# Automatic Sarcasm Detection in Blog Posts

Hamed Minaee and Ali A. Ghorbani

## Introduction

The rapid growth of social media, such as Twitter and blogs, has dramatically increased demand for analyzing the content of social media's posts.. All these information need to be analyzed, understood, and interpreted. Governments, corporations, and non-profits would all greatly benefit from the ability to better comprehend the conversations that are taking place in the digital sphere. Information retrieval is able to shed light on numerous domains such as text categorization, sentiment analysis, topic extraction and issue discovery. Sarcasm detection, which is closely related to all of the topics previously discussed, is an important topic in this field. Sarcasm is defined as the rhetorical process of intentionally using words or expressions for transmitting a meaning different from what literally has been said. In order to detect sarcasm in blog posts we developed a two-step process framework. in each step, we use different feature sets to capture the presence of sarcasm in the text and also rank the texts according to their level of sarcasm.
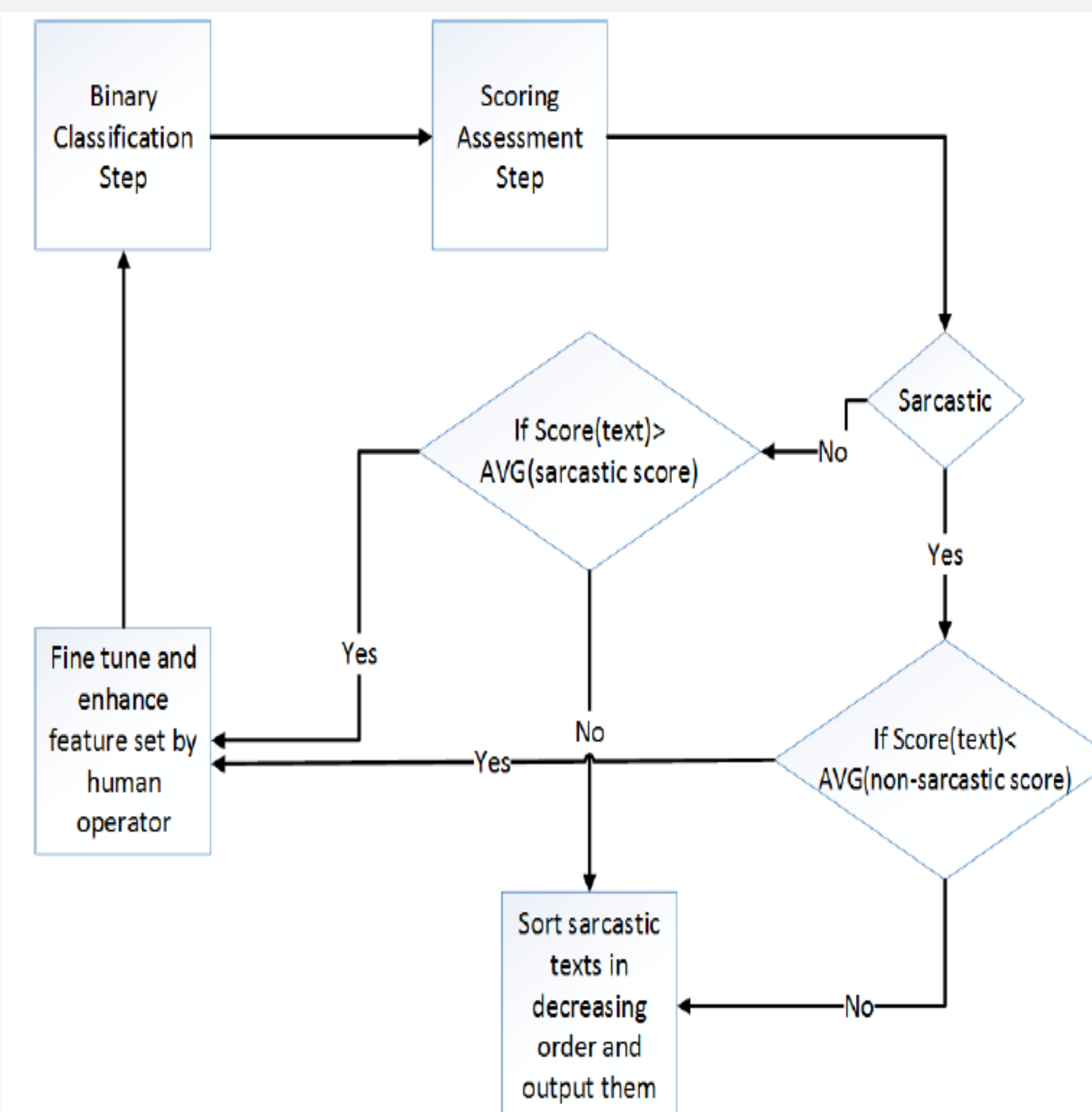
## Contributions

