

“EMPLOYABILITY A CORPUS OF PRE-ANNOTATED TWEETS IN THE EFFACACIOUS USE OF DATA MINING SENTIMENTS FROM SOCIAL MEDIA”

SHUBHAM CHANDNA

BAL BHARTI PUBLIC SCHOOL, PITAMPURA, NEW DELHI

ABSTRACT

Twitter could be a small blogging web site, wherever users will post messages in the short text referred to as Tweets. Tweets contain user opinion associate degreed sentiment towards an object or person. This sentiment data is incredibly helpful in numerous aspects for business and governments. During this paper, we tend to gift a way that performs the task of tweet sentiment identification employing a corpus of pre-annotated tweets. We tend to gift a sentiment grading operate that uses previous data to classify (binary classification) and weight many sentiment-bearing words/phrases in tweets. Victimization this grading operate we tend to succeed classification accuracy of eighty-seven on Stanford Dataset and half of 1 mile on Mejjaj dataset. Victimization supervised machine learning approach; we manage to achieve a classification accuracy of half of 1 mile on Stanford dataset.

1. BACKGROUND

With a massive increase in new technologies, a variety of individuals expressing their views and opinions via net square measure increasing. This data is incredibly helpful for businesses, governments, and people. With over 340+ million Tweets (short text messages) per day, Twitter is changing into a significant supply of knowledge. Twitter could be a micro-blogging website, that is in style owing to its short text messages popularly called "Tweets." Tweets have a limit of one hundred forty characters. Twitter features a user base of 140+ million active users¹ 1As on March twenty-one, 2012. Source: <http://en.wikipedia.org/wiki/Twitter> and so could be a necessary supply of knowledge. Users usually discuss current affairs and share their personals views on numerous subjects via tweets. Out of all the favored social media's like Facebook, Google+, Myspace, and Twitter, we elect Twitter as a result of 1) tweets square measure tiny long, so less ambiguous; 2) unbiased; 3) square measure simply accessible via API; 4) from various socio-cultural domains. During this paper, we tend to introduce associate degree approach which may be accustomed notice the opinion in associate degree aggregate assortment of tweets. During this approach, we tend to used two completely different datasets that square measure build victimization emoticons and list of suggestive words severally as clangorous labels. We tend to provide a new technique of grading

"Popularity Score," that permits determination of the recognition score at the extent of individual words of a tweet text. We tend to conjointly stress on numerous varieties and levels of pre-processing needed for higher performance. Roadmap for the remainder of the paper: connected work is mentioned in Section a pair of. In Section three, we tend to describe our approach to handle the matter of Twitter sentiment classification at the side of pre-processing steps. Datasets utilized in this analysis square measure mentioned in Section four. Experiments and Results square measure bestowed in Section five. InSection half dozen, we tend to lift the feature vector approach to Twitter sentiment classification. Section seven presents a discussion on the ways and that we conclude the paper with later adding Section eight.

2. RELATED WORK

Research in Sentiment Analysis of user-generated content may be categorized into Reviews (Turney, 2002; Pang et al., 2002; Hu and Liu, 2004), Blogs (Draya et al., 2009; Chesley, 2006; He et al., 2008), News (Godbole et al., 2007), etc. of these classes touch upon giant text. On the opposite hand, Tweets square measure shorter length text and square measure tough to analyses owing to its distinctive language and structure. (Turney, 2002) Worked on product reviews. Turney used adjectives and adverbs for playacting opinion classification on reviews. He used PMI-IR formula to estimate the linguistics orientation of the sentiment phrase. He achieved a mean accuracy of seventy-four on 410 reviews of various domains collected from Opinion. (Hu and Liu, 2004) Performed feature primarily based sentiment analysis. Victimization Noun-Noun phrases they knew the options of the merchandise and determined the sentiment orientation towards every element. (Pang et al., 2002) Tested various machine learning algorithms on flick Reviews. He achieved eighty-inaccuracies in unigram presence feature assault Naive Thomas Bayes classifier. (Draya et al., 2009) Tried to spot domain specific adjectives to perform web log sentiment analysis. They thought about the fact that opinions square measure principally expressed by articles and pre-defined lexicons fail to locate domain data. (Chesley, 2006) Performed topic and genre freelance weblog classification, creating novel use of linguisticoptions.

