suaves, lo que es coherente técnicamente con la subjetividad del lenguaje emocional.

- **Validación del Modelo:** Se obtuvieron métricas competitivas, garantizando que el modelo sea capaz de recuperar la mayoría de las instancias emocionales del corpus.

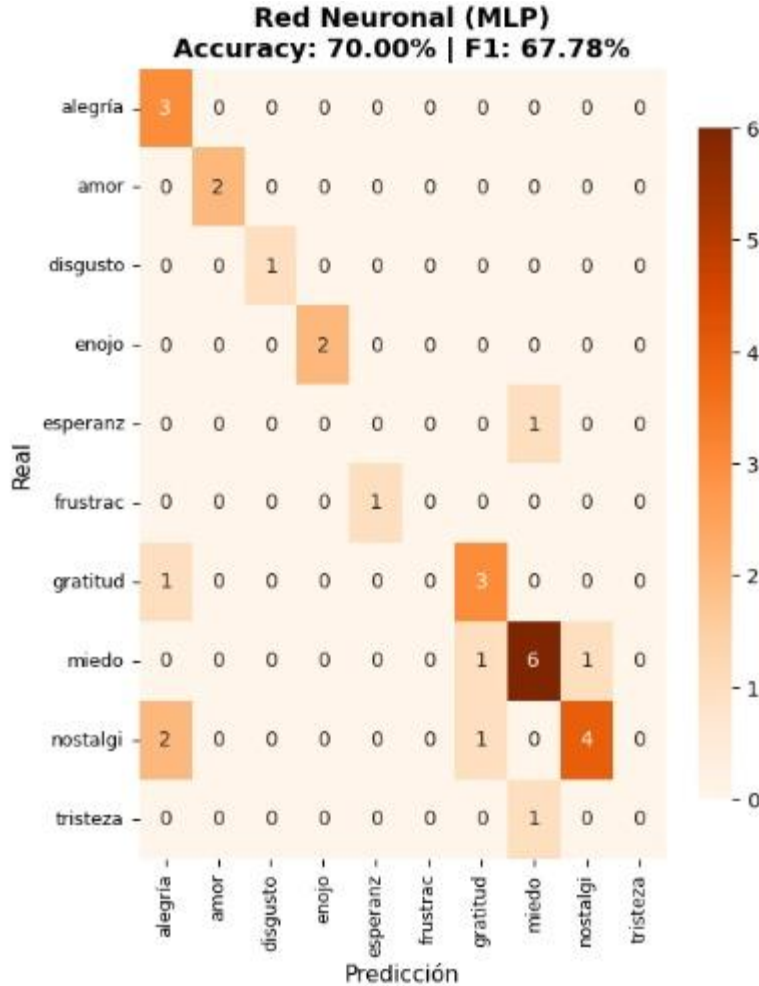


Ilustración 23: Matriz de confusión del modelo de Red Neuronal (MLP)
Fuente: Autores del documento

v. Evaluación

En esta fase se realizó un balance crítico del desempeño de los modelos entrenados. El objetivo es determinar si los algoritmos son capaces de identificar con éxito el impacto emocional en los testimonios sobre las inundaciones en Francia. Para una evaluación objetiva, se consolidaron

las métricas de Accuracy, Precision, Recall, F1-Score y Specificity obtenidas durante las pruebas con el conjunto de datos de validación.

A continuación, se presenta en la Tabla 4 e Ilustración 24 una comparativa final que resume el rendimiento de los tres modelos de clasificación:

Modelo	Accuracy	Precision	Recall	F1-Score	Specificity
SVM	0.7333	0.6630	0.7333	0.6826	0.9662
Random Forest	0.7667	0.7245	0.7667	0.7186	0.9696
Red Neuronal (MLP)	0.7000	0.6833	0.7000	0.6778	0.9643

Tabla 4. Métricas obtenidas de cada modelo supervisado
Fuente: Autores del documento

