# Twitter Sentiment and Emotion Analysis

**Oihane Cantero**
UPV/EHU
email@domain

**Julen Etxaniz**
UPV/EHU
email@domain

**Jose Javier Saiz**
UPV/EHU
email@domain

## Abstract

This corpus is a monolingual resource of unique tweets about product reviews, with each tweet paired with an adversarial sentence. Both original and adversarial tweets are manually and automatically annotated according to sentiment polarity (positive, negative, neutral). The objective is to compare and evaluate both types of annotation.

## 1   Introduction

This corpus is a monolingual resource of unique English tweets about technology, with each tweet paired with an adversarial sentence. Both original and adversarial tweets are manually and automatically annotated according to sentiment polarity (positive, negative, neutral) and emotions (joy, anger, sadness, optimism). For automatic annotation, the TweetEval (Barbieri et al., 2020) evaluation framework was used, which covers both sentiment polarity and stance detection for Twitter-specific data, among other classification tasks. The objective is to compare and evaluate both types of annotation.

The aim of this resource is to compare the annotations of tweets about technology and their adversarial sentences. It will also be interesting to compare these annotations when we do it by hand or automatically. The results will show whether the models are affected in their performance on classification tasks when faced with adversarial samples. If performance is significantly compromised by this phenomenon, improvements in the criteria used to generate the models may be proposed.

## 2   Related work

## 3   Material and methods

### 3.1   Data statement

A data statement is a document which serves as the contextualization and description of a data set, usually within the Natural Language Processing context (Bender and Friedman, 2018). It was proposed to make the data sets more understandable for any user and overall, to avoid bad practices within the research and professional practice. In the following sections, we are going to cover the data statement for this corpus so to explain the decisions and reasons behind our work.

The resource and code are publicly available at GitHub[1]. The code is licensed under the MIT open source license. All non-code materials provided are made available under the terms of the CC BY 4.0 license (Creative Commons Attribution 4.0 International license).

### A. CURATION RATIONALE

The corpus is a monolingual test set intended to evaluate text classification models for their ability to estimate sentiment and emotion. This is based on their performance against regular documents and adversarial phenomena. It contains a sample of 140 tweets about technology obtained from the Twitter API via the Tweepy package. We exclude retweets, quotes and replies to get better tweets. We also exclude tweets that have links to reduce spam tweets. The exact search query used is: *context:65.848920371311001600 lang:en -is:retweet -is:quote -is:reply -has:links*. The goal for using these parameters to collect the text samples was to obtain as many opinions with pronounced sentiments as possible, so that the manual annotation would be as accurate and homogeneous as possible.

[1] https://github.com/juletx/twitter-sentiment

## B. LANGUAGE VARIETY

Data extraction was made using the IETF language tag for English *en* (ISO 639-1). This language label does not refer to any geographic region or specific variety of English, but since the demographic group of speakers is mostly from the United States, as explained in the section 3.1-C, the language code is *en-US* (ISO 639-1/ ISO 3166-1 alpha-2), which refers to all American English varieties.

## C. SPEAKER DEMOGRAPHIC

Tweets are included upon language and semantic content. Any other type of information or demographic data cannot be controlled for selection, which means that this corpus is not based on specific speakers, but on randomly chosen instances from a larger collection. Therefore, we must rely on general Twitter demographics, which include a majority of speakers which are male (70.4%) and ranges between 18 and 49 years old (78.7%). Tweet metadata sometimes provides location information, which for our sample suggests that a significant number of the instances originate in the United States. This information is relevant for determining the prevailing language variety in the corpus.

## D. ANNOTATOR DEMOGRAPHIC

The annotation was made by three Spanish MSc students with training in linguistics and computer science.

## E. SPEECH SITUATION

Tweets are characterized by being short, spontaneous and conversational in nature. That is, opinions are concise and sentiments are condensed into a few words, which is important for the corpus because it guarantees that the annotation will be carried out and analyzed accurately. Tweets are also time-sensitive texts, which means that the content the instances is only relevant to the time in which they were written.

## F. TEXT CHARACTERISTICS

The sample includes tweets made before February 1. The Twitter context "Interests and Hobbies Vertical: Technology" (domain and entity ID: 65.848920371311001600) was used as the thematic parameter. This context includes "top-level interests and hobbies groupings" about "technology and computing", as indicated in the description of the metadata context.