Every post from the weblog is classed as positive, negative and objective. To the simplest of our information, there's less quantity of labor tired twitter sentiment analysis. (Go et al., 2009) Performed sentiment analysis on Twitter. They knew the tweet polarity victimization emoticons as clangorous labels and picked up a coaching dataset of one.6 million tweets. They according to associate degree accuracy of eighty-one.34% for his or her Naive Thomas Bayes classifier. (Davidov et al., 2010) Used fifty hashtags and fifteen emoticons as clangorous labels to make a dataset for twitter sentiment classification. They evaluate the result of various sorts of options for sentiment extraction. (Diakopoulos and Shamma, 2010) worked on political tweets to spot the

final sentiments of the folks on first U.S. presidential dialogue in 2008. (Bora, 2012) Conjointly created their dataset supported clangorous labels. They created a listing of forty words (positive and negative) that we're accustomed to determining the polarity of the tweet. They used a mix of a minimum word frequency threshold and Categorical Proportional distinction as a feature choice technique and achieved the best accuracy of eighty-three.33% on a hand-labeled check dataset. (Agarwal et al., 2011) Performed three categories (positive, negative and neutral) classification of tweets. They collected their dataset victimization Twitter stream API and asked human judges to annotate the information into three categories. They'd 1709 tweets of every class creating a complete of 5127 altogether. In their analysis, they introduced POS-specific previous polarity options at the side of twitter specific options. They achieved soap accuracy of seventy-five.39% for unigram + senti options. Our work uses (Go et al., 2009) and (Bora, 2012) datasets for this analysis. We tend to use Naive Thomas Bayes technique to choose the polarity of tokens within the tweets. At the side of that, we offer associate degree helpful insight on however preprocessing ought to be done on a tweet. Our technique of Senti Feature Identification and recognition Score perform well on each the datasets. In the feature vector approach, we tend to show the contribution of personal informatics and Twitter specific options. Three Approach Our approach may be divided into numerous steps. Every one of those steps squares measure freelances of the opposite however necessary at the identical time.

3.1 Baseline

In the baseline approach, we tend to initial clean the tweets. We tend to take away all the individual characters, targets (@), hashtags (#), URLs, emoticons, etc. and learn the positive & negative frequencies of unigrams in coaching. Each unigram token is given two likelihood scores: Positive likelihood (Pp) and Negative likelihood (Np) (Refer to Equation 1). We tend to follow the same cleanup method for the check tweets. Once cleanup the check tweets, we tend to type all the potential unigrams and check for his or her frequencies within the coaching model. We tend to add up the positive and negative likelihood innumerable all the constituent unigrams and use their distinction (positive-negative) to search out the general score of the tweet. If the tweet score is >zero, then it's positive otherwise contrary.

$$\begin{aligned}P_f &= \text{Frequency in Positive Training Set} \\N_f &= \text{Frequency in Negative Training Set} \\P_p &= \text{Positive Probability of the token.} \\&= P_f / (P_f + N_f) \\N_p &= \text{Negative Probability of the token.} \\&= N_f / (P_f + N_f)\end{aligned}$$

3.2 Emoticons and Punctuations Handling

We build slight changes within the pre-processing module for handling emoticons and punctuations. We tend to use the emoticons list provided by (Agarwal et al., 2011) in their analysis. This list² is constructed from Wikipedia list of emoticons³ and is hand labeled into five categories (extremely positive, positive, neutral, negative and intensely negative). during this experiment, we tend to replace all the emoticons that square measure labeled positive or extraordinarily positive with 'zzhappyzz' and rest all alternative emoticons with 'zzsadzz.' We tend to append and prepend 'zz' too happy and unhappy to stop them from the mixture into tweet text. In the end, 'zzhappyzz' is scored +1 and 'zzsadzz' is scored -1. Exclamation marks (!) and question marks (?) conjointly carry some sentiment. In general, '!' is employed after we have to be compelled to stress on a positive word and '?' is employed to focus on the state of confusion or disagreement. we tend to replace all the occurrences of '!' with 'zzexclaimzz' and of '?' with 'zzquestzz.' We add 0.1 to the overall tweet score for every '!' And take off zero.1 from the overall tweet score for every '?'. 0.1 is chosen by trial and error technique.

