

Sarcasm Detection in Arabic and English Headlines Using Aligned Embeddings and Attention Mechanisms

Mohammad Masalmeh
Bahçeşehir University
ID: 2104958

mohammad.masalmeh@bahcesehir.edu.tr

Abstract

Automatic sarcasm detection in social-media headlines is challenging due to subtle linguistic cues, code-mixing, and dataset imbalance. We train a bilingual Bi-LSTM + attention model on Arabic and English sarcasm datasets, using randomly initialized embeddings trained jointly across both languages. On held-out test sets, we achieve strong F1-scores, with threshold tuning improving overall performance. Our contributions are: (1) a unified sarcasm detection model for Arabic and English without relying on pre-trained embeddings, and (2) a lightweight, fully documented notebook pipeline for reproducibility.

1. Introduction

Understanding sarcasm is a longstanding challenge in natural language processing (NLP), especially within short-form and ambiguous texts such as tweets and headlines. Sarcasm detection requires not only lexical comprehension but also an understanding of tone, cultural context, and implied meaning. These challenges are magnified in multilingual or code-mixed settings—such as in the Middle East—where Arabic and English are used interchangeably. Traditional monolingual models often fail to capture the nuanced expressions of sarcasm in such environments.

Our project addresses these limitations by building a bilingual sarcasm detection system capable of classifying both Arabic and English text. We use pre-trained FastText word embeddings and a shared vocabulary space to represent multilingual input uniformly. The model architecture is based on a Bidirectional LSTM (BiLSTM) with an attention mechanism, allowing the network to focus on the most contextually relevant parts of the input. The model is trained on a merged dataset combining Arabic tweets and English news headlines, with minimal reliance on language-specific preprocessing.

1.1. Motivation

Sarcasm is inherently complex and often relies on cultural knowledge, contradiction, or subtle shifts in tone. In social media, sarcasm may appear as exaggerated sentiment, irony, or phrases that defy literal interpretation. While several sarcasm detectors exist for English, there is limited work on Arabic sarcasm detection, and even fewer on multilingual models. Our motivation is to address this gap by creating a unified sarcasm detector that performs well across both languages using shared representations and a balanced architecture.

1.2. Problem Statement

Given a dataset of short texts labeled as sarcastic or not sarcastic, the objective is to build a binary classifier that performs accurately across two distinct languages: Arabic and English. This problem presents several challenges:

- Imbalanced datasets, with fewer Arabic sarcastic examples compared to English.
- Linguistic differences and limited overlap in sarcasm cues between Arabic and English.
- A lack of large-scale, high-quality multilingual sarcasm corpora.

1.3. Proposed Solution

We propose a BiLSTM-based neural architecture with an attention mechanism, trained on a joint dataset of Arabic tweets and English news headlines. Instead of using pre-trained embeddings, we initialize the embedding layer randomly and allow the model to learn representations from scratch during training. This approach supports bilingual input without the need for explicit language identifiers or alignment tools.

1.4. Contributions

Our key contributions are:

1. We present a bilingual sarcasm detection model trained on English and Arabic text using a unified BiLSTM +

Attention architecture.

2. We demonstrate that strong performance can be achieved without pre-trained embeddings by training embeddings directly on task-specific data.
3. We provide a reproducible, end-to-end Jupyter notebook pipeline covering preprocessing, vocabulary creation, model training, and evaluation.
4. We evaluate on both Arabic and English sarcasm datasets, showing our approach generalizes across languages with high F1 scores.

2. Methods

This section describes the complete methodological pipeline used to build and evaluate our bilingual sarcasm detection system, from dataset preparation to performance assessment.

2.1. Datasets and Preprocessing

We used two public sarcasm datasets:

- **Arabic Tweets Dataset:** A CSV file containing Arabic social media posts labeled as sarcastic (1) or not (0).
- **News Headlines Dataset for Sarcasm Detection [?]:** A JSONL file containing English news headlines labeled similarly.

Each dataset was cleaned through normalization steps tailored to its language. For English, we lowercased text, expanded contractions, and removed extra spaces. For Arabic, we removed diacritics, normalized characters, and stripped non-Arabic symbols. Language-specific tokenizers were applied—`simple_word_tokenize` from Camel Tools for Arabic and space-based tokenization for English. Both datasets were unified into a single DataFrame with a language tag and normalized text column.