## 3.2 Guidelines

We defined very simple initial guidelines before starting the annotation process. The tweets were annotated according to the sentiment (positive, negative, neutral) and emotion (joy, anger, sadness, optimism) that was inferred by the semantic content. Only one label per tweet is allowed for both tasks. We read tweets one by one and annotate both tags.

After the first annotation, we realized that we needed to update the guidelines to increase the annotator agreement. On the one hand, there were some misunderstandings with words and symbols. For example, &lt;3 is used as a heart and &gt; equals greater than. On the other hand, there were some were some cases where specific criteria had to be decided to select a label.

When annotating sentiments, some tweets seemed to be ironic or sarcastic, and it was difficult to decide a tag. Since irony makes the "intended meaning [...] appear on the surface to express the opposite" (Irony, 2022), we decided for those cases expressing irony or sarcasm to label the opposite sentiment to the overall sentiment when taken literally.

Another problem for which a clear guideline was needed was tweets with contradicting sentiments. For this case, we decided to keep the overall sentiment to avoid using too many neutral tags.

When there is no clear emotion, we decided to select the *joy* tag because it is the term with the broadest meaning among the possible options. Also, if there is uncertainty between the *joy* and *optimism* sentiments, we also decided to select *joy* for the same reason.

These guidelines are summarized and exemplified in Table 1.

## 3.3 Inter Annotator Agreement

Initially, a sample of the first 20 instances was shared among the annotators. The annotations for these initial tweets were used to calculate Inter Annotator Agreement. Table 2 shows examples of these instances.

Table 3 shows that there is moderate agreement and a low Kappa score in the first results. The agreement results in the sentiment task are homogeneous between pairs, while in the emotion task the results were more varied. This is explained by the task labels: the sentiment task uses labels which are more clear, while the labels in the emotion task are

| Problem | Guideline | Example | |
|---------|-----------|---------|---------|
| | | Sentence | Label |
| Irony/ Sarcasm | Use opposite sentiment to the literal one | "oh yeah tesla well what about a car that just logs into your tiktok acct and drives you to starbucks" | negative |
| Contradicting sentiments | Prefer overall sentiment over the "neutral" tag | "Finally managed to move my business email from google hosting to another host. So stressful and difficult. The whole internet is so hard! At least the metaverse is coming, I'm confident that will make everything easy and good, phew." | positive |
| Unclear emotion | Prefer "joy" tag for being broader in sense | "All the software I create will be free and open source, but that doesn't necessarily mean I won't write cryptic software for some of my projects" | joy |

Table 1: Defined guidelines after the Inter Annotator Agreement with corresponding examples.

more subjective and open to interpretation. This is main reason why the guidelines were updated, after which the annotation was repeated along the Inter Annotator Agreement calculation. The results for the annotations with the updated guidelines show a clear improvement in the emotion task.

Thereafter each annotator was given a different set of 40 tweets for annotation. The same methodology was followed, where the manual annotations were used to evaluate the accuracy of the models.

### 3.4 Automatic annotation

These tweets were also labelled using automatic models. Models trained for the TweetEval (Barbieri et al., 2020) evaluation framework were used. This framework covers many classification tasks that include sentiment and emotion classification.

The base RoBERTa (Liu et al., 2019) model was pretrained using 58M English tweets. Then each model is finetuned using specific data for that task. SemEval 2017 Subtask A data (Rosenthal et al., 2019) is used for sentiments and SemEval 2018 Affects in Tweets (Mohammad et al., 2018) is used for emotions.

The fine-tuned models are available in Hugging-Face and that makes them very easy to use. They only require a small preprocessing step to substitute all mentions with @user and links with http.

### 3.5 Adversarial examples

We will create some adversarial tweets automatically to measure the impact in the metrics of the models. The aim of these adversarial examples is to confuse the models while maintaining the manual label. Therefore, if we generate the adversarial examples correctly, they should have the same meaning and sense, so there is no need to annotate them. As these attacks can easily be replicated, we decided not to include them in the final resource.

