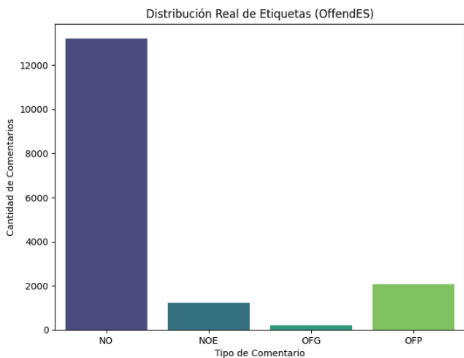
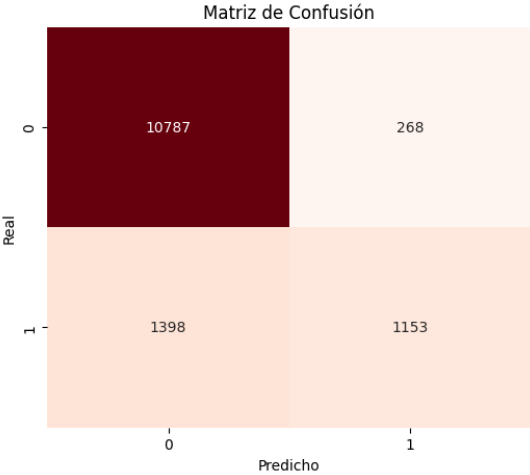


Fecha	Actividad	Solución y Decisiones Técnicas	Evidencias
Lunes 22 de Diciembre de 2025	<p>Definición del Alcance y Carga de Datos</p> <p>Originalmente, planteamos un modelo multclasificación para 4 categorías:</p> <ol style="list-style-type: none">1. Ofensivo (Insulto directo).2. Vulgar (Groserías sin ofensa).3. Ambos.4. Ninguno. <p>Problema Crítico: Al descargar el dataset (training_set.tsv) y cargarlo localmente, el DataFrame mostró dimensiones erróneas (4, 6). Parecía un archivo corrupto</p>	<p>Solución de Carga:</p> <p>Cambio de estrategia a lectura remota. Utilizamos la URL "Raw" del repositorio de GitHub para forzar la descarga íntegra.</p> <p>Ajuste de Formato:</p> <p>Detectamos que el uso de comas en los comentarios rompía el CSV. Se forzó el parámetro sep='\\t' (tabulaciones).</p> <p>Dimensiones finales correctas: (16710, 6).</p>	<p>Error de Carga:</p> <p>ParserError: Error tokenizing data. C error: Expected 2 fields in line 6, saw 17</p> <p>Carga Exitosa:</p> <pre>import pandas as pd # URLs del repositorio url_train = 'https://raw.githubusercontent.com/fmplaza/offendES/main/split_meOffendES/training_set.tsv' url_test = 'https://raw.githubusercontent.com/fmplaza/offendES/main/split_meOffendES/test_set.tsv' # Tuvimos un error en una versión anterior por no especificar sep='\\t' porque es un archivo separado # por tabulaciones y pensamos que era por comas. try: df_train = pd.read_csv(url_train, sep='\\t') df_test = pd.read_csv(url_test, sep='\\t') print(f'Datos cargados. Train: {df_train.shape}, Test: {df_test.shape}') except Exception as e: print(f'Error cargando los datos:', e) Datos cargados. Train: (16710, 6), Test: (13606, 6)</pre>