3.3 Stemming

We use Porter Stemmer⁴ to stem the tweet words. We tend to modify porter stemmer and limit it to step one solely. Step one gets eliminate plurals and -ed or -ing.

3.4 Stop Word Removal

Stop words assume a negative job in the errand of opinion order. Stop words happen in both positive and negative preparing set, in this way including greater vagueness in the model arrangement. And furthermore, don't convey any assessment data and subsequently are of no utilization to us. We make a rundown of stop words like he, she, at, on, a, the, and so forth and disregard them while scoring. We additionally dispose of words which are of length ≤ 2 for scoring the tweet to amend structure and spelling. Spell remedy is an essential part in conclusion examination of client created content. Clients compose certain characters' self-assertive number

of times to put more accentuation on that. We utilize the spell remedy calculation from (Bora, 2012). In their calculation, they supplant a word with any character rehashing more than twice with two words, one in which the rehashed character is put once and second in which the rehashed character is put twice. For instance, the word 'swwwееееtttt' is supplanted with eight words 'swet,' 'swwet,' 'sweet,' 'swett,' 'sweet,' et cetera. Another necessary sort of spelling botches happens as a result of avoiding some of the characters from the spelling. like "there" is for the most part composed as "thr." Such sorts of spelling botches are not right now dealt with by our framework. We propose to utilize phonetic level spell amendment strategy in future.

3.6 Senti features

At this progression, we attempt to decrease the impact of non-feeling bearing tokens on our order framework. In the standard technique, we considered all the unigram tokens similarly and of positive and negative words. We utilize the rundown of most usually employed positive and negative words given by Twitrratr5. When we go over a token in this rundown, rather than scoring it using the Naïve Bayes recipe (Refer Equation 1), we score the token +/- 1 relying upon the outline in which it exists. Every one of the symbols which are absent from this rundown went under stage 3.3, 3.4, 3.5 and were checked for their event after each progression.

3.7 Noun identification

In the wake of doing every one of the amendments (3.3 - 3.6) on a word, we take a gander at the decreased word if it is being changed over to a Noun or not. We distinguish the word as a Noun word by taking a gander at its grammatical feature tag in English WordNet (Miller, 1995). If the dominant part sense (most ordinarily utilized sense) of that word is Noun, we dispose of the word while scoring. Thing words don't convey feeling and in this manner are of no utilization in our trials.

3.8 Popularity Score

This scoring technique supports the scores of the most regularly utilized words, which are area particular. For instance, cheerful is utilized dominantly to express a positive assessment. In this strategy, we numerous its prevalence factor (pF) to the score of each unigram token which has scored in the past advances. We utilize the event recurrence of a token in the positive and negative dataset to settle on the heaviness of ubiquity score. Condition 2 indicates how the prominence factor is ascertained for every symbol. We chose an edge 0.01 min to bolster as the cut-off criteria and diminished it significantly at each level. Support of a word is characterized as the extent of tweets in the dataset which contain this token. The esteem 0.01 is picked to such an

extent that we cover a substantial number of tokens without missing critical tokens, in the meantime pruning less successive symbols.

$$\begin{aligned}
 P_f &= \text{Frequency in Positive Training Set} \\
 N_f &= \text{Frequency in Negative Training Set} \\
 & \text{if}(P_f - N_f > 1000) \\
 & \quad pF = 0.9; \\
 & \text{elseif}((P_f - N_f) > 500) \\
 & \quad pF = 0.8; \\
 & \text{elseif}((P_f - N_f) > 250) \\
 & \quad pF = 0.7; \\
 & \text{elseif}((P_f - N_f) > 100) \\
 & \quad pF = 0.5; \\
 & \text{elseif}((P_f - N_f < 50)) \\
 & \quad pF = 0.1;
 \end{aligned}$$

Figure 1 shows the flow of our approach.