We used the Adversary[2] library to generate 8 generic attacks. These attacks are not specifically prepared for our models, they are simple attacks that can be tested on any type of model. Words are selected with a probability of 0.3 and the selected attack is applied to those words. Attacks can also be combined, but we decide to apply them separately for simplicity. An example of each attack can be seen in Table 4. This tweet was manually annotated as negative sentiment and anger emotion.

## 4 Results

We will compare the results and extract some conclusions about the performance of the models. We will visualize the percentage of tweets from each class in the manual and automatic annotations. We will also visualize the most common words of each class in word clouds.

[2] https://github.com/airbnb/artificial-adversary

| Text | Sentiment | | | Emotion | | |
|---|---|---|---|---|---|---|
| | **julen** | **oihane** | **javier** | **julen** | **oihane** | **javier** |
| *just put a CD into my MacBook to burn it and my computer is literally trembling with reawakened recognition* | negative | positive | negative | anger | joy | joy |
| *oh yeah tesla well what about a car that just logs into your tiktok acct and drives you to starbucks* | neutral | positive | negative | joy | joy | anger |
| *#100DaysOfCode Haven't updated in a while due to not feeling well, just been reviewing some HTML/CSS &amp; JavaScript until I feel better to take on new concepts. Also been watching mock interviews :)* | negative | negative | negative | sadness | optimism | optimism |

Table 2: Instance examples used to calculate Inter Annotator Agreement.

| pairs | agreement_sentiment | kappa_sentiment | agreement_emotion | kappa_emotion |
|---|---|---|---|---|
| julen oihane | 70.0 → **65.0** | 53.8 → **45.1** | 55.0 → **70.0** | 33.3 → **49.2** |
| julen javier | 60.0 → **70.0** | 40.1 → **49.4** | 40.0 → **70.0** | 13.0 → **46.9** |
| javier oihane | 60.0 → **55.0** | 39.2 → **32.3** | 65.0 → **70.0** | 46.2 → **50.0** |
| average | 63.3 → **63.3** | 44.4 → **42.3** | 53.3 → **70.0** | 30.8 → **48.7** |

Table 3: Inter Annotator Agreement result improvements after updating the guidelines. Final results are in bold.

## 4.1 F1 scores

F-score was calculated for each label and macro-average for all labels in each task in order to compare the model results against the manual annotations. As can be seen in Table 5, the results with adversarial sentences output lower F-scores than with the original sentences. For the sentiment task, the performance with original sentences delivers a 65.8 F-score against a 61.7 mean F-score with the adversarial sentences. The same applies to the emotion task results, which bring a 53.3 F-score with original sentences and a 50.7 mean F-score with the adversarial sentences. Even if in overall results for adversarial sentences are worse than the original ones, we can remark that some of the adversarial attacks doesn't decrease F-score, they even improve it in some cases. The best results for sentiment analysis are for the *swap_words* attack and this same attack performs quite well in emotion detection. And for emotion detection the best results are for *remove_spacing*, that doesn't receive so bad results for sentiment analysis. So watching

these results, we may conclude that removing a space of the sentence or word order doesn't really matter for these tasks.

In Figure 1, we see the results of Table 5 into barplots, we see the score differences between each tag. The tag that is classified worse for both original and adversarial sentences is *positive*, and for emotions, we see that *anger* and *joy* are the tags that get better results, and *optimism* is clearly below all the others.

In Figure 2, we can observe that between the two tasks, sentiment analysis get better results for both original and adversarial sentences.

And in Figure 3, we can see the difference in scores between the two tasks and all the adversarial attacks.

## 4.2 Pieplots

In Figure 4 we can see the distribution of each of the tags of the sentiment analysis task. In the manual distribution, positive and negative tags are equally distributed, and are less neutral tags. In the

| Type | Sentence | Label | |
|---|---|---|---|
| | | Emotion | Sentiment |
| Original | having no sort of WiFi actually sucks | negative | sadness |
| Swap words | having no sort WiFi of actually sucks | negative | sadness |
| Remove space | having-no'sort of WiFi actually sucks | negative | anger |
| Replace letters with symbols | having no $or7 of WiFi actually $u{[}kS | neutral | sadness |
| Swap letters | having no srot of WiFi actually sucks | negative | anger |
| Insert punctuation | having no sor]{t of WiFi actually ]su!cks | neutral | sadness |
| Insert duplicate characters | having no sorttt of WiFi actually sucks | negative | sadness |
| Delete characters | having no sot of WiFi actually sucks | negative | sadness |
| Change case | having no SORT of WiFI actually sucks | negative | anger |

