

# Preguntas

June 5, 2025

## 1 Preguntas

---

### 1.1 Justificación y Contexto del Problema

**P:** ¿Por qué eligieron específicamente el reconocimiento emocional por voz en lugar de otras modalidades como reconocimiento facial o análisis de texto?

**R:** Elegimos el reconocimiento emocional por voz porque la voz contiene información emocional rica y compleja, algo que los profesionales de salud mental han utilizado durante décadas. Parámetros como el tono, la velocidad del habla, las pausas, la intensidad y los patrones prosódicos revelan estados emocionales específicos de manera objetiva. Además, es menos invasiva que el video, más accesible que los sensores fisiológicos, y puede funcionar en contextos donde otras modalidades no están disponibles.

**P:** ¿Cómo planean abordar las diferencias culturales en la expresión emocional, especialmente entre el español mexicano y el inglés?

**R:** Para abordar las diferencias culturales, incluimos el Mexican Emotional Speech Database (MESD) para capturar las características prosódicas y fonéticas del español mexicano. Combinaremos este dataset con otros en inglés, como RAVDESS y Speech Emotion Recognition (SER), para crear un modelo más robusto y generalizable. Durante el entrenamiento, evaluaremos el rendimiento por idioma y aplicaremos técnicas de adaptación de dominio si es necesario para asegurar un desempeño equitativo.

**P:** ¿Cuáles son las limitaciones éticas de implementar un sistema de reconocimiento emocional automático en salud mental?

**R:** Las principales limitaciones éticas son: \* **Privacidad y confidencialidad:** El manejo de datos sensibles (grabaciones de voz) requiere un consentimiento informado explícito y mecanismos robustos de anonimización y cifrado. \* **Precisión y malinterpretación:** Un diagnóstico erróneo basado en IA podría tener consecuencias graves. El sistema debe ser solo una herramienta de apoyo, no de diagnóstico definitivo. \* **Sesgo algorítmico:** Si los datasets de entrenamiento no son diversos, el modelo podría tener un rendimiento inferior o sesgado para ciertos grupos demográficos, perpetuando desigualdades. \* **Uso indebido:** La tecnología podría ser utilizada para vigilancia o manipulación si cae en manos equivocadas. \* **Deshumanización:** Dependiendo demasiado de la IA podría reducir la empatía y la conexión humana en la atención psicológica.

**P:** ¿Cómo asegurar la privacidad y seguridad de los datos de audio sensibles de los usuarios?

**R:** Aseguraremos la privacidad y seguridad mediante: \* **Anonimización:** Eliminar cualquier información de identificación personal de los audios. \* **Cifrado:** Cifrar los datos de audio tanto en tránsito como en reposo. \* **Consentimiento informado:** Obtener un consentimiento claro y explícito de los usuarios sobre cómo se usarán sus datos. \* **Acceso restringido:** Implementar controles de acceso estrictos a los datos y modelos. \* **Políticas de retención de datos:** Definir políticas claras sobre cuánto tiempo se almacenarán los datos. \* **Conformidad legal:** Adherirnos a regulaciones de protección de datos como GDPR o la Ley Federal de Protección de Datos Personales en Posesión de los Particulares en México.

---

## 1.2 Metodología y Técnicas

**P:** ¿Qué métricas de evaluación específicas se utilizarán y por qué son adecuadas para el reconocimiento emocional?

**R:** Utilizaremos: \* **Precisión (Accuracy):** La proporción de predicciones correctas sobre el total. Útil para una visión general. \* **Precisión por clase (Precision):** La proporción de verdaderos positivos respecto a todos los que el modelo clasificó como positivos para una clase. Mide falsos positivos. \* **Recall por clase (Sensibilidad):** La proporción de verdaderos positivos respecto a todos los que realmente pertenecen a esa clase. Mide falsos negativos. \* **F1-Score por clase:** La media armónica de precisión y recall. Es crucial cuando las clases están desbalanceadas y necesitas un equilibrio entre falsos positivos y falsos negativos. \* **Matriz de Confusión:** Visualiza el rendimiento del modelo en cada clase, mostrando dónde se confunde el modelo (ej. si confunde “miedo” con “sorpresa”).

Estas métricas son adecuadas porque no solo nos dan una visión global del rendimiento, sino que también nos permiten analizar el desempeño del modelo para cada emoción individualmente, lo cual es vital en un contexto donde el error en una emoción (ej. tristeza severa) puede ser más crítico que en otra (ej. neutralidad).

**P:** ¿Por qué es importante la fase de preprocesamiento de audio y cuáles son los desafíos clave?

**R:** El preprocesamiento es crucial porque: \* **Normaliza la calidad:** Los datasets provienen de diversas fuentes con distintas calidades de grabación (ruido, volumen, frecuencia de muestreo), y el preprocesamiento las unifica para que el modelo no aprenda artefactos. \* **Extrae características relevantes:** Ayuda a aislar las señales de voz del ruido y a resaltar las propiedades acústicas que realmente portan información emocional. \* **Optimiza el rendimiento del modelo:** Datos limpios y estandarizados resultan en modelos más precisos y generalizables.