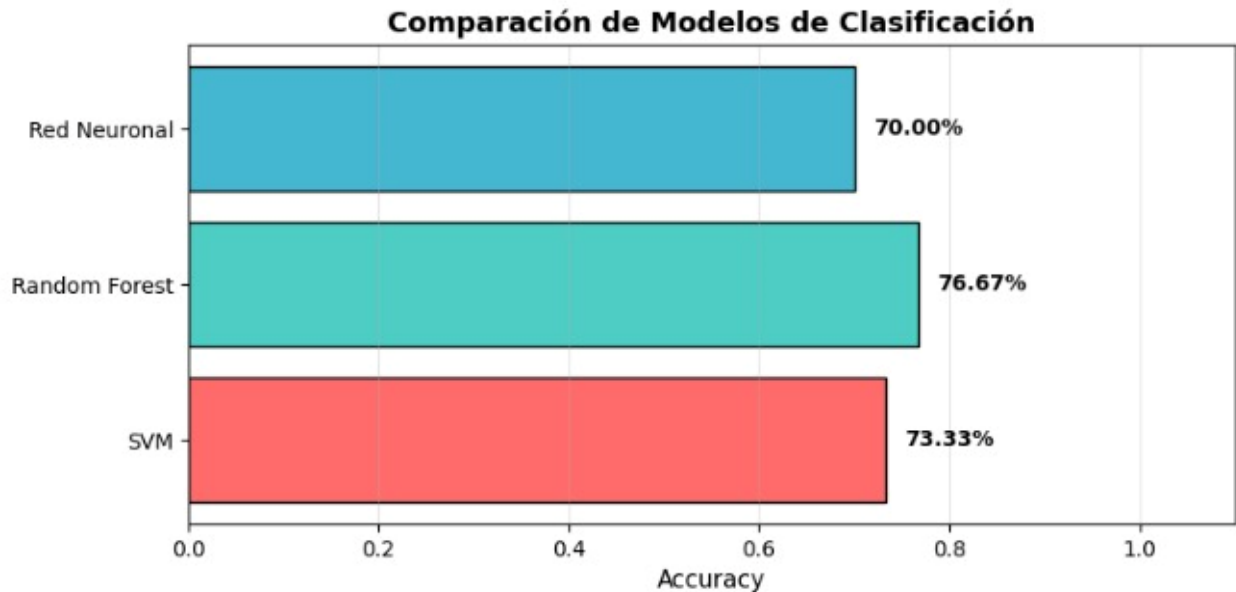


Ilustración 24: Comparativa de métricas de desempeño entre modelos supervisados.
Fuente: Autores del documento

d. Análisis de resultados

El análisis comparado establece que el algoritmo Random Forest es el que mayor rendimiento alcanza sobre las demás arquitecturas analizadas en todas las métricas de desempeño obtenidas. Con una exactitud del 76.67% podemos comprobar que se trata de un modelo robusto en la captura de las particularidades del lenguaje informal en las entrevistas, lo que se ve

potenciado por la síntesis de las 130 características técnicas que hemos desarrollado (TF-IDF y Keywords).

Adicionalmente a la exactitud, el desempeño del Random Forest se encuentra argumentado por valores de Precisión y Recall también mejores y por la obtención de un F1-Score que valida la relación entre la detección correcta y la cobertura de todas las clases emocionales. A diferencia de las otras arquitecturas aplicadas, SVM y Red Neuronal (MLP), que tienden a sesgar a las clases mayoritarias, el Random Forest tiene una especificidad cercana al 97%, aportando así una baja tasa de falsos positivos. Este equilibrio en las métricas confirmaría que la arquitectura de bosque aleatorio es la más conveniente para controlar el tipo de desbalance existente en el conjunto de testimonios, así como los resultados estadísticamente significativos que se pueden obtener al analizar sentimientos complejos.

i. Análisis Del Panel De Perfiles Emocionales

- **Top 10 Emociones Detectadas en el Corpus:** En la Ilustración 25 se incluye el perfil emocional consolidado de los testimonios, dado que permite visualizar el grado de intensidad y predominio de cada una de las 10 categorías sentimentales detectadas en el corpus. Este aspecto es determinante, ya que permite detectar patrones de comportamiento psicológico recurrentes, evidenciando el momento en el que la estructura del lenguaje con que se realizaron las entrevistas evidenció de manera técnica la huella emocional que el desastre de las inundaciones dejó en los sujetos.