Table 4: Examples of adversarial sentences and the predicted sentiment and emotion.

| text | neg | neu | pos | sen | ang | joy | opt | sad | emo |
|---|---|---|---|---|---|---|---|---|---|
| original | 71.0 | 71.0 | 70.3 | 65.8 | 66.7 | 66.1 | 27.9 | 52.6 | 53.3 |
| swap_words | **73.1** | **73.1** | **71.0** | **67.2** | 67.6 | 64.6 | 27.9 | 51.3 | 52.8 |
| remove_spacing | 68.2 | 68.2 | 60.2 | 61.2 | **69.6** | **68.2** | 33.3 | **55.0** | **56.5** |
| letter_to_symbol | 65.2 | 65.2 | 58.5 | 59.0 | 62.5 | 62.0 | 22.9 | 38.9 | 46.6 |
| swap_letters | 69.0 | 69.0 | 64.3 | 62.8 | 67.6 | 64.1 | 21.6 | 31.6 | 46.2 |
| insert_punctuation | 62.8 | 62.8 | 55.7 | 57.5 | 62.5 | 74.0 | 20.0 | 35.0 | 47.9 |
| insert_duplicate_characters | 72.5 | 72.5 | 63.6 | 63.9 | 66.7 | 66.7 | 26.3 | 47.4 | 51.8 |
| delete_characters | 68.2 | 68.2 | 59.5 | 59.9 | 68.5 | 66.7 | 27.8 | 41.0 | 51.0 |
| change_case | 71.1 | 71.1 | 63.6 | 62.6 | 63.2 | 65.6 | **37.2** | 44.4 | 52.6 |

Table 5: F1 scores for original text and adversarial attacks. F1 for each label and macro-average of all labels for sentiment and emotion.

automatic annotation, the distribution of positive and negative tags are also equally distributed, but the majority class is neutral. This can be because while manually annotating, we tend to avoid neutrality, and we try to put a positive or negative tag at all costs. Another reason may be that we decided to take the overall sentiment when there are positive and negative sentiments in the same tweets, and maybe the model makes an average and classifies them as neutral.

In Figure 5, we can see that something similar appends. In our guidelines, we decided that we will put the joy tag when we will have doubts. This increased our manual annotations' joy tags number. For optimism and anger, the results are quite similar but for anger, because the automatic annotation puts more than the manual one. This should be because our guidelines implies that there will be more joy, because we put it when there is no emotion and when we had doubts between joy and optimism, as we think that it is the more general emotion.

## 4.3 Wordclouds

## 5 Conclusions

Classification tasks in Twitter related to sentiment analysis have grown in importance over the last years, and with an increasing number of models and frameworks created around this tasks, it is worth evaluating the implementation of these methods.

Our evaluation focuses on the outcomes of a popular Twitter-specific framework across two tasks by comparing its results to human annotation. This evaluation method is useful for the purpose of this work because it allows us to examine the broader potential impact of these models and whether the guidelines used to create the training corpus for these tasks translate into adequate performance on sentiment analysis tasks.

We saw that good guidelines are essential to annotate manually, even if the task seems to be "easy". They allow annotators to annotate the same way, and are very important if we want to create a coherent resource.