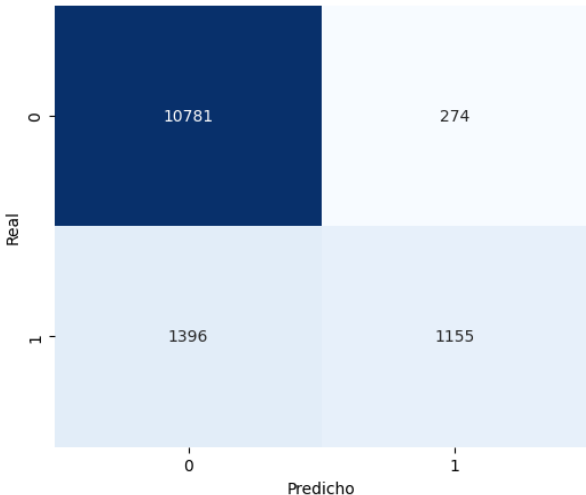
Fecha	Actividad	Solución y Decisiones Técnicas	Evidencias										
	o una "vista previa".												
Martes 23 de Diciembre de 2025	<p>Análisis Exploratorio (EDA) y Re-definición</p> <p>Hallazgo: Al graficar las 4 etiquetas originales, encontramos un desbalance inviable:</p> <p>- NO (No ofensivo): 13,212 muestras.</p> <p>- OFG (Ofensivo Grupal): Solo 212 muestras.</p> <p>Riesgo: Con tan pocos datos en OFG, el modelo podría ignorar esta clase para inflar el Accuracy.</p>	<p>Decisión:</p> <p>Simplificar el problema de Multiclase a Clasificación Binaria.</p> <p>Agrupación:</p> <p>- NO y NOE =Clase 0 (No Ofensivo).</p> <p>- OFP y OFG =Clase 1 (Ofensivo).</p> <p>Esto nos permite consolidar la minoría y mejorar el aprendizaje.</p>	<p>Gráfica de Desbalance:</p>  <table><caption>Distribución Real de Etiquetas (OffendES)</caption><tr><th>Tipo de Comentario</th><th>Cantidad de Comentarios</th></tr><tr><td>NO</td><td>13,212</td></tr><tr><td>NOE</td><td>1,000</td></tr><tr><td>OFG</td><td>212</td></tr><tr><td>OFP</td><td>2,000</td></tr></table> <p>Diccionario de Mapeo:</p> <pre># Mapeo de etiquetas y_train = df_train['label'].map({'NO': 0, 'NOE': 0, 'OFP': 1, 'OFG': 1}) y_test = df_test['label'].map({'NO': 0, 'NOE': 0, 'OFP': 1, 'OFG': 1})</pre>	Tipo de Comentario	Cantidad de Comentarios	NO	13,212	NOE	1,000	OFG	212	OFP	2,000
Tipo de Comentario	Cantidad de Comentarios												
NO	13,212												
NOE	1,000												
OFG	212												
OFP	2,000												

Fecha	Actividad	Solución y Decisiones Técnicas	Evidencias																														
	<p>Preprocesamiento</p> <p>Problema: Necesitábamos limpiar el texto sucio de redes sociales.</p> <p>Acción: Aplicamos una limpieza estándar con Regex, eliminando todo carácter que no fuera letra o número.</p> <p>Eliminamos los emojis.</p>	<p>Implementación :</p> <p>Se creó la función limpieza_basica y se vectorizó con TF-IDF.</p> <p><i>Nota: En este punto, creímos que eliminar caracteres especiales era la mejor práctica para reducir ruido.</i></p>	<pre>import re # Definimos la función de limpieza def limpieza_basica(text): text = str(text).lower() # Mantenemos solo letras y espacios. text = re.sub(r'^a-záéíóúñ\s', '', text) text = re.sub(r'\s+', ' ', text).strip() return text # Aplicamos la limpieza df_train['comment_clean'] = df_train['comment'].apply(limpieza_basica) df_test['comment_clean'] = df_test['comment'].apply(limpieza_basica)</pre>																														
Sábado 27 de diciembre de 2025	<p>Entrenamiento Modelo Base (SVM v1)</p> <p>Configuración: Utilizamos</p>	Por los resultados, vemos que el modelo está sesgado hacia la clase mayoritaria (No Ofensivo)	<p>Reporte de Clasificación:</p> <table><tr><th>Resultados</th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>No Ofensivo</td><td>0.89</td><td>0.98</td><td>0.93</td><td>11055</td></tr><tr><td>Ofensivo</td><td>0.81</td><td>0.45</td><td>0.58</td><td>2551</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.88</td><td>13606</td></tr><tr><td>macro avg</td><td>0.85</td><td>0.71</td><td>0.75</td><td>13606</td></tr><tr><td>weighted avg</td><td>0.87</td><td>0.88</td><td>0.86</td><td>13606</td></tr></table>	Resultados	precision	recall	f1-score	support	No Ofensivo	0.89	0.98	0.93	11055	Ofensivo	0.81	0.45	0.58	2551	accuracy			0.88	13606	macro avg	0.85	0.71	0.75	13606	weighted avg	0.87	0.88	0.86	13606
Resultados	precision	recall	f1-score	support																													
No Ofensivo	0.89	0.98	0.93	11055																													
Ofensivo	0.81	0.45	0.58	2551																													
accuracy			0.88	13606																													
macro avg	0.85	0.71	0.75	13606																													
weighted avg	0.87	0.88	0.86	13606																													

Fecha	Actividad	Solución y Decisiones Técnicas	Evidencias
	<p>LinearSVC por defecto.</p> <p>Resultado: Accuracy engañoso del 88%.</p> <p>Falla: Al revisar la matriz de confusión, notamos que el modelo casi nunca predecía la clase 'Ofensivo' (Recall ~0.45).</p>	<p>debido al desbalance que detectamos en la sesión anterior, además, nos estamos cuestionando si mantener los emojis porque creemos que la limpieza les quitó contexto a los insultos.</p> <p>Necesitamos arreglar la limpieza y el balanceo.</p>	<p>Matriz de Confusión Roja:</p> 