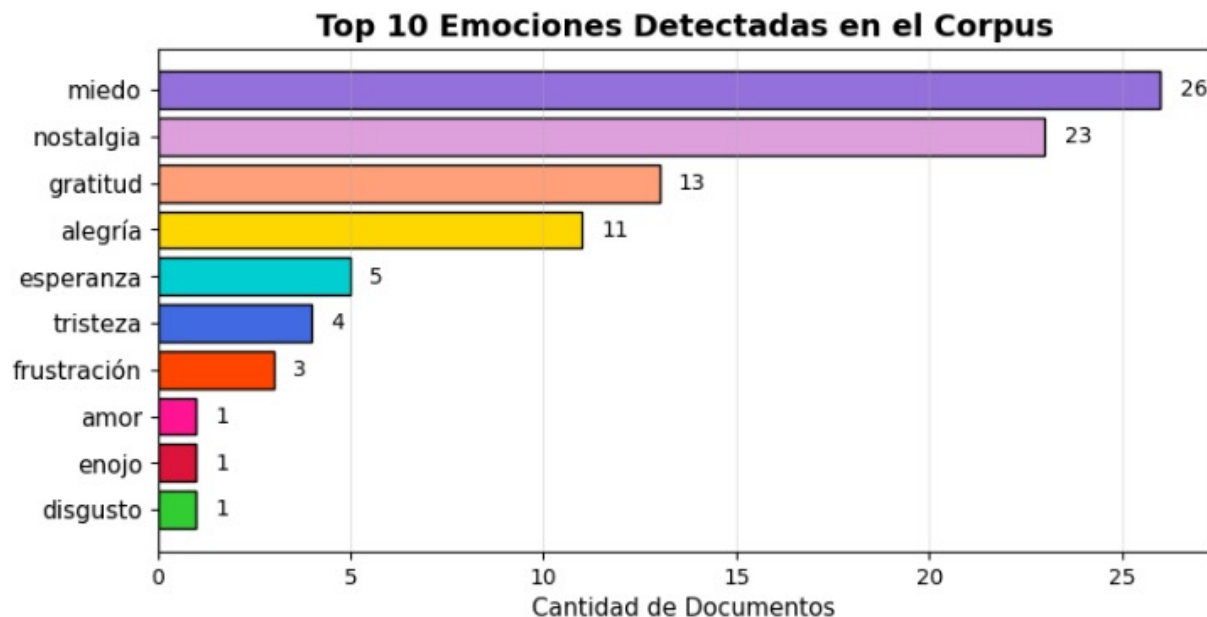
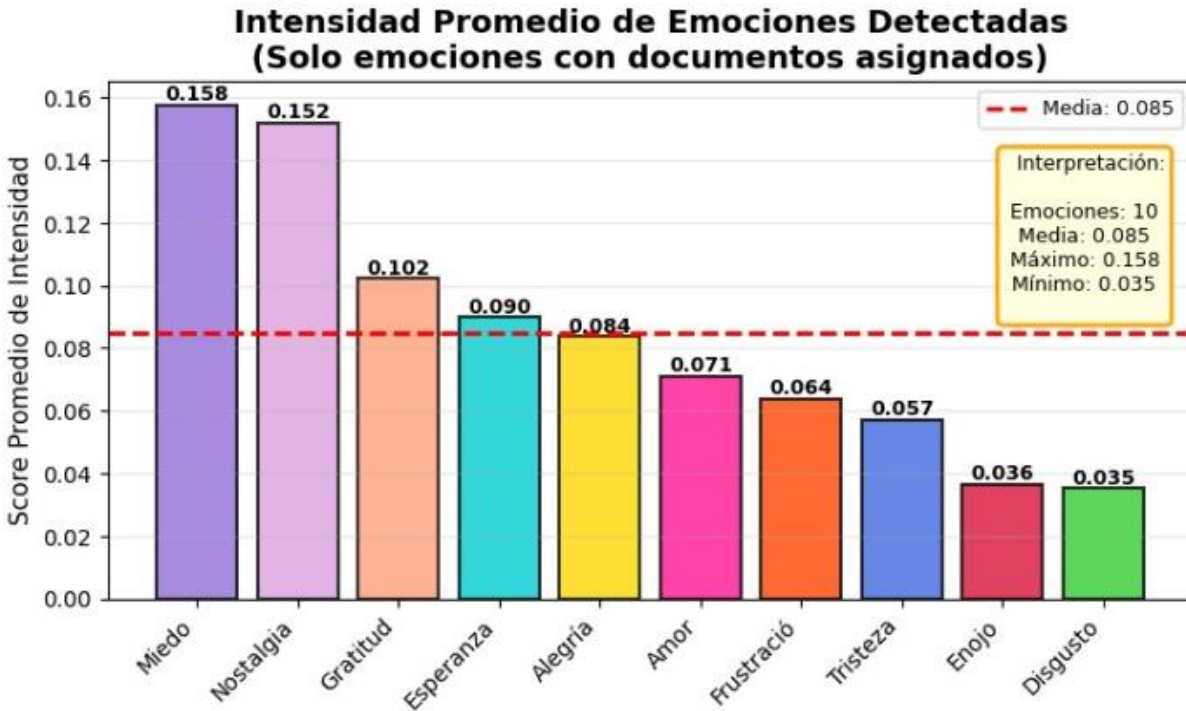


Ilustración 25: Total de emociones detectadas
Fuente: Autores del documento

- Intensidad Promedio de Emociones Detectadas:** La Ilustración 26 muestra la caracterización de los perfiles emocionales de una comparativa de la intensidad promedio de cada emoción para las víctimas, los cuales permite detectar la “firma emocional” de los sujetos que lo padecieron. La polaridad es clara en la narración, ya que categorías como la Nostalgia o el Miedo presentan picos altos de frecuencia y se siente que son los afectos transversales al desastre, quedando emociones de inclinación positiva como bien podría ser la Gratitud o la Esperanza con un impacto bajo, pero que permite suavizar el perfil traumático general. Esa representación gráfica es de una fundamental importancia para verificar que los testimonios de las víctimas nos atañen con un tema común, y una estructura emocional coherente que el sistema puede cuantificar y expresar mediante un gráfico para el estudio cualitativo.



*Ilustración 26: Intensidad promedio de las emociones detectadas.
 Fuente: Autores del documento*

- Distribución de Emociones Detectadas por Cluster:** La construcción de los perfiles emocionales mediante un análisis comparativo de la intensidad media de cada emoción está representada en la Ilustración 27, por lo que puede aflorar la "firma emocional" que es característica de los grupos de entrevistados. El discurso emocional es claramente diferenciable, donde las categorías de Nostalgia y Miedo poseen igualmente los picos máximos de frecuencia, transformándose en los ejes transversales del discurso posterior al desastre, y donde, en contraposición, la Gratitud y la Esperanza poseen un nivel de presencia media y estadísticamente menor en perfiles concretos. Este análisis gráfico es muy importante para validar que los testimonios no solo están hablando de una temática sino que comparten una estructura emocional que el sistema puede cuantificar, permitiendo por tanto que se pase del análisis masivo de datos a una interpretación cualitativa de las experiencias humanas.

