

“Myanmar Sarcasm Detection in Social Media Using Deep Learning”

1. Problem Statement

Sarcasm is a subtle and complex form of expression that often conveys meaning opposite to the literal wording. In Myanmar social media, sarcastic comments are common, especially in discussions about politics, entertainment, and current events. However, sarcasm is notoriously difficult for machines to detect because it relies on tone, context, cultural nuances, and sometimes humour.

The goal of this project is to build a deep learning model that can automatically detect sarcasm in Myanmar-Language social media posts. Such a system could help improve content moderation, enhance sentiment analysis accuracy, and contribute to the growth of Burmese NLP research.

2. Input and Output

Input: A Myanmar-language social media posts or comments. Text will be normalized to Unicode, cleaned (removing links, excessive symbols), and tokenized while preserving important sarcasm cues such as emojis and punctuation.

Samples:

- “ဒီအစိုးရကတော့ အရမ်းကောင်းတာပေါ့ 😂”
- “ဒီနေ့မိုးရွာလို့လမ်းပိတ်နေတယ်”

Output: Predicted Label:

- **Sarcastic or Not Sarcastic**
- Along with a confidence score indicating the model’s certainty.

Example Output Table:

Input Text	Label	Confidence
ဒီအစိုးရကတော့ အရမ်းကောင်းတာပေါ့ 😏	Sarcastic	0.94
“ဒီနေ့မိုးရွာလို့လမ်းပိတ်နေတယ်”	Not Sarcastic	0.87

3. Dataset

Type: Binary Burmese text classification dataset (Sarcastic vs. Not Sarcastic).

Sources:

- Public Myanmar social media content from:
 - Public Facebook meme and entertainment pages
 - YouTube comments on Myanmar content
- Adapted text from existing Burmese comment datasets (e.g., offensive/ hate speech datasets on Hugging Face) with manual sarcasm relabeling.
- GitHub Repositories:
 - https://huggingface.co/datasets/simbolo-ai/burmese-hatespeech?utm_source=chatgpt.com
 - https://github.com/ye-kyaw-thu/myHateSpeech?utm_source=chatgpt.com

Planned Size: At least 1,000 labeled comments, balanced across both classes (500 Sarcastic, 500 Not Sarcastic).

Labeling Process:

- Sarcasm is marked if the intended meaning contradicts the literal meaning, or if the tone indicates irony or mockery.
- Label Studio: to define two labels such as Sarcastic (1) vs. Not Sarcastic (0).

4. Expected Performance

With a clean and balanced dataset, and by fine-tuning a multilingual transformer such as XLM-RoBERTa, the aim is to achieve:

- Accuracy: 70% or higher
- Macro F1-Score: 70% or higher

Performance will be evaluated on a held-out test set to ensure generalization.

5. Motivation

This project addresses a gap in Myanmar NLP by focusing on sarcasm detection and it is challenging yet also important task for improving AI understanding of social media discourse.

Accurate sarcasm detection will:

- Enhance sentiment analysis accuracy for Myanmar-language content.
- Support safer and more content-aware content moderation on social platforms.
- Add to the limited pool of Burmese NLP datasets and models, encouraging further research in low-resource languages.

In addition, this project is also about bridging the gap between cutting-edge NLP research and Myanmar language applications. By creating and sharing a labeled sarcasm dataset and a trained detection model, this work can serve as a foundation for future projects in sentiment analysis, humor detection, and online discourse analysis. This project will deliver:

- Help AI systems better understand Myanmar Social Media tone.
- Reduce misclassification errors in sentiment and moderation systems.
- Contribute a rare Burmese sarcasm dataset to the research community.
- Demonstrate the potential of deep learning for nuanced low-resource language tasks.