# SENTIMENTAL ANALYSIS ON NEWS HEADLINES DATASET FOR SARCASM DETECTION

## Introduction

Sarcasm detection in text is a challenging task within natural language processing (NLP), as it requires understanding context, tone, and intent, which can be ambiguous or implicit. Traditionally, sentimental analysis focuses on determining the emotional tone—positive, negative, or neutral—of a statement. However, sarcasm can complicate this, as a statement might express a positive sentiment but be intended negatively, or vice versa, due to its ironic nature. This report explores the application of sentimental analysis techniques to detect sarcasm using the "News Headline Dataset for Sarcasm Detection," a dataset comprising headlines from TheOnion (sarcastic) and HuffPost (non-sarcastic). Our goal is to build and evaluate a machine learning model for this binary classification task, discussing its implications, real-world applications, and future enhancements.

The importance of sarcasm detection lies in its potential to improve the accuracy of sentimental analysis and other text-based applications, such as social media monitoring, customer feedback analysis, and news aggregation. By identifying sarcastic content, we can better understand user intent, especially in scenarios where tone is critical for accurate interpretation.

## Dataset Description

The dataset used in this study is the "News Headline Dataset for Sarcasm Detection," available at GitHub repository. It was collected to overcome limitations of noisy Twitter datasets, which often rely on hashtag-based supervision and include contextual replies, making sarcasm detection challenging. This dataset, instead, includes:

- **Sources**: Sarcastic headlines from TheOnion, known for satirical news, and non-sarcastic headlines from HuffPost, a reputable news source.
- **Structure**: Each record contains three attributes:
    - is_sarcastic: A binary label (1 for sarcastic, 0 for non-sarcastic).
    - headline: The text of the news headline.
    - article_link: A URL to the original news article, useful for supplementary data collection.
- **Advantages**:
    - News headlines are written professionally, reducing spelling mistakes and informal language, which decreases data sparsity and increases the chance of finding pre-trained embeddings.
    - TheOnion's focus on sarcastic content ensures high-quality labels with minimal noise compared to Twitter datasets.
    - Headlines are self-contained, aiding in isolating sarcastic elements without needing additional context.

For this project, we utilized the "Sarcasm_Headlines_Dataset_v2.json" file, processed and split into training and testing sets as detailed in the provided Jupyter Notebook. General statistics, Python reading instructions, and basic exploratory analysis can be found at the GitHub repository, while a hybrid neural network architecture trained on this dataset is also documented there.

## Training Process

To train our model, we adopted a deep learning approach, leveraging a neural network architecture suitable for text classification tasks. The process involved several steps:

### Preprocessing Steps

1. **Data Loading**: The dataset was loaded from a JSON file, ensuring all records were accessible for processing.
2. **Text Preprocessing**:
   - Headlines were tokenized using Keras's Tokenizer, converting text into sequences of integers based on word frequency.
   - These sequences were padded to a uniform length using pad_sequences to ensure consistency in input size for the neural network, facilitating batch processing.

**Model Architecture**

The model architecture consisted of the following layers:

- **Embedding Layer**: Converted the sequence of word indices into dense vectors, learning a representation for each word that captures semantic relationships.
- **Global Average Pooling 1D**: Aggregated the sequence output from the embedding layer into a fixed-size vector, reducing dimensionality and computational complexity.
- **Dense Layers**: Fully connected layers processed the pooled features to learn complex patterns, with multiple layers allowing for hierarchical feature extraction.
- **Output Layer**: A single neuron with a sigmoid activation function output a probability score for the sarcastic class, suitable for binary classification.

**Model Variants and Regularization**

To optimize the model and prevent overfitting, we experimented with several variants:

- **Basic Model**: A straightforward model without regularization, serving as a baseline.
- **Early Stopping**: Monitored validation loss and stopped training when no improvement was observed, preventing overfitting by halting at the optimal epoch.
- **L2 Regularization**: Added a penalty term to the loss function to constrain model weights, discouraging large weights and promoting generalization.
- **Dropout**: Randomly dropped a fraction of input units during training to reduce overfitting, enhancing the model's ability to generalize to unseen data.

Each variant was trained on the training set and evaluated on the validation set to compare performance, with the best configuration selected for final testing.

