## Emotion Recognition in Arabic Tweets

Mariam ElOraby<sup>1</sup>, Rawan Reda<sup>2</sup>, and Salma El-Shafey<sup>3</sup>

 $^{1}46-3294$   $^{2}46-13489$   $^{3}46-17002$ 

May 6, 2023

#### Abstract

While sentiment analysis of Arabic tweets has been studied, emotion recognition, which identifies specific emotions expressed in the text, has received less attention. This project aims to explore the field of emotion recognition of Arabic tweets. The proposed approach will leverage BERT-based models, to detect and classify seven emotions: joy, sadness, anger, fear, love, sympathy, and surprise, or none. However, working with Arabic tweets presents several challenges due to the complexity and uniqueness of the Arabic language, including the lack of standardization in spelling and punctuation, the presence of dialects and regional variations, and the rich vocabulary and complex syntax. By preparing the data in an attempt to these gaps and challenges, this project aims to provide valuable insights into the emotions and attitudes of Arabic-speaking users on social media platforms.

### 1 Introduction

Social media platforms have become a rich source of information for studying human behavior and emotions. With millions of users expressing their thoughts and feelings on Twitter, sentiment and emotion analysis of tweets have become a popular area of research. However, most of the existing studies on emotion recognition of tweets have been focused on English tweets, while the analysis of other languages is relatively less explored. Arabic is one of the most widely used languages on social media platforms <sup>1</sup>, and the analysis of Arabic tweets can provide valuable insights into the emotions and attitudes of Arabic-speaking users.

Our project aims to explore the field of emotion recognition of Arabic tweets. Emotion recognition refers to the automatic detection and classification of emotions expressed in the text. The difference between sentiment analysis and emotion recognition is that sentiment analysis focuses on identifying the polarity of a text, whether it is positive, negative, or neutral, while emotion recognition aims to identify the specific emotion expressed in the text. We are focusing on the recognition of seven emotions: joy, sadness, anger, fear, love, sympathy, and surprise, and none.

Despite the growing interest in emotion recognition of tweets, there are still several gaps in the existing literature, particularly in the analysis of Arabic tweets. Most of the existing studies on Arabic tweets have focused on sentiment analysis rather than emotion recognition.

Our motivation is to address these gaps and contribute to the field of emotion recognition of Arabic tweets by proposing an automatic approach that can detect and classify emotions expressed in Arabic tweets. Our proposed approach will leverage deep learning techniques, including several transformer-based models (BERT-based models in particular) to capture the semantic and contextual information of the text. The models will be evaluated on an Arabic tweet dataset to assess their effectiveness and robustness.

Working with Arabic tweets datasets in general has several challenges due to the complexity and uniqueness of the Arabic language. One of the significant challenges is the lack of standardization in spelling and punctuation in Arabic tweets, specially in dialectal Arabic. It is a highly inflected

<sup>&</sup>lt;sup>1</sup>https://www.extradigital.co.uk/articles/social-media/arabic-social-media.html

language, meaning that words can have multiple forms depending on their position in the sentence and their relationship with other words. Therefore, Arabic tweets may contain misspellings, abbreviations, or variations in spelling and punctuation, making it challenging to preprocess and normalize the data. Additionally, Arabic language has a rich vocabulary and complex syntax, making it challenging to develop effective models for emotion recognition. Another challenge is the presence of dialects and regional variations in Arabic tweets, which can lead to variations in meaning and expression of emotions.

## 2 Background

Transformers are a state-of-the-art model in natural language processing (NLP) that have shown remarkable success in various tasks, including sentiment analysis and emotion recognition. These models, such as BERT, use self-attention mechanisms to capture contextual information and have been shown to outperform traditional NLP models on a range of NLP tasks [DCLT18]. One of the main advantages of transformers is their ability to learn from large amounts of data and transfer knowledge across different tasks, making them a powerful tool for emotion recognition in Arabic tweets. By leveraging pre-trained transformer-based models, we can achieve high accuracy in recognizing emotions in Arabic tweets, even when dealing with challenges such as non-standard spelling and regional variations. Overall, transformers represent a promising approach for advancing Arabic NLP research and improving our understanding of people's emotions and attitudes in the Arab world.