Los desafíos clave incluyen: \* **Ruido de fondo:** Eliminar el ruido ambiental sin perder información emocional valiosa. \* **Variabilidad en la voz:** Diferencias en el tono, acento, género y edad de los hablantes. \* **Duración variable:** Gestionar grabaciones de diferentes longitudes. \* **Silencios y pausas:** Identificar y eliminar segmentos sin habla activa. \* **Desbalance de clases:** Si algunas emociones tienen menos muestras, el modelo puede sesgarse. \* **Complejidad emocional:** Las emociones pueden ser mixtas o sutiles, dificultando su categorización.

**P:** ¿Cómo justificarían la elección de un modelo de aprendizaje profundo (CNN) sobre modelos más tradicionales (SVM, Random Forest) para este problema?

**R:** Justificamos la elección de CNN por: \* **Extracción automática de características:** Las

CNNs pueden aprender directamente las características discriminativas de los espectrogramas de audio sin necesidad de ingeniería manual de características, lo cual es una ventaja significativa sobre SVM o Random Forest que requieren características extraídas explícitamente. \* **Capacidad para capturar patrones complejos:** Las emociones se manifiestan como patrones espaciotemporales intrincados en el audio. Las CNNs son excelentes para identificar estos patrones locales y sus relaciones a lo largo del tiempo y la frecuencia. \* **Rendimiento superior en datos de audio/imagen:** En tareas de clasificación basadas en datos de secuencia o imagen (como los espectrogramas de audio), las CNNs han demostrado consistentemente un rendimiento superior a los modelos tradicionales. \* **Escalabilidad:** Las CNNs escalan mejor con grandes volúmenes de datos y pueden aprovechar la aceleración por GPU, lo cual es importante para nuestros datasets combinados.

Aunque SVM y Random Forest son buenos modelos de referencia y se usarán para comparación, las CNNs son más adecuadas para la complejidad y la naturaleza de los datos de audio emocional.

**P: ¿Qué técnicas de reducción de dimensionalidad planean usar y cuándo las aplicarían en el pipeline?**

**R:** Planeamos usar: \* **Análisis de Componentes Principales (PCA):** Principalmente en la fase de **análisis exploratorio**. Nos ayudará a visualizar la estructura de los datos en un espacio de menor dimensión y a entender la separabilidad de las clases emocionales. \* **Análisis Discriminante Lineal (LDA):** Será la técnica principal para reducción de dimensionalidad en la fase de **entrenamiento del modelo**. LDA es supervisado y busca maximizar la separación entre clases, lo cual es ideal para problemas de clasificación.

Aplicación temporal: \* **PCA:** Se aplicaría después de la extracción de características y la normalización, pero antes de cualquier división de datos para visualización. \* **LDA:** Se aplicaría después de todo el preprocesamiento de características y normalización, pero **antes de dividir los datos en conjuntos de entrenamiento, validación y prueba**. Esto asegura que la transformación se aprenda solo de los datos de entrenamiento para evitar *data leakage*.

---

### 1.3 Implementación y Evaluación

**P: ¿Cómo manejarían el desbalance de clases si algunas emociones están subrepresentadas en los datasets?**

**R:** Si encontramos desbalance de clases, lo manejaríamos con las siguientes estrategias: \* **Sobremuestreo de la clase minoritaria (Oversampling):** \* **SMOTE (Synthetic Minority Over-sampling Technique):** Crearía nuevas muestras sintéticas de las clases subrepresentadas basadas en las existentes. \* **Duplicación aleatoria:** Simplemente copiar algunas muestras de las clases minoritarias. \* **Submuestreo de la clase mayoritaria (Undersampling):** \* **Random Undersampling:** Eliminar aleatoriamente algunas muestras de las clases sobrerrepresentadas. \* **Tomek Links o NearMiss:** Métodos más sofisticados que eliminan muestras de las clases mayoritarias que son “difíciles de aprender” o están cerca de los límites de las clases minoritarias. \* **Aumento de datos (Data Augmentation):** Aplicar transformaciones al audio existente (ej. añadir ruido ligero, cambiar el pitch, variar la velocidad, aplicar filtros) para crear nuevas muestras sin desvirtuar la emoción. \* **Técnicas a nivel de algoritmo:** \* **Ponderación de clases:** Asignar un peso mayor a las clases minoritarias durante la función de pérdida del modelo, haciendo que el

modelo se “preocupe” más por clasificarlas correctamente. \* **Umbral de decisión ajustado:** Modificar los umbrales de probabilidad para la clasificación en la salida del modelo.

La elección final dependerá del grado de desbalance y del rendimiento observado durante las pruebas iniciales.