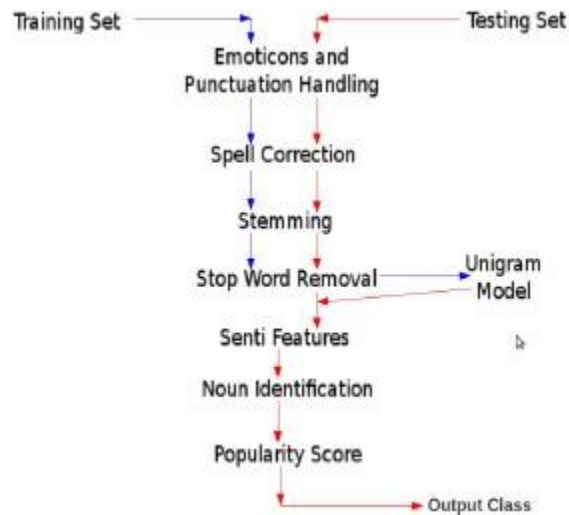


Figure 1: Flow Chart of our Algorithm

4. DATASETS

In this section, we explain the two datasets used in this research. Both of these datasets are built using noisy labels.

4.1 Stanford Dataset

This dataset (Go et al., 2009) was constructed consequently utilizing emojis as uproarious names. Every one of the tweets which contain ':')' was stamped positive and tweets containing ':(' were checked negative. Tweets that did not have any of these marks or had both were disposed of. The preparation dataset has ~1.6 million tweets, parallel number of positive and negative tweets. The preparation dataset was commented on into two classes (positive and negative) while the testing information was hand explained into three categories (positive, negative and nonpartisan). For our experimentation, we utilize just positive and negative class tweets from the testing dataset for our experimentation. Table 1 gives the points of interest of dataset.

Training Tweets	
Positive	800,000
Negative	800,000
Total	1,600,000
Testing Tweets	
Positive	180
Negative	180
Objective	138
Total	498

Table 1: Stanford Twitter Dataset

4.2 Mejaj

Mejaj dataset (Bora, 2012) was constructed utilizing full marks. They gathered an arrangement of 40 words and physically sorted them into positive and negative. They mark a tweet as confident if it contains any of the positive estimation words and as negative if it includes any of the contrary assumption words. Tweets which don't include any of these boisterous marks and tweets which have both positive and negative words were disposed of. Table 2 gives the rundown of words which were utilized as uproarious marks. This dataset contains just two class information. Table 3 presents the points of interest of the dataset.

Positive Labels	Negative Labels
amazed, amused, attracted, cheerful, delighted, elated, excited, festive, funny, hilarious, joyful, lively, loving, overjoyed, passion, pleasant, pleased, pleasure, thrilled, wonderful	annoyed, ashamed, awful, defeated, depressed, disappointed, discouraged, displeased, embarrassed, furious, gloomy, greedy, guilty, hurt, lonely, mad, miserable, shocked, unhappy, upset

Table 2: Noisy Labels for annotating Mejaj Dataset

Training Tweets	
Positive	668,975
Negative	795,661
Total	1,464,638
Testing Tweets	
Positive	198
Negative	204
Total	402

Table 3: Mejaj

5 EXPERIMENT

In this section, we explain the tests carried out using the above-proposed approach.

5.1 Stanford Dataset

On this dataset (Go et al., 2009), we play out a progression of trials. In the first arrangement of experiments, we prepare on the given qualifying information and test on the testing information. In the second arrangement of investigations, we perform 5-overlap cross approval utilizing the preparation information. Table 4 demonstrates the aftereffects of every one of these investigations

on steps which are clarified in Approach (Section 3). In table 4, we give results for

each progression emojis and accentuations dealing with, spell rectification, stemming and stop word expulsion referred to in Approach (Section 3). The Baseline + All Combined outcomes alludes to the blend of these means (emojis, accentuations, spell rectification, Stemming and stop word expulsion) performed together. Arrangement 2 results are normal of exactness of each crease.