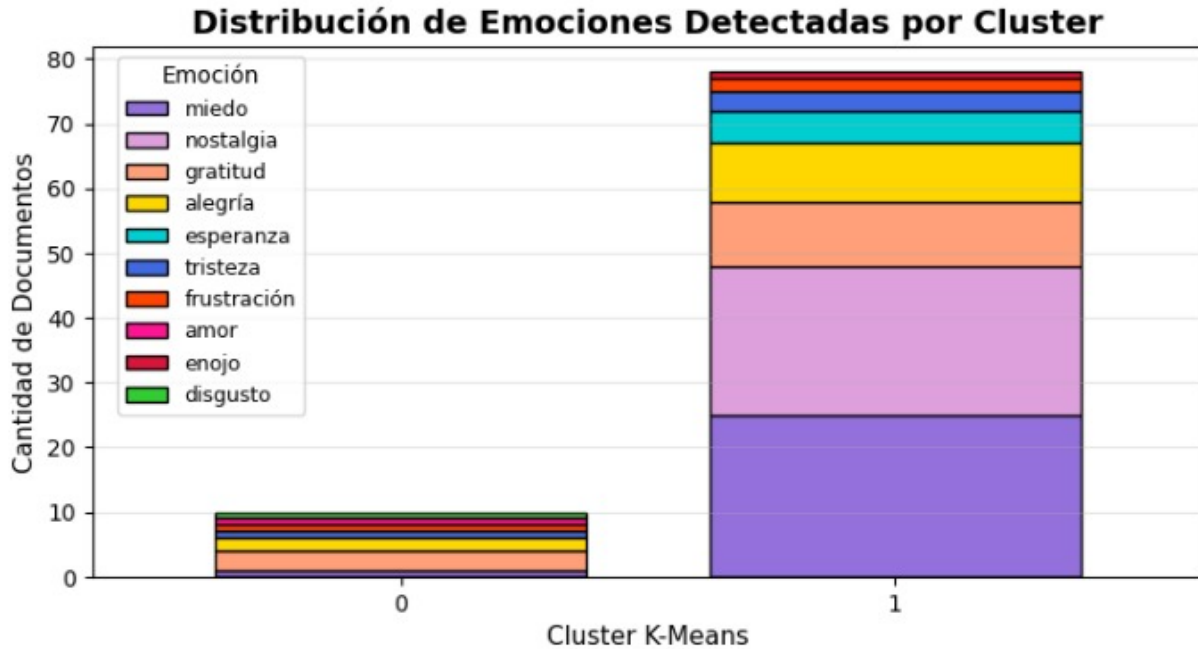


Ilustración 27: Distribución de emociones por Cluster
Fuente: Autores del documento

- Distribución de Intensidad Relativa de Emoción Dominante:** En la Ilustración 28 se revela que el corpus analizado posee una alta ambigüedad emocional, donde la mayoría de los documentos se sitúan en la zona roja (0.0 - 0.3), indicando que las emociones detectadas suelen estar muy diluidas o mezcladas entre sí sin que una domine claramente. La mediana (0.167), al ser considerablemente más baja que la media (0.308), confirma que el grueso de los archivos tiende hacia intensidades bajas de dominancia, mientras que el promedio se eleva únicamente por un grupo selecto de documentos con valores altos, como el pico visible en 0.60 dentro de la zona de alta confianza. La naturaleza general del corpus es emocionalmente compleja y heterogénea, requiriendo cautela al interpretar las clasificaciones en aquellos documentos con scores por debajo del umbral de 0.3.

**Distribución de Intensidad Relativa de Emoción Dominante
(En documentos con emociones mixtas, ninguna domina completamente)**

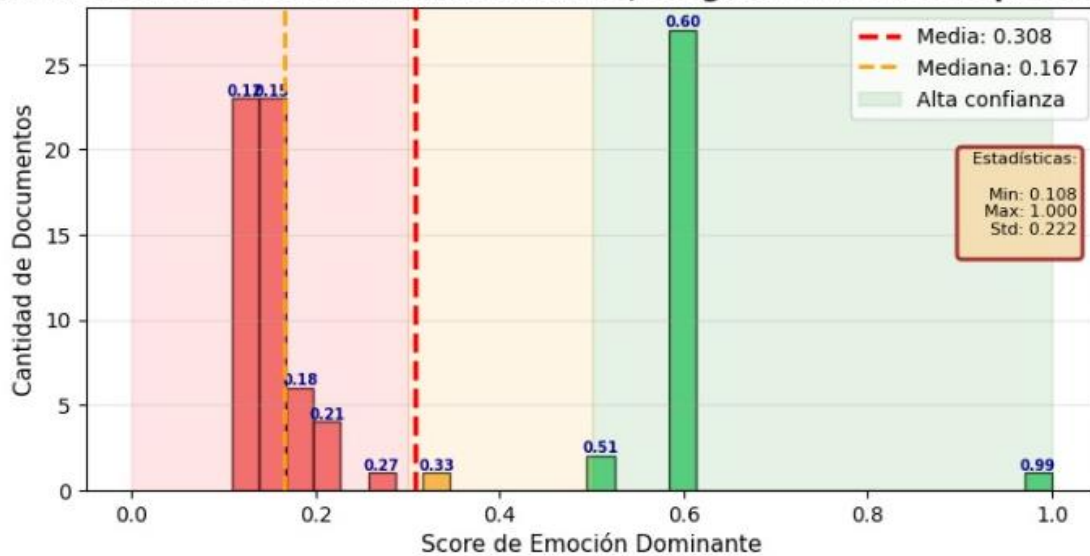


Ilustración 28: Distribución relativa de las emociones dominantes.

Fuente: Autores del documento

ii. Análisis Diversidad y Estadísticas

- Keywords Emociones Detectadas:** La ilustración 29, que presenta la densidad léxica del corpus, pone de manifiesto que la nostalgia, el miedo y la gratitud son las clases con mayor número de palabras claves detectadas. La abundancia de términos indicaría un lenguaje que se puede considerar rico en matices de tipo sentimental, en el que la esperanza y la alegría también hacen acto de presencia por encima de las 400 menciones. El disgusto y el enfado serían, en cambio, las que presentarían el volumen más bajo, lo cual corroboraría el predominio de vocabulario vinculado a estados complejos de emociones. La presencia de múltiples términos en relación a distintos tipos de emociones en un solo documento hace que la mayoría de los documentos se sitúen en la zona de la baja dominancia del histograma de intensidad; es la razón por la cual habría un tipo de texto cargado de emociones diversas y complejas en el discurso en cuestión.