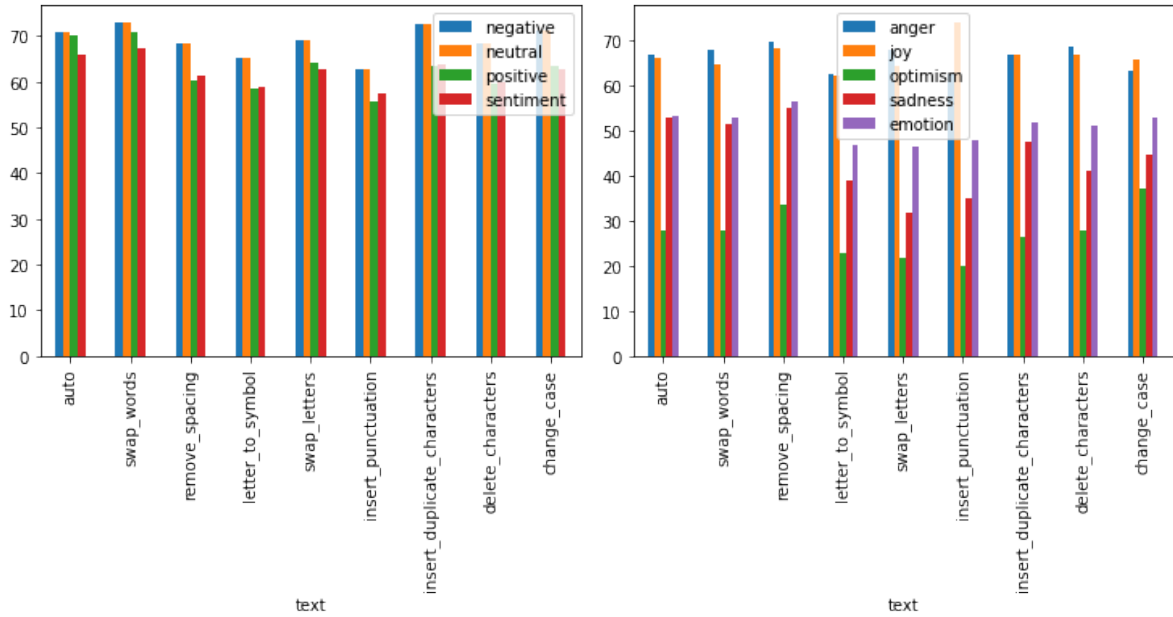The outcomes of this work have shown that clas-

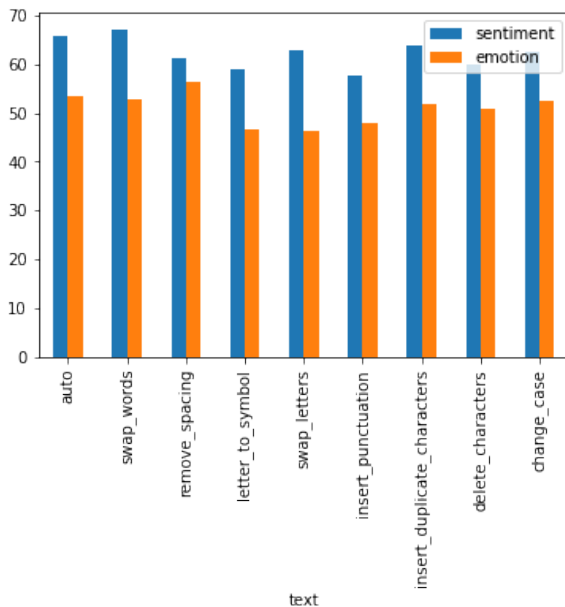Figure 1: F1 for each label and macro-average of all labels for sentiment and emotion.



Figure 2: Sentiment and emotion barplot.



Figure 3: Sentiment and emotion scatterplot.

and other classification frameworks with different guidelines.

sification models have multiple sources of error and tend to under-perform when facing different scenarios such as unknown words, lexicon gaps or unspecified sentence contexts. This shows there is room for improvement in this task. For future work, it would be useful to analyze the performance of the models in multiple other related tasks, such as irony/ sarcasm detection or hate detection, instead of a pair of tasks. Other improvements on this evaluation would be to explore additional languages
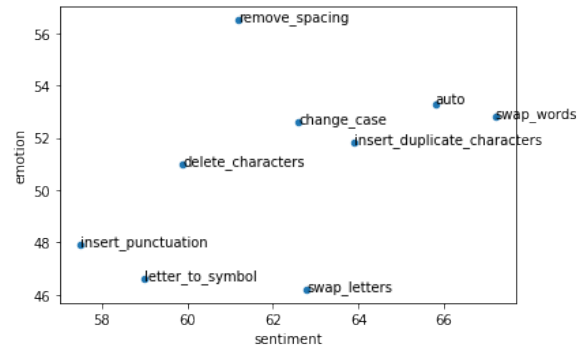
## References

Francesco Barbieri, Jose Camacho-Collados, Luis Espinosa Anke, and Leonardo Neves. 2020. TweetEval: Unified benchmark and comparative evaluation for tweet classification. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1644–1650, Online. Association for Computational Linguistics.