In this project we use BERT-based models to recognize emotions in Arabic tweets. BERT (Bidirectional Encoder Representations from Transformers) is a pre-trained transformer model that has achieved state-of-the-art results on a variety of natural language processing tasks. Unlike traditional NLP models that process text in a linear manner, BERT uses a multi-layer bidirectional transformer encoder that captures contextual information from both left and right contexts of the input text. Additionally, BERT utilizes a self-attention mechanism, which allows it to focus more on important parts of the input text and less on irrelevant information, leading to improved performance on tasks that require understanding long-range dependencies in the input text. The aforementioned properties allow BERT to learn rich semantic representations of text that can be fine-tuned for various downstream tasks, including sentiment analysis, question answering, and language translation. BERT-based models are pre-trained on massive amounts of data, and can be fine-tuned on specific tasks with much smaller labeled datasets. By leveraging BERT's pre-trained knowledge, we can achieve higher accuracy on several Arabic NLP tasks such as emotion recognition or classification.

## 3 Dataset

We use a publicly available dataset on HuggingFace<sup>2</sup>. The collection of the dataset went through several phases to make it balanced because emotions such as fear and love were not very common in text unlike happiness and sadness [AKEB18]. The authors of the dataset conducted some standard preprocessing steps on the tweets, such as normalizing Arabic characters, removing diacritics, and eliminating links, mentions, and Retweet indicators ("RT"). Then the processed dataset was made available publicly for research purposes.

#### 3.1 Analysis

The dataset consists of 10,065 tweets, mostly tweets written in the Egyptian dialect. The collected tweets spanned 8 emotions: neutral/none, joy, sadness, anger, sympathy, fear, love and surprise, with counts as shown in Table 1

We wanted to analyse the most common words of each emotion category, but the dataset was noisy as expected. Thus, we first focused on cleaning and preprocessing it. The cleaning process began with

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/datasets/emotone\_ar

Emotion Label	Count
none	1550
anger	1444
joy	1281
sadness	1256
love	1220
sympathy	1062
surprise	1045
fear	1207

Table 1: Breakdown of emotions in the dataset. From: [AKEB18]

قاب	ضحك	سعادة	حزن	مفاجأة	محرج	ملل	غضب
♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥ ♥	<b>ॐ</b> :P:- P	© ← ← ⊕ ← ⊕ ← ⊕ ← ⊕ ← ⊕ ← ⊕ ← ⊕ ←		◎ 😱 🥸	© <b>©</b> <u>R</u>	<b>□</b> ♥ ♥	

Table 2: Emojis and the words they were replaced with.

removing any repeated characters like the "ين in "إيسيه", this is to reduce input size and an attempt to normalize the dataset further. We further removed stopwords, punctuations, and replaced some emojis with words that conveys the same mood of the emoji, and removed the rest. The complete list of emojis and emoticons considered in our project are shown in table 2. We also removed any other weird patterns through their uni-codes. After this cleaning process, we tallied the most frequently occurring words within each emotion category. This is presented in Figure 1. The counts in more details are documented in Table 3 in the Appendix.

It is worth noting that the dataset authors pointed out that a section of the dataset was obtained by filtering tweets from Egypt between July 31st, 2016 and August 20th, 2016 using the search term "أولبياد" (Olympics). This was because during that period, people were anticipated to tweet about the Olympics using a wide range of emotions, as explained in the dataset paper. Therefore, the frequent appearance of the term in the word clouds can be attributed to this.

## 4 Solution Approach

Our proposed method begins with some slight preprocessing and cleaning on the tweets in the dataset to eliminate any noise or special characters. We won't need to remove the stopwords here because they provide valuable contextual information, such as negation words.