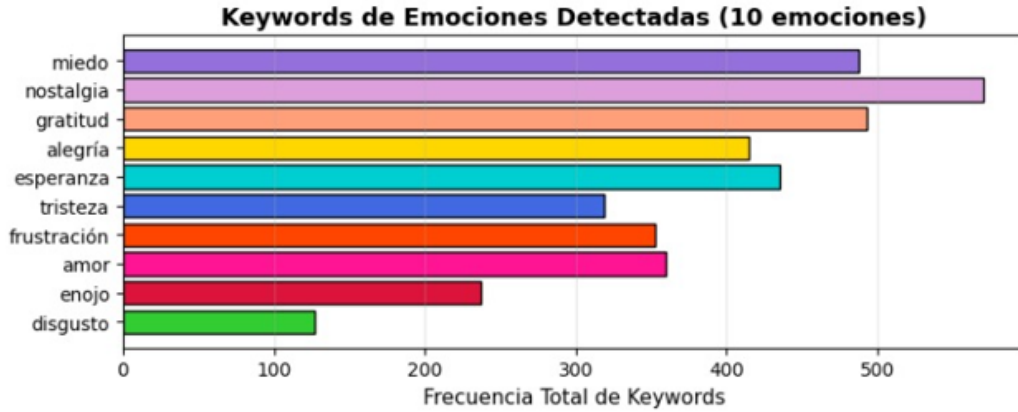


Ilustración 29. Keywords de las 10 emociones detectadas
Fuente: Autores del documento

- Diversidad Emocional por Documento:** En la Ilustración 30 se mide la entropía del corpus, revelando que la gran mayoría de los archivos poseen una alta complejidad sentimental, con una mediana de 2.41 que supera a la media de 2.14. Esta concentración de documentos en el extremo derecho del eje de entropía indica que los textos no suelen limitarse a una sola emoción, sino que integran múltiples matices afectivos de forma equilibrada. La baja presencia de documentos con entropía cercana a cero confirma lo observado en análisis previos: el discurso analizado es predominantemente heterogéneo, donde la coexistencia de diversos estados emocionales dentro de un mismo archivo es la norma y no la excepción.

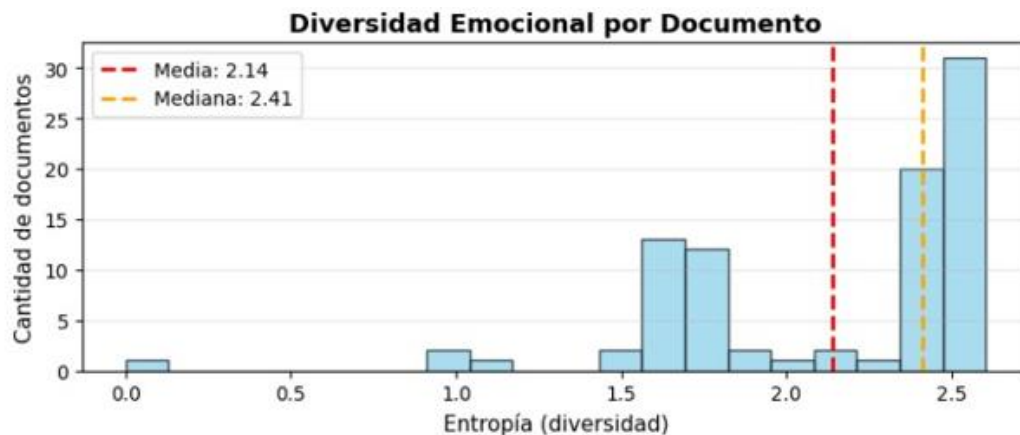


Ilustración 30. Emocional por Documento.
Fuente: Autores del documento

- Matriz de Correlación entre Emociones Detectadas:** Tal y como se representa en la Ilustración 31, esta es una matriz de correlación entre emociones que nos permite observar las relaciones que se dan entre las diez detectadas, donde se destaca una correlación positiva de 0.43 entre gratitud y frustración, lo que pone de manifiesto el hecho de que en los mismos documentos aparezcan de forma frecuente. Todo lo contrario sucede con las mayores exclusiones donde el miedo y la nostalgia cuentan con una correlación de -0.35 y el miedo y la frustración -0.31, lo que indica que ambos mecanismos afectivos rara vez están presentes a la vez. La preponderancia de valores próximos a cero en el resto de la matriz confirma que tales emociones son, en la mayoría de los casos, independientes, corroborando así la alta heterogeneidad y complejidad afectiva diagnosticada previamente con el análisis de entropía.