➢ Binary classification of posts along with scoring the posts regarding to their level of sarcasm
➢ Automatic evaluation of binary classification by applying the scoring component

## Dataset

From our preliminary observations, one of the major difficulties when it comes to sarcasm detection in blog posts is the infrequent use of sarcasm. In comparison to tweets, the sheer amount of sarcasm in blogs is limited. Moreover, the medium of text does not accurately reflect tone of the voice. Therefore, it is of utmost importance to have a robust dataset to ascertain sarcasm to compensate for this discrepancy. We used a WordPress API with some predefined search terms from a variety of "trending issues" to collect the data. For instance, search terms such as 'Obama" or 'Trump" are good candidates for our intended purpose. Once the search terms were defined, we collected 4000 blog posts, 400 of which were identified as sarcastic. These results were obtained by three annotators who labelled the data set into two categories: sarcastic and non-sarcastic. In case of a disagreement, the blog posts were categorized as sarcastic or non-sarcastic by the annotators through a majority vote. Furthermore, in addition to capturing the body of the text, we also crawled data such as the blog post title and the blog post writer's username. In the interest of posterity, most duplicate blogs were automatically removed. However, given the presence of syntactic differences, as well as web links, hashtags and other minor differences, the data set contained a small number of duplicate blogs which made duplication removal increasingly difficult. As a result, we added a manual screening inspection to the process to ensure all duplicates were removed.

## Framework and feature set

we propose a novel framework that consists of a two-step process, which uses different feature sets to analyze the entire blog post and the isolated sentences.