5.2 Mejaj Dataset

The comparable provision of trials was performed on this dataset (Bora, 2012) as well. In the first arrangement of examinations, preparing and testing was done on the particularly given datasets. In the second arrangement of investigations, we perform 5-overlap cross approval on the preparation information. Table 5 demonstrates the aftereffects of every one of these trials. In table 5, we give results for each progression emojis and accentuations taking care of, spell remedy, stemming and stop word expulsion referred to in Approach (Section 3). The Baseline + All Combined outcomes allude to mix of these means (emojis, accentuations, spell revision, Stemming and stop word expulsion) performed together. Arrangement 2 results are normal of precision of each overlap.

5.3 Cross Dataset

To approve the vigor of our methodology, we tried different things with cross-dataset preparing and testing. We prepared our framework on one dataset and worked on the other dataset. Table 6 reports the consequences of cross-dataset assessments.

6	Feature	Vector	Approach

In this element vector approach, we shaped highlights utilizing Unigrams, Bigrams, Hashtags (#), Targets (@), Emoticons, Special Symbol (!) and used a semi-directed SVM classifier. Our component vector contained 11 highlights. We partition the highlights into two gatherings, NLP highlights, and Twitter particular highlights. NLP highlights incorporate recurrence of positive

Method	Series 1 (%)	Series 2 (%)
Baseline	78.8	80.1
Baseline + Emoticons + Punctuations	81.3	82.1
Baseline + Spell Correction	81.3	81.6
Baseline + Stemming	81.9	81.7
Baseline + Stop Word Removal	81.7	82.3
Baseline + All Combined (AC)	83.5	85.4
AC + Senti Features (wSF)	85.5	86.2
wSF + Noun Identification (wNI)	85.8	87.1
wNI + Popularity Score	87.2	88.4

Table 4: Results on Stanford Dataset

Method	Series 1 (%)	Series 2 (%)
Baseline	77.1	78.6
Baseline + Emoticons + Punctuations	80.3	80.4
Baseline + Spell Correction	80.1	80.0
Baseline + Stemming	79.1	79.7
Baseline + Stop Word Removal	80.2	81.7
Baseline + All Combined (AC)	82.9	84.1
AC + Senti Features (wSF)	86.8	87.3
wSF + Noun Identification (wNI)	87.6	88.2
wNI + Popularity Score	88.1	88.1

Table 5: Results on Mejaj Dataset

Method	Training Dataset	Testing Dataset	Accuracy
wNI + Popularity Score	Stanford	Mejaj	86.4%
wNI + Popularity Score	Mejaj	Stanford	84.7%

Table 6: Results on Cross Dataset evaluation

NLP	Unigram (f1) Bigram (f2)	# of positive and negative unigram # of positive and negative Bigram
Twitter Specific	Hashtags (f3) Emoticons (f4) URLs (f5) Targets (f6) Special Symbols (f7)	# of positive and negative hashtags # of positive and negative emoticons Binary Feature - presence of URLs Binary Feature - presence of Targets Binary Feature - presence of '!

Table 7: Features and Description

Feature Set	Accuracy (Stanford)
f1 + f2	85.34%
f3 + f4 + f7	53.77%
f3 + f4 + f5 + f6 + f7	60.12%
f1 + f2 + f3 + f4 + f7	85.89%
f1 + f2 + f3 + f4 + f5 + f6 + f7	87.64%

Table 8: Results of Feature Vector Classifier on Stanford Dataset

Unigrams matched, negative unigrams matched, positive bigrams matched, negative bigrams matched, etc., and Twitter specific features included Emoticons, Targets, HashTags, URLs, etc. Table 7 shows the features we have considered. HashTags polarity is decided based on the constituent words of the hashtags. Using the list of positive and negative words from Twitter, we try to find if hashtags contain any of these words. If so, we assign the polarity of that to the hashtag. For example, "#imsohappy" contains a positive word "happy," thus this hashtag is considered as the positive hashtag. We use the emoticons list provided by (Agarwal et al., 2011) in their research. This list is built from Wikipedia list of emoticons and is hand labeled into five classes (to a great degree positive, positive, nonpartisan, negative and a great degree negative). We decrease this five class rundown to two class by merging extremely positive class to single positive class and rest other classes (extremely negative, negative and neutral) to single negative class. Table 8 reports the accuracy of our machine learning classifier on Stanford dataset.