**P: ¿Cómo garantizarían la reproducibilidad de sus resultados?**

**R:** Para garantizar la reproducibilidad, seguiríamos estos pasos: \* **Control de versiones (Git/GitHub):** Todo el código, scripts y configuraciones se mantendrían en un repositorio de Git, lo que permite rastrear cada cambio. \* **Gestión de entornos (Conda/venv):** Se crearía un entorno virtual o `conda` con versiones específicas de todas las librerías (`requirements.txt` o `environment.yml`). \* **Semillas aleatorias fijas:** Estableceríamos semillas aleatorias fijas para todas las operaciones que involucren aleatoriedad (ej. inicialización de pesos de redes neuronales, división de datos, generación de números aleatorios). \* **Documentación clara:** Documentaríamos detalladamente cada paso de la metodología, incluyendo preprocesamiento, extracción de características, arquitectura del modelo, hiperparámetros y estrategias de evaluación. \* **Datasets accesibles:** Asegurar que los datasets utilizados sean públicos y se proporcionen enlaces claros a sus fuentes, o si se crean datasets procesados, cómo recrearlos. \* **Resultados y métricas detalladas:** Publicaríamos todas las métricas de evaluación, matrices de confusión y resultados de la validación cruzada. \* **Uso de validación cruzada k-fold:** Para asegurar la robustez de la evaluación. \* **Monitoreo y visualización del entrenamiento:** Con herramientas como TensorBoard o Weights & Biases para registrar métricas y graficar el progreso.

**P: ¿Cómo planean validar su modelo para asegurar que generalice bien a nuevas voces y situaciones?**

**R:** Para asegurar una buena generalización, implementaremos: \* **Validación cruzada k-fold estratificada:** Dividiremos el dataset en “k” subconjuntos (folds), y en cada iteración, un fold diferente será el conjunto de prueba, mientras que los restantes se usarán para entrenamiento y validación. Estratificada significa que se mantiene la proporción de clases emocionales en cada fold. Esto garantiza que el modelo se prueba en diferentes subconjuntos de datos. \* **Conjunto de prueba independiente:** Una porción del dataset (ej. 10-20%) se reservará desde el principio y no se usará para entrenamiento ni para ajuste de hiperparámetros. Este conjunto se utilizará solo una vez, al final, para la evaluación final del modelo y obtener una estimación imparcial de su rendimiento en datos no vistos. \* **Diversidad del dataset:** La combinación de MESD (español mexicano), RAVDESS (inglés, profesional) y SER (inglés, diversas fuentes) ya nos proporciona una base de datos más diversa en términos de hablantes, acentos, calidades de grabación y contextos. \* **Aumento de datos:** Aplicar técnicas de aumento de datos (ruido, cambios de pitch/velocidad) para hacer el modelo más robusto a variaciones y mejorar su generalización. \* **Early stopping:** Para prevenir el *overfitting* y detener el entrenamiento cuando el rendimiento en el conjunto de validación deja de mejorar.

---

## 1.4 Aplicaciones y Escalabilidad

**P: ¿Cómo implementar el sistema en tiempo real?**

**R:** Para implementar el sistema en tiempo real, se requiere: \* **Procesamiento en streaming:** Dividir el audio de entrada en *chunks* (segmentos) pequeños y procesarlos de forma continua. \*

**Buffer circular:** Mantener una ventana deslizante de audio para asegurar que el modelo siempre tenga suficiente contexto temporal para la emoción. \* **Optimización del modelo:** Utilizar un modelo ligero y optimizado para una inferencia rápida (baja latencia), posiblemente con técnicas de cuantificación o *pruning*. \* **Hardware adecuado:** Aprovechar el poder de las GPUs o NPUs para un procesamiento paralelo eficiente, o usar dispositivos *edge* para inferencia local. \* **Integración con APIs:** Exponer el modelo como un servicio web a través de una API RESTful para que otras aplicaciones puedan consumirlo fácilmente.

**P: ¿Cómo escalar para múltiples idiomas?**

**R:** Para escalar el sistema a múltiples idiomas, consideramos varias estrategias: \* **Transfer learning:** Reutilizar las capas base del modelo entrenadas en un idioma o en un dataset multilingüe grande y luego realizar *fine-tuning* (ajuste fino) con datos específicos del nuevo idioma. \* **Características universales:** Los MFCCs y otras características acústicas son relativamente universales en su capacidad para representar el contenido de voz, lo que facilita la adaptación inter-idioma en las primeras capas del modelo. \* **Datasets multilingües:** Entrenar conjuntamente el modelo en datasets que contengan datos de múltiples idiomas para que aprenda representaciones más generales e idiomáticamente independientes. \* **Adaptación de dominio:** Si el rendimiento decae significativamente en un nuevo idioma, se pueden aplicar técnicas de adaptación de dominio para alinear los espacios de características entre el idioma fuente y el objetivo. \* **Modelos específicos por idioma:** En casos donde las diferencias son muy grandes, podría ser necesario entrenar un modelo específico para cada idioma, aunque esto es menos eficiente.

---