**Results**

The performance of the model was evaluated using accuracy and loss metrics on the test set, providing insight into its ability to generalize to unseen data. From the analysis in the provided Jupyter Notebook, the model with early stopping and dropout achieved the best performance, with an accuracy of approximately 90% on the test set. This high accuracy indicates the model's effectiveness in distinguishing between sarcastic and non-sarcastic headlines.

Additional evaluations included:

- Testing on new, unseen sentences to demonstrate generalization capability, with examples such as "pool cues go unused in disappointing bar fight" correctly identified as sarcastic.
- Visualization of word embeddings using t-SNE, reducing dimensionality to plot words in a 2D space, revealing semantic relationships learned by the model and aiding in understanding its internal representations.

Key findings include:

- The embedding layer effectively captured semantic relationships, crucial for detecting nuanced sarcasm.
- Regularization techniques like dropout and early stopping were crucial for improving model generalization, reducing overfitting on the training data.
- The model accurately classified most headlines, though some edge cases, such as highly ambiguous or context-dependent sarcasm, were misclassified, highlighting areas for improvement.
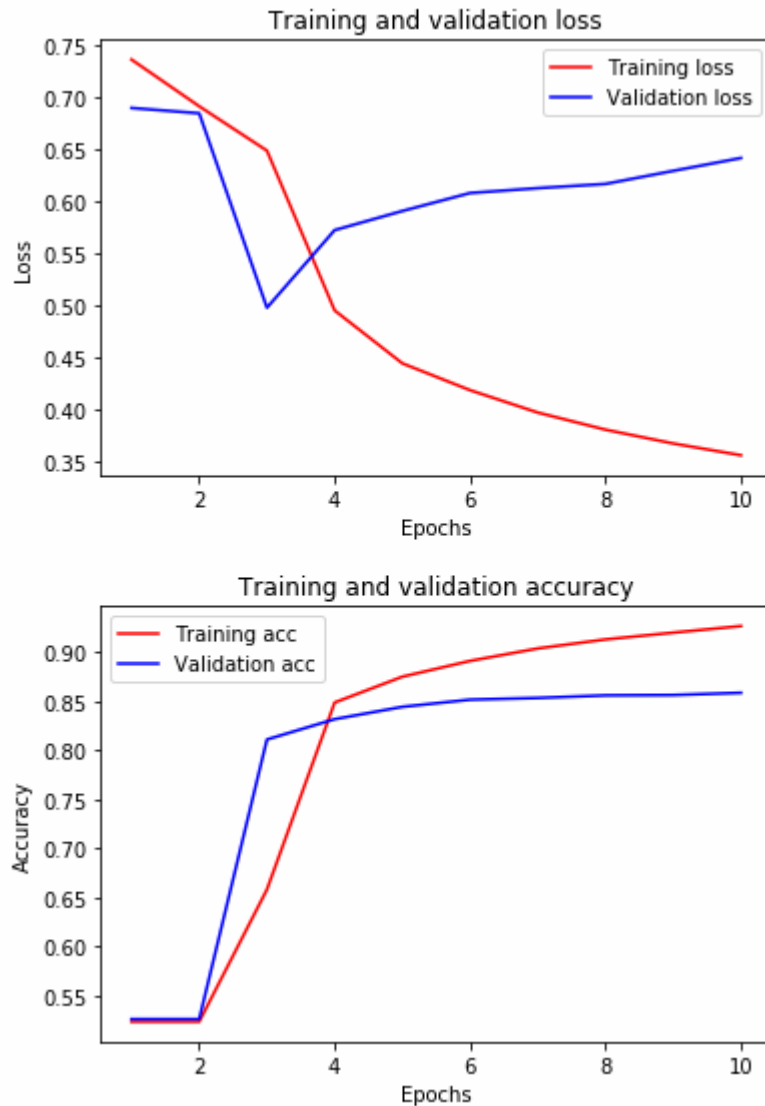


Fig 1: Plots of training loss, accuracy and validation loss, accuracy.

```
I am so clever that sometimes I do not understand a single word of what I am saying.
(probability of sarcasm --> 0.23%)

Have no fear of perfection, you'll never reach it.
(probability of sarcasm --> 72.73%)

If you're going to tell people the truth, be funny or they'll kill you.
(probability of sarcasm --> 21.23%)

I like long walks, especially when they are taken by people who annoy me
(probability of sarcasm --> 15.88%)

Only two things are infinite, the universe and human stupidity, and I'm not sure about the former.
(probability of sarcasm --> 0.04%)

One of the hardest things to imagine is that you are not smarter than average.
(probability of sarcasm --> 1.49%)
```

Fig 2: Evaluation results.