Emily M. Bender and Batya Friedman. 2018. Data statements for natural language processing: Toward mitigating system bias and enabling better science. *Transactions of the Association for Computational Linguistics*, 6:587–604.

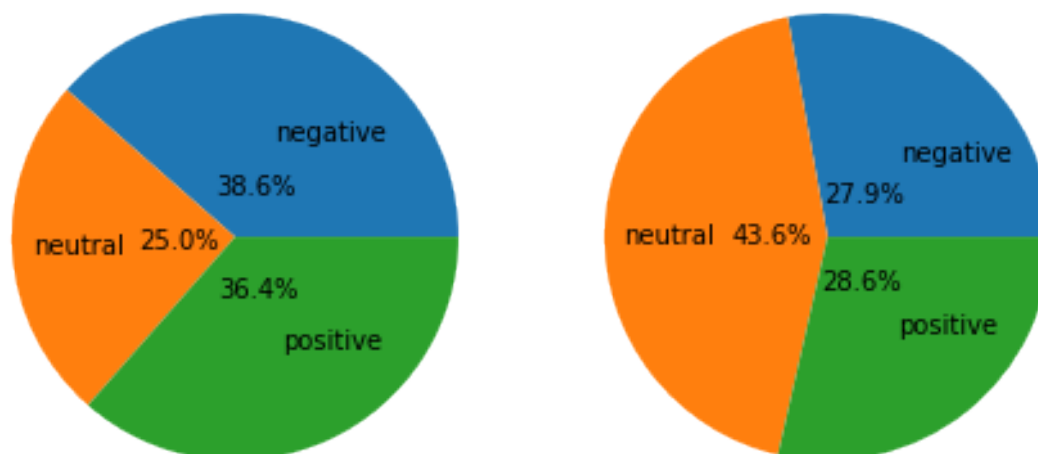Irony. 2022. *Oxford Reference*. Oxford University Press.

Figure 4: Sentiment label percentages for manual and automatic annotation.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. SemEval-2018 task 1: Affect in tweets. In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 1–17, New Orleans, Louisiana. Association for Computational Linguistics.

Sara Rosenthal, Noura Farra, and Preslav Nakov. 2019. Semeval-2017 task 4: Sentiment analysis in twitter. *arXiv preprint arXiv:1912.00741*.
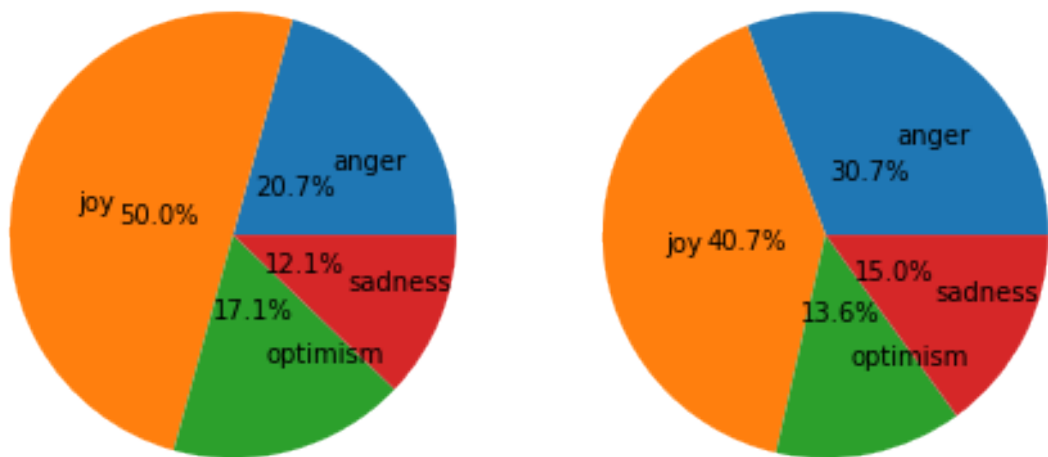
Figure 5: Sentiment label percentages for manual and automatic annotation.