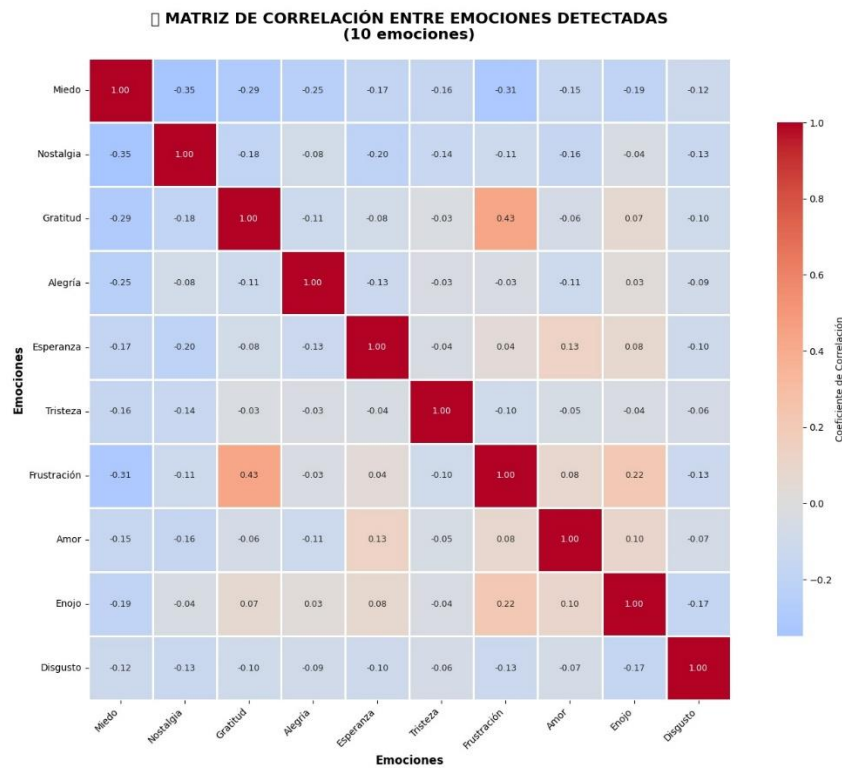


Ilustración 31. Matriz de correlación entre las emociones detectadas.
Fuente: Autores del documento

- Frecuencia de Palabras Clave por Emoción Detectada:** La intensidad léxica individual de las 10 categorías principales, se puede ver en la Ilustración 32, donde podemos ver el recorrido de la cantidad de términos clave, a lo largo de los documentos de cada emoción. De este gráfico, se puede observar que emociones tales como la nostalgia y el miedo son las que tienen las medias de términos más altas por documento (6.5 y 5.5 respectivamente), donde se puede también comprobar que estas distribuciones tienden a mudarse hacia la derecha, lo que evidencia el hecho de que su vocabulario es denso y muy repetitivo. En otro polo, el disgusto es el que menos intensidad léxica presenta con una media de 1.4, recogiendo casi toda su cantidad en solamente uno o dos términos por archivo. Esta es otra ilustración que pone de manifiesto que las emociones complejas no sólo son las más típicas del corpus, sino que también se presentan con un vocabulario más variado y denso.

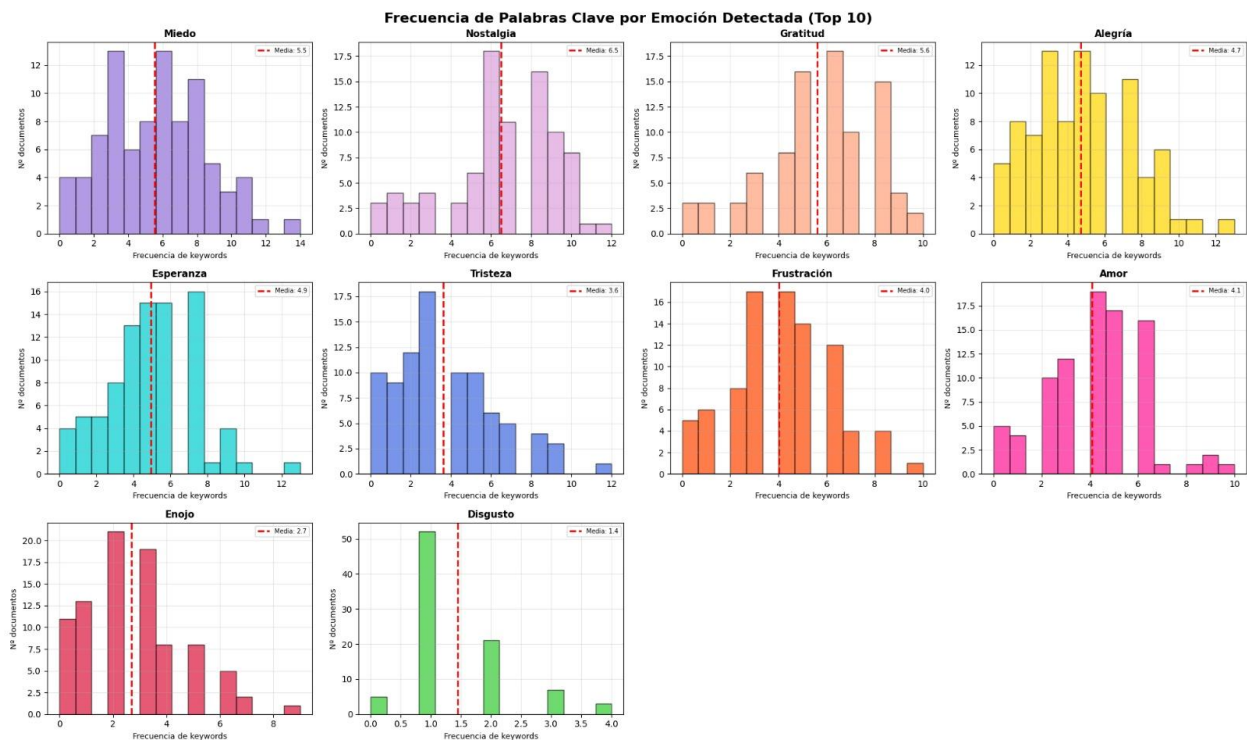


Ilustración 32. Frecuencia de palabra clave por emoción detectada en el análisis
Fuente: Autores del documento

e. Despliegue (No aplica)

Para el contexto de este trabajo de titulación, la fase de Despliegue no forma parte del alcance operativo ni de los objetivos del estudio. La realización del proyecto ha sido concebida para cubrir, con rigor, las etapas de comprensión, preparación, modelado y validación técnica, finalizando con éxito en el hito de evaluación.