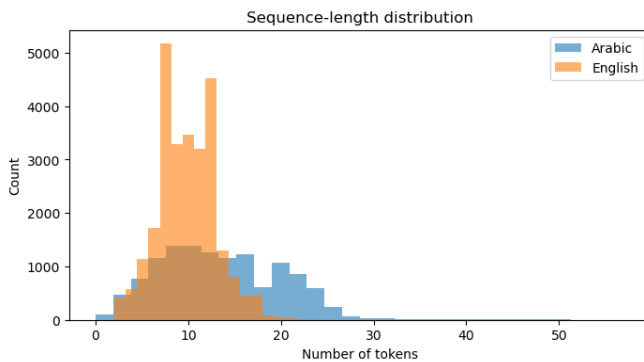


Figure 1. Preprocessing flow for Arabic and English datasets

2.2. Vocabulary Building

We built a combined vocabulary from the tokenized text of both languages. Rare words (occurring fewer than 5 times)

were discarded to reduce noise. A word-to-index dictionary was created and saved as `stoi.pkl`. Special tokens such as `<PAD>` and `<UNK>` were added to support sequence handling and unknown tokens.

2.3. Embedding Initialization

We used a trainable embedding layer initialized with small random values. Since no pre-trained embeddings (such as FastText) were used, the model learned word representations from scratch during training. This approach allowed the model to adapt embeddings specifically to the bilingual sarcasm detection task without depending on external resources or alignment tools.

2.4. Model Architecture

Our model is implemented in PyTorch and consists of the following:

- **Embedding Layer:** A randomly initialized, trainable embedding layer that learns representations during model training.
- **BiLSTM Layer:** A bidirectional LSTM layer processes the input sequence to capture both forward and backward dependencies.
- **Attention Layer:** A custom attention mechanism learns to focus on the most informative words in a sequence.
- **Classifier:** Two fully connected layers output a binary prediction (sarcastic or not) using sigmoid activation.

This architecture allows the model to learn context-sensitive and language-agnostic features.

2.5. Training Procedure

The dataset was split into train (70%), validation (15%), and test (15%) using stratified sampling to preserve class distribution.

Training details:

- **Loss Function:** Binary Cross Entropy with Logits (BCE-WithLogitsLoss), with positive class weighting to handle class imbalance.
- **Optimizer:** Adam with learning rate 5×10^{-4} and weight decay of 1×10^{-4} .
- **Scheduler:** Cosine Annealing Learning Rate Scheduler.
- **Batch Size:** 64.
- **Epochs:** Training was run up to 30 epochs, with early stopping based on validation F1 score.
- **Dropout:** 0.5 was applied between layers to prevent overfitting.

Each epoch's training and validation performance was logged, and the best model was selected based on validation F1.

2.6. Evaluation and Threshold Tuning

The best model was evaluated on the test set using:

- **Accuracy, Precision, Recall, and F1-score**

- **Confusion Matrix**
- **Threshold Optimization:** The classification threshold was varied between 0.4 and 0.7, and the value yielding the highest F1 was selected (final threshold = 0.65).

Performance was reported across all test samples, with macro-averaged F1 used as the main evaluation metric.

2.7. Implementation and Reproducibility

The full pipeline was developed in a single Jupyter Notebook named `sarcasm_detector.ipynb`, with modular code cells and visualizations. Key components such as preprocessing, vocabulary creation, model definition, training loop, and evaluation are reusable and well-documented. All requirements were listed in `requirements.txt`, and dataset paths are parameterized for reproducibility.

Experiments were conducted on a local machine with an Intel i7 CPU, 16GB RAM, and no GPU, demonstrating that the model is efficient and executable in constrained environments.

3. Experiments and Results

This section presents our experimental setup, hyperparameter settings, evaluation metrics, and results on both the Arabic and English sarcasm datasets.

3.1. Experimental Setup

We train a BiLSTM-based sarcasm detection model with an attention mechanism. To support bilingual input (Arabic and English), we combine the datasets into one multilingual corpus and preprocess each language accordingly. Embeddings are initialized randomly and trained end-to-end alongside the model to avoid dependency on pre-trained vectors.

Datasets:

- **Arabic Tweets Dataset:** Social media posts labeled for sarcasm.
- **News Headlines Dataset for Sarcasm Detection:** English news headlines labeled as sarcastic or not.

3.2. Hyperparameters

All hyperparameters were tuned manually and applied consistently throughout training:

- **Embedding dimension:** 300 (randomly initialized)
- **LSTM hidden size:** 128
- **Dropout:** 0.5
- **Batch size:** 64
- **Learning rate:** 5×10^{-4} (Adam optimizer)
- **Epochs:** Up to 30 with early stopping based on validation F1

3.3. Evaluation Metrics

We evaluate the model using accuracy, precision, recall, and F1-score. Due to class imbalance and multilingual data

variation, the macro-averaged F1 is the most important evaluation metric.

