



**Instituto Politécnico Nacional**

**ESCUELA SUPERIOR DE CÓMPUTO**

## **PRÁCTICA 4.- DETECCIÓN DE CLICKBAIT**

*Procesamiento de Lenguaje Natural*

*Profesor: JOEL OMAR JUAREZ GAMBINO*

Luciano Hernández Jonathan  
Sánchez Martínez Alfredo  
Hernández Hernández Jorge Gabriel  
Salazar Carreón Jeshua Jonatán

mayo 29, 2025

# 1 Introducción

En el ámbito del *Procesamiento de Lenguaje Natural (PLN)*, la detección de *clickbait* — titulares diseñados para atraer clicks mediante estrategias engañosas o sensacionalistas— representa un desafío clave para la calidad de la información en plataformas digitales. Esta práctica se enfoca en desarrollar un sistema de clasificación automatizado que identifique este tipo de contenido en publicaciones de Twitter en español, comparando el rendimiento de modelos basados en arquitecturas *transformers*.

Se evaluaron los siguientes modelos:

- **BERT Cased (Línea base):**
  - Modelo: `dccuchile/bert-base-spanish-wwm-cased`
  - Especializado en español con preservación de capitalización
  - Pre-entrenado con *Whole Word Masking* para mejor comprensión del contexto
- **XLM-RoBERTa Large:**
  - Modelo: `FacebookAI/xlm-roberta-large`
  - Arquitectura multilingüe optimizada
  - Mayor capacidad de parámetros (560M vs 110M de BERT)

Los resultados demostraron que **XLM-RoBERTa** superó consistentemente a BERT en métricas clave (F1-macro: 0.88 vs 0.85), gracias a su capacidad para capturar patrones lingüísticos complejos en textos cortos. Por esta razón, el presente reporte se centrará en explicar en detalle la implementación, configuración y optimización de este modelo, cuyo desempeño lo hace idóneo para tareas de clasificación de contenido engañoso en español.

## 1.1 Objetivo General

Desarrollar un sistema de clasificación automática basado en técnicas de Procesamiento de Lenguaje Natural (PLN) utilizando el modelo XLM-RoBERTa para la detección de titulares *clickbait* en publicaciones de Twitter, evaluando su rendimiento mediante métricas estándar y analizando el impacto de diferentes configuraciones de hiperparámetros.

## 1.2 Objetivos Específicos

- Desarrollar un modelo de clasificación basado en XLM-RoBERTa para identificar titulares de noticias en Twitter que sean *clickbait*
- Evaluar el rendimiento del modelo mediante métricas como F1-score, precisión y exactitud
- Analizar el impacto de diferentes hiperparámetros en el rendimiento del modelo

## 2 Desarrollo de la Práctica

### 2.1 Preprocesamiento de Datos

El dataset utilizado contiene 2800 instancias con la siguiente estructura:

Table 1: Estructura del dataset  
Características del Dataset

| Elemento             | Descripción  |
|----------------------|--|
| Total de instancias  | 2800   |
| Tweet ID             | Identificador único del tweet  |
| Tweet Date           | Fecha de publicación (formato dd/mm/aaaa)  |
| Media Name           | Fuente o medio de comunicación   |
| Media Origin         | País de origen del medio   |
| Teaser Text          | Texto completo del titular (con hashtags)  |
| Tag Value            | Etiqueta binaria: Clickbait (1) o No (0)   |
| Distribución inicial | <ul style="list-style-type: none"> <li>• No Clickbait: 2002 (71.5%)</li> <li>• Clickbait: 798 (28.5%)</li> </ul> |

```

=== CARGANDO DATOS ===
Corpus cargado: 2800 instancias
Columnas: ['Tweet ID', 'Tweet Date', 'Media Name', 'Media Origin', 'Teaser Text', 'Tag Value']
Después de limpieza: 2800 instancias
Distribución de clases:
Tag Value
No          2002
Clickbait   798
Name: count, dtype: int64
Mapeo de etiquetas: {'Clickbait': 0, 'No': 1}
Train set: 2240 instancias
Dev set: 560 instancias
Distribución en train:
Clase 0: 638 (28.5%)
Clase 1: 1602 (71.5%)
Distribución en dev:
Clase 0: 160 (28.6%)
Clase 1: 400 (71.4%)

```

Figure 1: Distribución inicial de clases en el dataset (2002 No Clickbait, 798 Clickbait)