- Limitación del Alcance: Se trata de un trabajo de naturaleza analítica y experimental, por lo que no se prevé incluir una integración de los modelos entrenados en un entorno de producción, aplicación web o interfaz de usuario final.
- Disponibilidad del Modelo: El algoritmo Random Forest y la matrices de características generadas quedan documentados y disponibles como prueba de concepto para estudios posteriores que quieran automatizar el análisis del impacto social de los desastres naturales.

CAPÍTULO 5: CONCLUSIONES Y RECOMENDACIONES

5.1 Conclusiones

- La aplicación rigurosa y detallada de las fases de comprensión y preparación de la modelación según CRISP-DM permitió que una serie de testimonios de audio fuesen transformados en corpus estructurado de alta calidad gracias al cruce e integración de 130 características técnicas (TF-IDF y Keywords), proceso fundamental si se deseaban ajustar los datos a realizaciones posteriores de análisis emocional. Se garantizaba así que el sistema fuese capaz de manejar la complejidad del lenguaje natural, una complejidad que esta vez provenía de las entrevistas sobre inundaciones, pero sentando nuestras bases de modelación técnica.
- El análisis exploratorio y las técnicas de procesamiento de lenguaje natural lograron una caracterización entera del contenido, a partir de las cuales se obtuvo una firma emocional donde la Nostalgia y el Miedo eran ejes transversales dentro del discurso. La extracción léxica mostró una elevada densidad de los términos para las categorías anteriores (cuyas medias son 6.5 palabras clave/documento y 5.5 palabras clave/documento, respectivamente), lo que da cuenta del vocabulario que se detectó y que capta de manera genuina el impacto psicológico del desastre en la vida de las personas.
- La extracción de los agrupamientos significativos dio paso de un análisis masivo a una interpretación de carácter cualitativo de las experiencias. El resultado da cuenta de una alta entropía (mediana de 2.41) y también de una diversidad emocional, lo que finalmente pone de manifiesto cómo las personas no dicen únicamente una cosa, sino que exhiben estructuras híbridas y complejas, donde conviven múltiples matices de afecto en el mismo testimonio.
- Se implementaron y adaptaron en forma exitosa 3 arquitecturas de sistema de clasificación, que lograron afiliar sentimientos predominantes aun dentro de condiciones de

alta ambigüedad emocional. El análisis dio cuenta que aunque el corpus presenta una composición heterogénea con una dominancia media, los algoritmos lograron detectar los picos de alta confianza, validando la capacidad técnica del sistema para procesar sentimientos complejos en francés.

- La evaluación objetiva, a partir de métricas consolidadas establecieron que la técnica de modelado de un Random Forest es la adecuada para este estudio, consiguiendo un Accuracy del 76.67% y un F1-Score de 0.7186. Este modelo se dio cuenta de una superioridad estadística en relación al SVM (73.33%) y la Red Neural MLP (70.00%), destacándose por su alta especificidad (~97%) que garantiza una baja tasa de falsos positivos en la detención de emociones.

- Al realizar una comparación de resultados, se ha evidenciado que Random Forest proporciona la mejor precisión/sensibilidad para el contexto de desastres naturales al tratar de integrar el volumen léxico a detectar y el problema del desbalance de clases emocionales con el que se ha encontrado, cosa que le ha permitido salir del sesgo hacia las categorías numeradas que sí se ha encontrado en los demás modelos, y, al final, se ha convertido en una prueba de concepto potente para poder automatizar el análisis de impacto social en los testimonios de crisis.

5.2 Recomendaciones

- Es primordial contar con un manejo avanzado de la lengua original de las fuentes del audio (en este caso, el francés) ya que el Análisis de Sentimientos depende muy significativamente de matices como el sarcasmo, la doble negación, los modismos regionales, etc. Esto hace recomendable contar con un experto lingüista o con hablantes nativos para validar las etiquetas emocionales y supervisar el proceso de tokenización.

- En futuros análisis, se recomienda implementar una etapa previa de pre-procesamiento de audio que incluya reducción de ruido y normalización previa a la transcripción (*Speech-to-Text*), que sea igual o mejor a la realizada en este trabajo de titulación. Un error del 5% en la transcripción automática puede conducir a una clasificación errónea de la emoción, sobre todo en categorías tan cercanas léxicamente como son "Nostalgia" y "Tristeza".

- Dada la gran cantidad de recursos de procesamiento que requieren las librerías de Procesamiento de Lenguaje Natural (NLP) y los modelos de Deep Learning (como los de Transformers), se sugiere utilizar infraestructura con aceleración por GPU. Dicha infraestructura hará que no solo mejoren los tiempos de ejecución de las tareas de inferencia, sino también que se puedan llevar a cabo procesos de entrenamiento y fine-tuning más eficaces y en escalado.