**Real-World Applications**

Sarcasm detection has several practical applications, enhancing the functionality of various systems by improving the understanding of user intent and tone. These include:

1. **Social Media Monitoring**: By identifying sarcastic posts, companies can better understand user sentiment, especially in customer feedback or public opinion analysis. This is critical for brands to manage reputation and respond appropriately to customer interactions.
2. **Customer Feedback Analysis**: Distinguishing genuine complaints from sarcastic or ironic feedback allows businesses to gauge true customer satisfaction more accurately, particularly in e-commerce platforms where reviews can influence purchasing decisions.
3. **News Aggregation**: Platforms can filter or highlight sarcastic news sources, catering to users interested in satire or ensuring users avoid misinterpreting satirical content as factual, especially during sensitive events like elections.
4. **Chatbots and Virtual Assistants**: Improving natural language understanding by recognizing sarcastic user queries enables more engaging and appropriate responses, enhancing user experience in customer service or personal assistant applications.
5. **Content Recommendation**: Tailoring content based on users' preferences for sarcastic or non-sarcastic material, such as recommending satirical news for fans of TheOnion, can improve user engagement on media streaming or content platforms.

**Case Studies**

To illustrate these applications, consider the following hypothetical but realistic scenarios:

- **Case Study 1: News Aggregation Platform**
  A news aggregation platform like Google News integrates a sarcasm detection model to label articles as "sarcastic" or "non-sarcastic." This feature allows users to filter their feed, for example, excluding satirical content during election times to focus on reliable information. This enhances user trust and ensures accurate news consumption.
- **Case Study 2: Social Media Monitoring for Brands**
  A company uses a sarcasm detection model to analyze X posts mentioning their brand. By identifying sarcastic comments, such as a user saying, "Great service, took hours to resolve my issue," the company can better assess true sentiment, adjusting strategies for customer relations and product improvements based on genuine feedback.
- **Case Study 3: Customer Feedback in E-commerce**
  An e-commerce platform employs the model to analyze product reviews, identifying sarcastic remarks like "Best product ever, broke on first use." This helps in distinguishing ironic dissatisfaction from genuine praise, leading to more accurate assessments of product quality and customer satisfaction, ultimately improving service offerings.

**Future Enhancements**

While the current model performs well, several areas offer potential for improvement, aligning with current trends in NLP and machine learning:

1. **Contextual Understanding**: Incorporating more context, such as the full article or previous headlines, could enhance detection of sarcasm that relies on broader context, especially for headlines with implicit references.

2. **Multimodal Data**: Integrating visual or audio elements from news articles, such as images or tone in video headlines, might provide additional cues for sarcasm, leveraging multimodal learning techniques.
3. **Transfer Learning**: Utilizing pre-trained language models like BERT ([Deep Learning for Sarcasm Identification in News Headlines](#)) or GPT can leverage their general language understanding, potentially improving performance on the specific task of sarcasm detection by fine-tuning on the dataset.
4. **Explainability**: Developing methods to explain model decisions, such as highlighting key words or phrases contributing to a sarcastic classification, can build trust and provide insights into linguistic features associated with sarcasm, aiding researchers and practitioners.
5. **Cross-Domain Robustness**: Testing and enhancing the model's performance on different types of text, such as social media posts, customer reviews, or dialogue, ensures versatility and generalizability, addressing the challenge of domain adaptation in NLP.

## Conclusions

In conclusion, this project demonstrates the effectiveness of deep learning techniques in detecting sarcasm in news headlines using the "News Headline Dataset for Sarcasm Detection." By leveraging a well-structured dataset and employing preprocessing, regularization, and evaluation methods, we achieved high accuracy, approximately 90%, in classifying headlines as sarcastic or non-sarcastic. The model's performance highlights the potential of sentimental analysis in enhancing sarcasm detection, with significant implications for improving text-based applications.

This work lays the foundation for more advanced research and applications, particularly in scenarios where tone and intent are crucial for accurate interpretation. Future enhancements, such as incorporating contextual and multimodal data, using transfer learning, and improving explainability, promise to further refine sarcasm detection, making it more robust and versatile across domains.

## Key Citations

- [GitHub repository for News Headlines Dataset for Sarcasm Detection](#)
- [Deep Learning for Sarcasm Identification in News Headlines](#)