- **Binary Classification Step:** This step is organized based on four types of features: tone, adverb intensifier, mood imbalance, and readability of text. Each feature, save for tone and title, is represented with one dimension.
- **Scoring Assessment Step:** In this step, sarcasm detection is identified on a spectrum. Once the first phase is completed, all the blog posts are labeled as sarcastic or non-sarcastic. However, we have yet to assess the varying levels of sarcasm amongst the posts labelled as sarcastic. Thus, we incorporated a number of features in order to accurately score each of our posts as follows:
  1. Textual features
     - contextual imbalance
     - tone
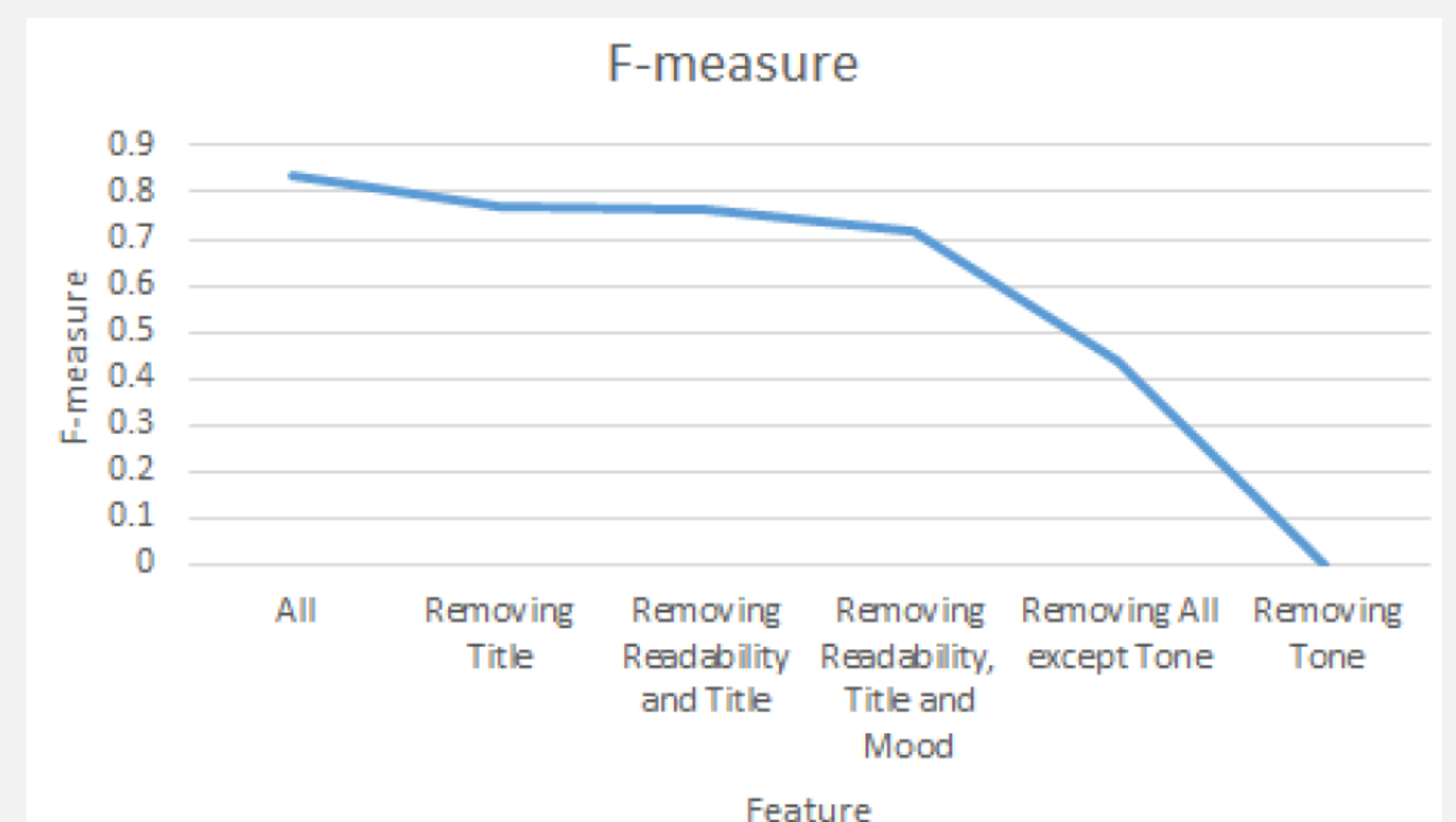     - signature
  2. Patterns

## Classification

The first step of our framework is to classify posts into groups of sarcastic and non-sarcastic based on defined feature set. This classifier is responsible for taking a post and identifying which class it belongs to.
We used 4 different techniques including Decision Tree, Ada Boost, SVM and Random Forest Tree and all classifiers were trained by a set of predefined labeled posts

## Experiments and Results

- **Step 1:** In the first step for binary classification we have used 4 different techniques including Decision Tree, Ada Boost, SVM and Random Forest Tree. The following table shows the result of different classification techniques:

|  | Precision | Recall | F-Measure |
|---|---|---|---|
| Decision Tree | 0.87 | 0.79 | 0.83 |
| Ada Boost | 0.82 | 0.63 | 0.72 |
| SVM | 0.88 | 0.45 | 0.60 |
| Random Forest | 0.84 | 0.74 | 0.79 |

In order to analyze the impact of each group of feature on the classification task in the first run, we used feature removal method and the following figure shows the feature importance analysis:



- **Step 2:** As was mentioned previously, since sarcasm is dependent on context, it will become very subjective to just do the binary classification and say if the text is sarcastic or not. As a result, once the binary classification is completed in the first phase, we opted to assign each post a degree of sarcasm; Hence, we can compare and sort all posts according to their sarcastic level. These results allowed to clearly identify on scale the level of sarcasm of each post; A high score indicates the post is very sarcastic, while a low score indicates the post is non sarcastic. The following is the result of the experiment:

|  | Average Score | Score >0.71 | Score <0.48 |
|---|---|---|---|
| Sarcastic posts | 0.71 | 51% | 14% |
| Non-sarcastic posts | 0.48 | 13% | 77% |

- As is clear from the table above, the average score for sarcastic posts is much higher. Also, 77 percent of the non-sarcastic texts are below the sarcastic average which is expected. Furthermore, 51 percent of sarcastic texts have a score above 0.71 and 49 percent are below 0.71. This proves that sarcastic texts have a varying degree of sarcasm.