BIBLIOGRAFÍA

Nair, V., Zhang, N., Mohan, S., & Cheung, A. (2022). Explainable AI for pre-trained code models: What do they learn? When do they not work? arXiv. <https://arxiv.org/abs/2211.12821>

Brown, S. (2021, 21 abril). Machine learning, explained. MIT Sloan. Rexcuperado de <https://mitsloan.mit.edu/ideas-made-to-matter/machine-learning-explained>

Bergman, D. (s.f.). What is machine learning? IBM. Recuperado de <https://www.ibm.com/think/topics/machine-learning>

Shearer, C. (2000). The CRISP-DM model: the new blueprint for data mining. *Journal of Data Warehousing*, 5(1), 13–22.

Talib, R., Hanif, M. K., Ayesha, S. y Fatima, F. (2016). Text Mining: Techniques, Applications and Issues. *International Journal of Advanced Computer Science and Applications (IJACSA)*, 7(11), 414-419. <https://doi.org/10.14569/IJACSA.2016.071153>

Zhang, X., Li, Y., & Zhao, Q. (2023). *Natural Language Processing: Recent Development and Applications*. *Applied Sciences*, 13(20), 11395. <https://doi.org/10.3390/app132011395>

Jiang, T. Gradus, J. y Rosellini, A. (2020). *Supervised machine learning: A brief primer*. *Current Opinion in Psychology*, 34, 1-6. <https://doi.org/10.1016/j.copsyc.2020.03.001>

Sarker, I. H. (2021). *Machine Learning: Algorithms, real-world applications and research directions*. *SN Computing Science*, 2(1). <https://doi.org/10.1007/s42979-021-00592-x>

Guido, R., Ferrisi, S., Lofaro, D., & Conforti, D. (2024). An Overview on the Advancements of Support Vector Machine Models in Healthcare Applications: A Review. *Information*, 15(4), 235. <https://doi.org/10.3390/info15040235>

Breiman, L. Random Forests. *Machine Learning* **45**, 5–32 (2001).
<https://doi.org/10.1023/A:1010933404324>

Lee, C. (2024). Artificial Neural Networks (ANNs) and Machine Learning (ML) Modeling Employee Behavior with Management Towards the Economic Advancement of Workers. *Sustainability*, 16(21), 9516. <https://doi.org/10.3390/su16219516>

Mao, Y. Liu, Q. y Zhang, Y. (2024). *Sentiment analysis: Methods, applications, and challenges*. *Journal of King Saud University – Computer and Information Sciences*.
<https://doi.org/10.1016/j.jksuci.2024.102048>

Google Developers. (2025). *Classification: Accuracy, recall, precision, and related metrics*. Recuperado de <https://developers.google.com/machine-learning/crash-course/classification/accuracy-precision-recall>

StatisticsByJim. (s. f.). *Sensitivity vs Specificity: Definition, Formulas & Interpreting*. Recuperado de <https://statisticsbyjim.com/basics/sensitivity-vs-specificity/>

Ting, K.M. (2011). Confusion Matrix. In: Sammut, C., Webb, G.I. (eds) *Encyclopedia of Machine Learning*. Springer, Boston, MA. https://doi.org/10.1007/978-0-387-30164-8_157

Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.

National Severe Storms Laboratory. (s.f.). *Severe weather 101: Thunderstorm basics*. National Oceanic and Atmospheric Administration.
<https://www.nssl.noaa.gov/education/svrwx101/thunderstorms/>

United Nations Office for Disaster Risk Reduction. (s.f.). *Fluvial (Riverine) Flooding (MH0604)*. Hazard Information Profiles. <https://www.undrr.org/understanding-disaster-risk/terminology/hips/mh0604>

Comisión Económica para América Latina y el Caribe. (2014). *Manual para la evaluación de desastres*. Naciones Unidas. <https://repositorio.cepal.org/handle/11362/35894>

David Ortiz Haro, Patrick Laclemece, Audrey Morel Senatore, Guillaume Delatour. Mutual aid: Collective adaptive behaviour for disaster response and social resilience. Qualitative research after Storm Alex. Behavioural and Social Sciences in Security (BASS) 2024 Conference, Centre for research and evidence on security threats, Jul 2024, Saint-Andrews, United Kingdom. (hal-04718169)

David Ortiz Haro. GESTION DE CRISE ET INTÉGRATION DES POPULATIONS : L'ÉLAN SOLIDAIRE QUI PERDURE SUITE À LA TEMPÊTE ALEX. *Risques Infos*, 2023, Réduire les vulnérabilités face aux risques industriels, 46, pp.36-38. (hal-04508361)

David Ortiz Haro, Patrick Laclemece, Audrey Morel Senatore, Guillaume Delatour. L'intégration des populations : une nouvelle perspective pour les acteurs de secours dans les catastrophes. Congrès Lambda Mu 22 « Les risques au cœur des transitions » (e-congrès) - 22e Congrès de Maîtrise des Risques et de Sécurité de Fonctionnement, Institut pour la Maîtrise des Risques, Oct 2020, Le Havre (e-congrès), France. (hal-03454718v2)

Ortiz Haro, D. (2025). Les comportements collectifs adaptatifs prosociaux dans les situations de catastrophe [Tesis de doctorado, Université de Technologie de Troyes]. *Theses.fr*. <http://www.theses.fr/2025TROY0018/document>