7 DISCUSSION

In this segment, we present a couple of models assessed utilizing our framework. The accompanying precedent indicates the impact of joining the commitment of emojis on tweet characterization. Precedent "Ahhh I can't move it, however, hello w/e it is on damnation I'm elated right now:- D." This tweet contains two conclusion words, "hellfire" and "elated." Utilizing the unigram scoring technique, this tweet is arranged unbiased however it is positive. If we consolidate the impact of emoji ":- D," at that point, this tweet is labeled positive. ":- D" is a solid positive emoji. Think about this model, "Bill Clinton Fail - Obama Win?". In this precedent, there is two estimation bearing words, "Fizzle" and "Win." In a perfect world, this tweet ought to be unbiased, yet this is labeled as a positive tweet in the dataset and utilizing our framework. In this tweet, if we ascertain the ubiquity factor (pF) for "Win" and "Fall flat," they turn out to be 0.9 and 0.8 individually. In light of the prevalence factor weight, the positive score dominates the negative score, and along these lines, the tweet is labeled as positive. It is critical to distinguish the setting stream in the content and furthermore how every one of these words alter or rely upon alternate expressions of the tweet. For computing the framework execution, we accept that the dataset which is utilized here is right. The vast majority of the occasions this suspicion is valid however there are a couple of situations where it comes up short. For instance, this tweet "My wrist still stings. I need to get it took a gander. I HATE the dr/dental practitioner/terrifying spots. :(Time to watch Eagle eye. On the off chance that you need to join, txt!" is labeled as positive, all things considered, this ought to have been labeled negative. Such mistaken tweets additionally impact the framework execution. There are a couple of constraints with the flow proposed approach which are additionally open research issues.

1. Spell Correction: In the above-proposed approach, we gave an answer for spell rectification which works just when the client enters additional characters. It comes up short when clients avoid a few characters like "there" is spelled as "thr." We propose the utilization of phonetic level spell rectification to deal with this issue.

2. Hashtag Segmentation: For taking care of hashtags, we searched for the presence of the positive or negative words⁹ in the hashtag. Be that as it may, there can be a few situations where it may not work effectively. For instance, "#thisisnotgood," in this hashtag on the off chance that we think about the nearness of positive and negative words, at that point this hashtag is labeled positive ("great"). We neglect to catch the nearness and impact of "not" or, in other words, hashtag as negative. We propose to devise and utilize some rationale to fragment the hashtags to get the right constituent words.

3. Setting Dependency: As talked about in one of the precedents above, even tweet content

which is constrained to 140 characters can have setting reliance. One conceivable technique to deliver this issue is to distinguish the articles in the tweet and after that discover the supposition towards those items.

8 CONCLUSION AND FUTURE WORK

Twitter notion examination is a vital and testing errand. Twitter is microblog experiences different semantic and syntactic blunders. In this exploration, we proposed a technique which fuses the ubiquity impact of words on tweet supposition characterization and furthermore accentuation on the most proficient method to preprocess the Twitter information for greatest data extraction out of the little substance. On the Stanford dataset, we accomplished 87% exactness utilizing the scoring technique and 88% employing SVM classifier. On Mejaj dataset, we demonstrated a change of 4.77% when contrasted with their (Bora, 2012) precision of 83.33%. In future, this work can be reached out through consolidation of better spell amendment components (might be at phonetic level) and word sense disambiguation. Additionally, we can recognize the objects and elements in the tweet and the introduction of the client towards them.