Fecha	Actividad y Problema Detectado	Solución y Decisiones Técnicas	Evidencia
Lunes 29 de Diciembre e de 2025	<p>Análisis de Errores</p> <p>Al revisar manualmente los "Falsos Negativos", notamos que muchos insultos solo podían ser considerados</p>	<p>Implementación de Librería emoji:</p> <p>Decidimos no borrar los emojis, sino traducirlos a texto para que el modelo pudiera interpretarlos semánticamente.</p>	<pre># Nueva función de limpieza # Ahora usamos emoji.demoji antes de borrar caracteres. def clean_text(text): text = str(text) #Traducimos emojis a texto text = emoji.demoji(text, language='es', delimiters=(" ", " ")) # Pasamos a minúsculas text = text.lower() # Limpieza de caracteres (permitimos guiones bajos _ por los emojis) text = re.sub(r'^a-záéíóüñ0-9\s_', '', text) text = re.sub(r'\s+', ' ', text).strip() return text # Aplicamos la limpieza df_train['comment_clean'] = df_train['comment'].apply(clean_text) df_test['comment_clean'] = df_test['comment'].apply(clean_text)</pre>

Fecha	Actividad y Problema Detectado	Solución y Decisiones Técnicas	Evidencia																														
	<p>insultos con los emojis.</p> <p>El bajo rendimiento no es culpa del algoritmo, sino de que estamos borrando la información más valiosa (los emojis) con nuestra limpieza.</p>	<p>Función clean_text:</p> <p>Usamos emoji.demojize(text, language='es').</p> <p>Esto debería dar al modelo las palabras clave que le faltaban.</p>																															
	<p>Entrenamient o y Persistencia del Error</p> <p>Entrenamos el mismo modelo (LinearSVC) con los nuevos datos enriquecidos con emojis.</p> <p>Resultado:</p>	<p>El problema raíz es el Desbalance de Clases (detectado en la Sesión 2), no solo la limpieza.</p> <p>Necesitamos forzar matemáticamente al modelo a prestar atención a la clase minoritaria.</p>	<div>Resultados:</div> <table><thead><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr></thead><tbody><tr><td>No Ofensivo</td><td>0.89</td><td>0.98</td><td>0.93</td><td>11055</td></tr><tr><td>Ofensivo</td><td>0.81</td><td>0.45</td><td>0.58</td><td>2551</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.88</td><td>13606</td></tr><tr><td>macro avg</td><td>0.85</td><td>0.71</td><td>0.75</td><td>13606</td></tr><tr><td>weighted avg</td><td>0.87</td><td>0.88</td><td>0.86</td><td>13606</td></tr></tbody></table>		precision	recall	f1-score	support	No Ofensivo	0.89	0.98	0.93	11055	Ofensivo	0.81	0.45	0.58	2551	accuracy			0.88	13606	macro avg	0.85	0.71	0.75	13606	weighted avg	0.87	0.88	0.86	13606
	precision	recall	f1-score	support																													
No Ofensivo	0.89	0.98	0.93	11055																													
Ofensivo	0.81	0.45	0.58	2551																													
accuracy			0.88	13606																													
macro avg	0.85	0.71	0.75	13606																													
weighted avg	0.87	0.88	0.86	13606																													

Fecha	Actividad y Problema Detectado	Solución y Decisiones Técnicas	Evidencia									
	<p>Aunque el <i>Accuracy</i> se mantuvo alto (~88%), la métrica clave Recall (Sensibilidad) de la clase 'Ofensivo' apenas mejoró, pues siguió estancada alrededor de 0.45.</p> <p>Mejorar los datos no fue suficiente. El modelo optimiza sus aciertos prediciendo siempre la clase mayoritaria ('No Ofensivo') porque no tiene penalización por fallar en la minoritaria.</p>	<p>Investigamos y encontramos el hiperparámetro <code>class_weight='balance d'</code> para SVM y Regresión Logística.</p>	<div><p>Matriz de Confusión</p><table><tr><th></th><th>0</th><th>1</th></tr><tr><th>0</th><td>10781</td><td>274</td></tr><tr><th>1</th><td>1396</td><td>1155</td></tr></table></div>		0	1	0	10781	274	1	1396	1155
	0	1										
0	10781	274										
1	1396	1155										