We then proceed to tokenize the cleaned tweets using BERT's tokenizer and fine-tune two or three BERT-based models for the task of emotion classification. After training, we evaluate the performance of the classification models and compare their results. Additionally, we experiment with different preprocessing techniques to observe their impact on the accuracy of the models.

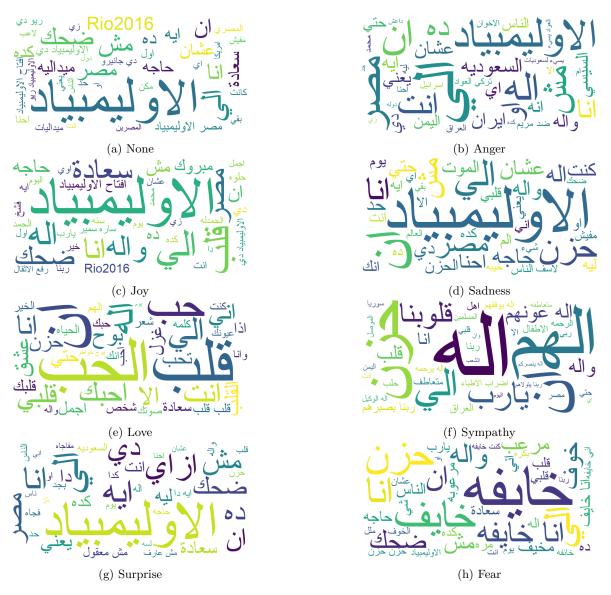


Figure 1: Visualizing emotions: wordclouds of the most frequently used words.

# References

- [AKEB18] Amr Al-Khatib and Samhaa R. El-Beltagy. Emotional tone detection in arabic tweets. In Alexander Gelbukh, editor, *Computational Linguistics and Intelligent Text Processing*, pages 105–114. Springer International Publishing, 2018.
- [DCLT18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. CoRR, abs/1810.04805, 2018.

# Appendix

Emotion Label	Most Common Words and their Counts
None	1512 : الاوليمبياد
	933 : الي
	: مصر
	: 166
	: دي : 156
	: 148 ضحك
	: 132
	UI: 123
	84 : ده
	83 : سعادة
Anger	: 1لي : 207
	اله : 156
	144 : الاوليمبياد
	131 : ان
	: 118 مصر
	110 : مش
	81 : انت
	os: 72
	72 : السعوديه
	UI: 68
Joy	392 : الأوليمبياد
	ال : 124
	114 : الي
	98 : قلب
	86 : سعادة
	78 : مصر
	UI: 77
	73 : دي
	69 : مبروك
- C 1	: 64 واله
Sadness	317 : الأوليمبياد
	: 200
	161 : الي
	: 159 - 107
	UI: 127
	: 113 مش
	: 103 مصر دار دوم
	ચી∶ 69
	: 55 عاجه

	55 : عشان
Love	: 352 علب
	الحب: 296
	: 1ن : 280
	: 131
	اله : 112
	93 : الى
	79 : بو <del>ح</del>
	15: 75: انت
	<b>Y</b> 1: 72
	UI: 72
Sympathy	اله : 485
	130 : الهم
	125 : <b>ح</b> زن
	106 : ان
	93 : الى
	91 : ربنا
	: 83 علب
	80 : يارب
	ს1: 79
	77 : قلو بنا
Surprise	302 : الاوليمبياد
	171 : مش
	142 : الي
	ايه : 136
	ป: 128
	89 : ازاي
	: 82
	81 : دي
	78 : دا
	73 : مصر
Fear	675 : خايفه
	انا : 318
	: 236 خايف
	228 : حزن
	179 : الي
	ان: 106
	94 : واله
	88 : خوف
	. 86 مش
	84 : اني

Table 3: Most frequent 10 words per each emotion category, and their corresponding counts.