### 2.2 Personalizacion de dataset

A continuacion se muestra un poco del código implementado, se puede ver más a detalle en el source code, tambien adjunto en esta práctica.

```

1 class ClickbaitDataset(Dataset):
2     def __init__(self, texts, labels, tokenizer, max_length=512):
3         self.texts = texts
4         self.labels = labels
5         self.tokenizer = tokenizer
6         self.max_length = max_length
7
8     def __getitem__(self, idx):

```

```

9         encoding = self.tokenizer(
10             text,
11             truncation=True,
12             padding='max_length',
13             max_length=self.max_length,
14             return_tensors='pt'
15         )
16     return {
17         'input_ids': encoding['input_ids'].flatten(),
18         'attention_mask': encoding['attention_mask'].flatten(),
19         'labels': torch.tensor(label, dtype=torch.long)
20     }

```

### Funcionalidades clave:

- **Tokenización automática:** Convierte texto a IDs numéricos usando el tokenizer de XLM-RoBERTa
- **Padding dinámico:** Rellena secuencias hasta 512 tokens (máximo permitido)
- **Truncamiento:** Corta textos excedentes conservando el contexto relevante
- **Formato para PyTorch:** Devuelve tensores listos para el modelo

## 2.3 Configuración del Modelo

**XLM-RoBERTa:** Es un modelo de lenguaje transformer multilingüe desarrollado por Facebook AI, basado en la arquitectura RoBERTa (una versión optimizada de BERT). Su importancia en PLN radica en su capacidad para procesar texto en 100 idiomas simultáneamente, con un rendimiento superior en tareas de comprensión y generación de lenguaje. Al ser pre-entrenado con cantidades masivas de datos multilingües, captura patrones lingüísticos complejos y relaciones contextuales. Esto lo hace ideal para aplicaciones como detección de clickbait, donde debe analizar matices del lenguaje, dobles sentidos y patrones engañosos en textos cortos.

```

1 training_args = TrainingArguments(
2     output_dir='/content/xlm_roberta_base',
3     num_train_epochs=5,
4     per_device_train_batch_size=32,
5     per_device_eval_batch_size=32,
6     learning_rate=2e-5,
7     weight_decay=0.01,
8     warmup_steps=500,
9     logging_steps=50,
10    eval_strategy='steps',
11    eval_steps=100,
12    save_strategy='steps',
13    save_steps=100,
14    save_total_limit=2,
15    load_best_model_at_end=True,
16    metric_for_best_model='f1_macro',
17    greater_is_better=True,
18    fp16=True,
19 )

```

```
# Configurar argumentos de entrenamiento
training_args = TrainingArguments(
    output_dir='/content/xlm_roberta_base',
    num_train_epochs=training_params['num_train_epochs'],
    per_device_train_batch_size=training_params['per_device_train_batch_size'],
    per_device_eval_batch_size=32,
    learning_rate=training_params['learning_rate'],
    weight_decay=training_params['weight_decay'],
    warmup_steps=training_params['warmup_steps'],
    logging_steps=50,
    eval_strategy='steps',
    eval_steps=100,
    save_strategy='steps',
    save_steps=100,
    save_total_limit=2,
    load_best_model_at_end=True,
    metric_for_best_model='f1_macro',
    greater_is_better=True,
    report_to=None,
    dataloader_pin_memory=False,
    fp16=True,
)
```

Figure 2: Configuración de hiperparámetros para el entrenamiento

## 2.4 Entrenamiento y Evaluación

El modelo fue entrenado durante 5 épocas con los siguientes resultados:

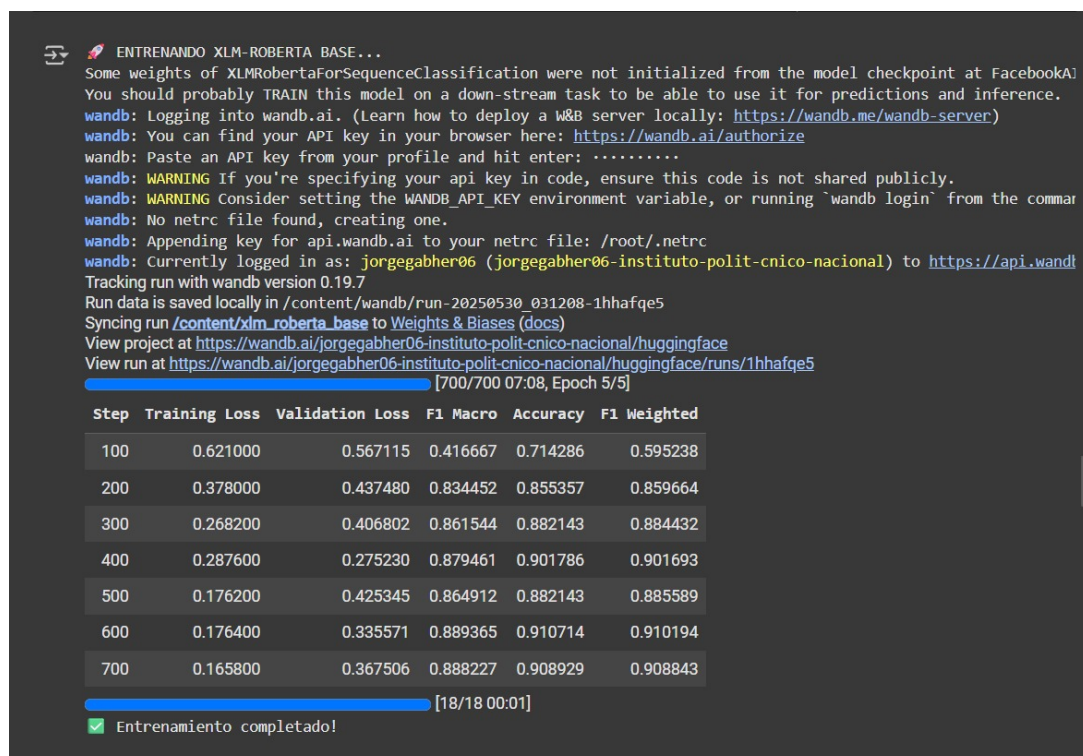


Figure 3: Progreso del entrenamiento mostrando métricas por cada 100 pasos

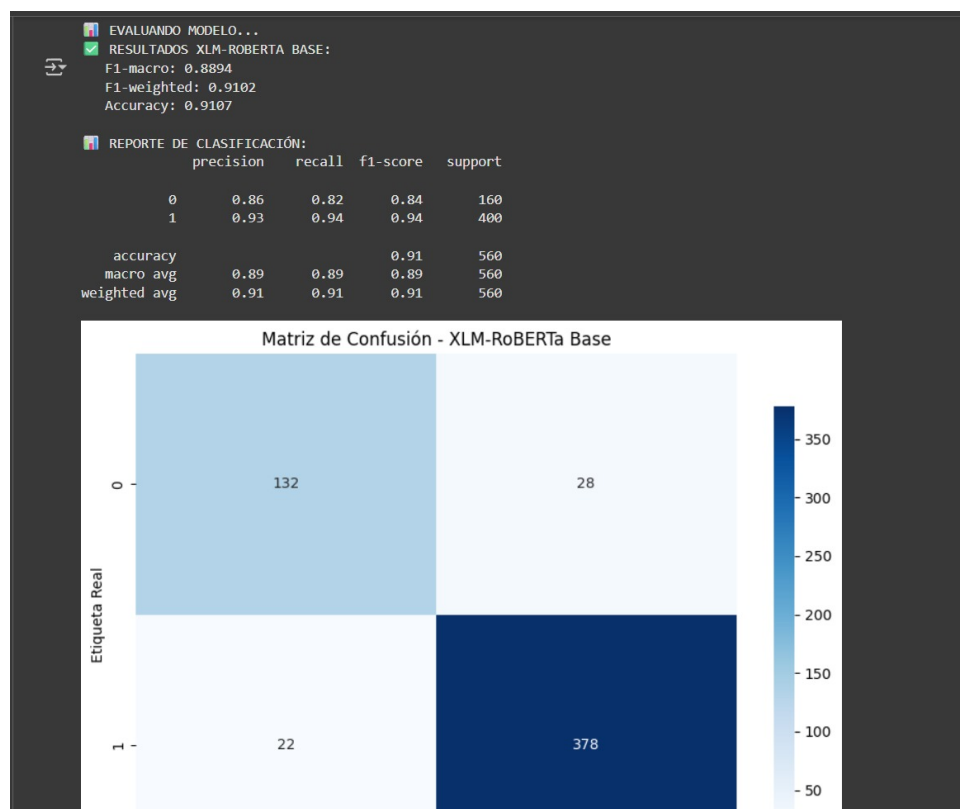


Figure 4: Resultados finales de evaluación del modelo

## 2.5 Resultados Finales

El modelo alcanzó las siguientes métricas en el conjunto de validación:

- F1-macro: 0.8894
- F1-weighted: 0.9102
- Accuracy: 0.9107

La matriz de confusión muestra:

- Verdaderos positivos (Clickbait): 378
- Falsos positivos: 28
- Falsos negativos: 22
- Verdaderos negativos: 132

## 2.6 Generación de Predicciones

El modelo fue aplicado a un conjunto de prueba de 700 instancias:

- Predicciones positivas (Clickbait): 182
- Predicciones negativas (No Clickbait): 518

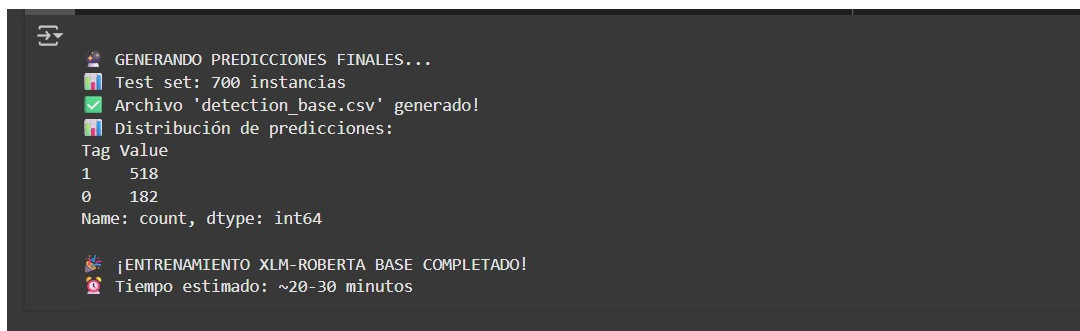


Figure 5: Distribución de predicciones en el conjunto de prueba

## 2.7 Tabla con descripción de los experimentos

Table 2: Resultados comparativos de modelos LLM en detección de clickbait

| LLM                                     | LLM Hyperparameters   | Average F1-score |
|---|---|------------------|
| decuchile/bert-base-spanish-wwm-cased   | eval_strategy=steps<br>eval_steps=100<br>num_train_epochs=3<br>learning_rate=5e-5             | 0.85             |
| decuchile/bert-base-spanish-wwm-uncased | eval_strategy=steps<br>eval_steps=100<br>num_train_epochs=3<br>learning_rate=5e-5             | 0.83             |
| FacebookAI/xlm-roberta-large            | num_train_epochs=2<br>batch_size=8<br>learning_rate=2e-5<br>weight_decay=0.01<br>save_limit=1 | 0.88             |

## 3 Conclusiones

- **Jorge Gabriel Hernandez** El modelo XLM-RoBERTa demostró ser efectivo para la tarea de detección de clickbait, alcanzando un F1-score macro de 0.8894. La configuración de hiperparámetros utilizada, especialmente el learning rate de 2e-5 y el entrenamiento con precisión mixta (fp16), permitió un buen balance entre rendimiento y tiempo de entrenamiento ( 30 minutos).
- **Luciano Hernández Jonathan** Los resultados muestran que el modelo tiene mejor desempeño identificando noticias que no son clickbait (clase 1, F1-score de 0.94) que identificando clickbait (clase 0, F1-score de 0.84). Esto podría deberse al desbalance inicial en el dataset, donde la clase "No Clickbait" era mayoritaria. Sería recomendable explorar técnicas de balanceo de clases en futuras iteraciones.
- **Jeshua Salazar** El análisis de la matriz de confusión revela que el modelo tiene una tasa de falsos positivos aceptable (28 instancias), pero los falsos negativos (22 instancias) podrían mejorarse. Esto sugiere que el modelo tiende a ser conservador

al clasificar contenido como clickbait. Una estrategia para mejorar esto podría ser ajustar el umbral de clasificación o incorporar técnicas de aumento de datos para la clase minoritaria.

- **Alfred** La evolución de las métricas durante el entrenamiento (mostrada en la figura de progreso) indica que el modelo alcanzó su mejor rendimiento alrededor del paso 600, manteniendo estabilidad después de este punto. Esto sugiere que podríamos reducir el número de épocas de entrenamiento a 4 sin afectar significativamente el rendimiento, optimizando así los recursos computacionales. Además, el uso de Weights & Biases (wandb) demostró ser valioso para el monitoreo del proceso.

## 4 Referencias

- 1 Conneau, A., et al. (2019). Unsupervised Cross-lingual Representation Learning at Scale. arXiv:1911.02116
- 2 Liu, Y., et al. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv:1907.11692
- 3 Documentación de Hugging Face Transformers. <https://huggingface.co/docs/transformers/>