Sesión / Fecha	Actividad y Problema Detectado	Solución y Decisiones Técnicas	Evidencia																													
Viernes 2 de enero de 2026	<p>Corrección del Desbalance</p> <p>Acción: Implementamos el parámetro <code>class_weight='balanced'</code> en el SVM y probamos también con Regresión Logística para comparar.</p>	<p>Validación Cruzada:</p> <p>Confirmamos que la mejora no fue suerte usando validación simple. La métrica F1-Score subió considerablemente.</p>																														
	<p>Resultado:</p> <p>El modelo dejó de ignorar los insultos. La sensibilidad (Recall) de la clase ofensiva saltó de ~0.48 a 0.68 en la primera prueba.</p> <p>Observación: Al fin el modelo está castigando sus errores en la clase minoritaria.</p>	<p>Decisión:</p> <p>Nos quedamos con SVM (LinearSVC) como modelo principal, pero necesitamos afinar la regularización para evitar <i>overfitting</i>.</p>	<table><thead><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr></thead><tbody><tr><td>No Ofensivo</td><td>0.93</td><td>0.88</td><td>0.91</td><td>11055</td></tr><tr><td>Ofensivo</td><td>0.59</td><td>0.73</td><td>0.65</td><td>2551</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.85</td><td>13606</td></tr><tr><td>macro avg</td><td>0.76</td><td>0.81</td><td>0.78</td><td>13606</td></tr><tr><td>weighted avg</td><td>0.87</td><td>0.85</td><td>0.86</td><td>13606</td></tr></tbody></table>		precision	recall	f1-score	support	No Ofensivo	0.93	0.88	0.91	11055	Ofensivo	0.59	0.73	0.65	2551	accuracy			0.85	13606	macro avg	0.76	0.81	0.78	13606	weighted avg	0.87	0.85	0.86
	precision	recall	f1-score	support																												
No Ofensivo	0.93	0.88	0.91	11055																												
Ofensivo	0.59	0.73	0.65	2551																												
accuracy			0.85	13606																												
macro avg	0.76	0.81	0.78	13606																												
weighted avg	0.87	0.85	0.86	13606																												

Sesión / Fecha	Actividad y Problema Detectado	Solución y Decisiones Técnicas	Evidencia																														
	<p>Afinación de Hiperparámetros (GridSearch)</p> <p>Problema: ¿Cuál es el mejor valor de regularización (C)?</p> <p>Acción: Ejecutamos una búsqueda exhaustiva (GridSearchCV) probando valores de C: [0.01, 0.1, 1, 10, 100].</p> <p>Resultado:</p> <p>La mejor configuración fue C=0.1. Valores más altos de C causaban ruido y valores más bajos no aprendían bien.</p>	<p>Optimización Final:</p> <p>Se re-entrenó el modelo con C=0.1 y class_weight='balanced'.</p> <p>Métricas Finales en Test:</p> <ul style="list-style-type: none">- Accuracy: 86%- Recall (Ofensivo): 0.73- Precision (Ofensivo): 0.60 <p>Consideramos esto un equilibrio aceptable para un sistema de detección de odio.</p>	<p>Mejor configuración encontrada: {'C': 0.1}</p> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>No Ofensivo</td><td>0.93</td><td>0.89</td><td>0.91</td><td>11055</td></tr><tr><td>Ofensivo</td><td>0.60</td><td>0.73</td><td>0.66</td><td>2551</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.86</td><td>13606</td></tr><tr><td>macro avg</td><td>0.77</td><td>0.81</td><td>0.78</td><td>13606</td></tr><tr><td>weighted avg</td><td>0.87</td><td>0.86</td><td>0.86</td><td>13606</td></tr></table>		precision	recall	f1-score	support	No Ofensivo	0.93	0.89	0.91	11055	Ofensivo	0.60	0.73	0.66	2551	accuracy			0.86	13606	macro avg	0.77	0.81	0.78	13606	weighted avg	0.87	0.86	0.86	13606
	precision	recall	f1-score	support																													
No Ofensivo	0.93	0.89	0.91	11055																													
Ofensivo	0.60	0.73	0.66	2551																													
accuracy			0.86	13606																													
macro avg	0.77	0.81	0.78	13606																													
weighted avg	0.87	0.86	0.86	13606																													

Sesión / Fecha	Actividad y Problema Detectado	Solución y Decisiones Técnicas	Evidencia
Sábado 4 de enero de 2026	<p>Pruebas Manuales y Documentación</p> <p>Probamos el modelo con frases propias (slang mexicano y emojis) para verificar robustez.</p>	<p>Observamos que el modelo demostró ser capaz de entender jerga de internet y contextos con emojis que en la fallaban.</p> <p>Compilación del Notebook final y exportación de requerimientos (requirements.txt).</p>	<p>Pruebas del modelo</p> <p>Frase: 'Me das cringe' Limpieza: 'me das cringe' Resultado: OFENSIVO</p> <p>Frase: 'Eres un inútil, no sirves para nada' Limpieza: 'eres un inútil no sirves para nada' Resultado: OFENSIVO</p> <p>Frase: 'Ojalá te mueras pronto maldito' Limpieza: 'ojalá te mueras pronto maldito' Resultado: OFENSIVO</p> <p>Frase: 'No eres tonta' Limpieza: 'no eres tonta' Resultado: OFENSIVO</p> <p>Frase: 'Me gustan tus videos' Limpieza: 'me gustan tus videos' Resultado: No Ofensivo</p>