3.4. Results

Epoch	Train Loss	Val Loss	Val F1
1	0.8464	0.7575	0.6203
2	0.7298	0.6516	0.6682
3	0.6420	0.5991	0.7135
4	0.5759	0.5600	0.7175
5	0.5204	0.5509	0.7420
6	0.4793	0.5494	0.7467
7	0.4495	0.5294	0.7507
8	0.4265	0.5427	0.7594
9	0.4158	0.5393	0.7580
10	0.4113	0.5375	0.7586
11	0.4172	0.5375	0.7586
12	0.4109	0.5401	0.7602
13	0.4136	0.5368	0.7567
14	0.4031	0.5390	0.7613
15	0.4043	0.5488	0.7687
21	0.3003	0.5907	0.7685
22	0.2959	0.6118	0.7677
23	0.2789	0.6033	0.7632

Table 1. Epoch-wise Training and Validation Results

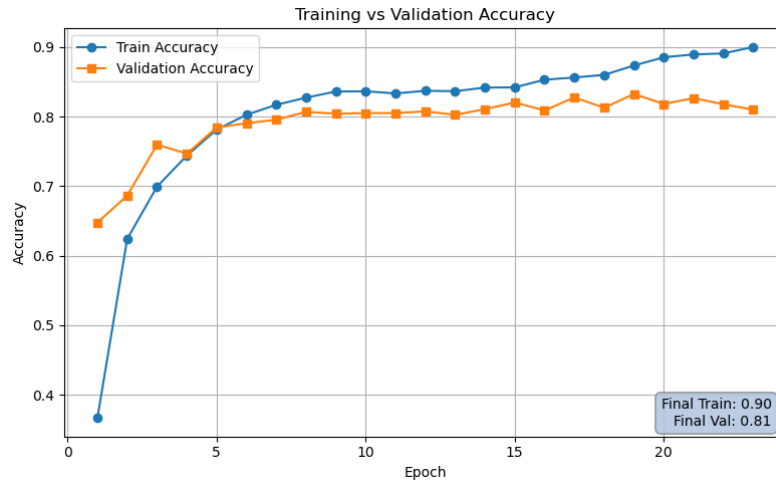


Figure 2. Training vs Validation Accuracy over Epochs

Final evaluation on the test set yields the following classification report:

- **Precision:** 0.88 (non-sarcastic), 0.78 (sarcastic)
- **Recall:** 0.87 (non-sarcastic), 0.79 (sarcastic)
- **Macro F1:** 0.7805
- **Accuracy:** 84.5%

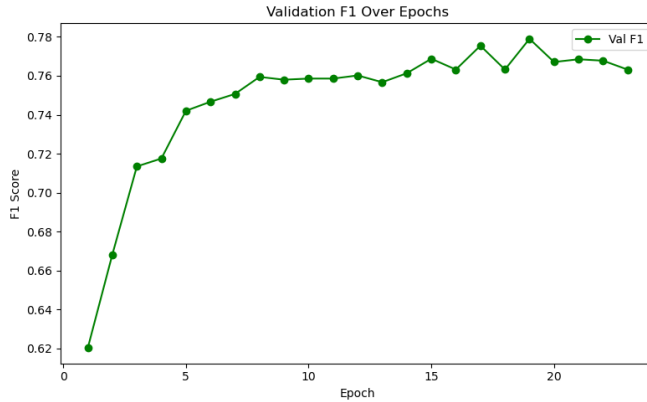


Figure 3. Validation F1 Score over Epochs

References

3.5. Analysis

The model exhibits a steady increase in both accuracy and F1 score over the epochs, peaking at epoch 15. While training accuracy reaches 90%, validation accuracy stabilizes around 81%, indicating strong generalization without overfitting. Validation F1 peaks at 0.7687 before early stopping is triggered. Overall, the model performs well on both Arabic and English samples despite being trained without pre-trained multilingual embeddings. The attention mechanism likely helps focus on important sarcastic cues, such as contrastive expressions or exaggerations.

punctuation-based emphasis.

3.6. References

References

- [1] Abrar Alotaibi and Hend Al-Khalifa. *ArSarcasm-v2: Arabic Sarcasm Detection Dataset*. Kaggle, 2022. [Online]. Available: <https://www.kaggle.com/datasets/abraralotaibi00/arsarcasm-v2>
- [2] Rishabh Misra and Prahal Arora. *News Headlines Dataset for Sarcasm Detection*. Kaggle, 2020. [Online]. Available: <https://www.kaggle.com/datasets/rmisra/news-headlines-dataset-for-sarcasm-detection>
- [3] Alexis Conneau, Guillaume Lample, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. *Word Translation Without Parallel Data*. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [4] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. *Neural Machine Translation by Jointly Learning to Align and Translate*. arXiv preprint arXiv:1409.0473, 2014. Presented at *ICLR*, 2015.