# Performing Text Modeling and Sentiment Analysis using LDA and deep learning

[1] Ishrat Shaikh, [2] Asma Basit Shahin

[1] Mtech Student [2] Assistant Professor

[1][2]Computer Science and Engineering, JNEC Aurangabad.

## Abstract

Emotional messages can be extracted from reviews of products, journals, forums, networking sites, excerpts from novels, and more. Although many relational data bases have bases impacts on validation checks. One type of product or service, the source of the message, can occur in stock market analysis or political discussions or news articles. Every place where people talk and share opinions freely may be the source. We intend to propose a Multigram (MMM) mixing model that can learn vocabulary about emotional feelings from document archives using NLP techniques. Second, we instigate the basis quality of the emotional model for english language (topic) created by the method that offers an Enhanced Latent Dirichlet Allocation (ELDA) using standard measurements such as consistency and frequency rating. Initially, it separates text from various articles in various sources. All the sentences are not emotional cues that we do, so it is possible that a high proportion of the sentence "neutral" is more than a sentence with an emotional component. In order to understand the impact of this, sentences will be sent to Sense Analyzer to create a basic label of "positive", "neutral" or "negative". The proposed method of detecting two emotions: classification of Words - Emotions and Documentation Ranking of Emotions.

## I Introduction

Social media provides access to the emotional information of weakly labeled users, including symbolism display called as emoticons along with emotional hashtags, that can be utilized to learn various emotions. In particular, emotion detection provides useful knowledge in designing diff docs displays utilizing binary numbers to plain frequency numbers to complex emotional concepts. It can also be used to search and index content through many emotions. To show diff emotions of social communities and groups emotions are divided into 6 types viz widely used to describe basic human emotions based on expressions [1]: disgusting, fearful, happiness, sadness, angry and surprise. Most of these are related to negative feelings, which "Surprise" is the most vague because it can be associated with positive or negative feelings.

Interestingly, the number of basic human emotions has been "reduced" or divided into 4 categories; sadness, Happiness, anger / disgust and fear / surprise [2] is a matter of surprising for many of us with only 4 basic emotions.

In the task recognition process, the most common approaches are based on rules, on statistics and on hybrids, and their use depends on factors such as data availability, domain experience and domain specificity. In the case of sentiment analysis, this task can be approached using lexical-based methods, machine learning or a concept-level approach [3]. We intend to explore on how we can work successfully using automatic learning methods using deep learning techniques.

Emotional messages can be found in product reviews and more. Although many databases focus on checking products or types of services, text sources may come from news articles, stock market analysis or political discussions in any places where people talk and share their opinions. Initially, messages from various articles are drawn from various sources. There is no assurance that the review will provide a balanced mix of words with all the required emotions.. In other words, since not every sentence has an interesting emotional signal, we have a high probability that the proportion of "neutral" sentences is beyond the emotional component. In order to understand the impact of this, sentences will be sent to the Sense Analyzer to create a basic label of "positive", "neutral" or "negative".

**I LATENT DIRICHLET ALLOCATION**

LDA is popularly used context modeling algorithm that finds hidden topics from document collection. Here, each topic that is discovered is represented by the word LDA. Find the hidden topic from the document, utilizing the words that occur in each document. Documents D = {d1, d2, …, dm} is the collection. And the total number of documents in the collection is 'm'. The LDA is applied to all documents to be divided into a specified number of topics. The main possible course of action behind the LDA is under the assumption that each document has several topics and each topic can be defined as cross-word distribution.

The LDA model is shown using 2 levels, the collection level and the document level. At the document level, each di document from the document set is entitled by the distribution of topics θ di = (ϑdi, 1, ϑdi, 2, …., ϑdi, V), V is the nos. of topics. At the assemblage level, documents are displayed as D. Each document is entitled by a likelihood distribution on words, คำ_j for topic j. Overall, we have รวม = {ϕ1, ϕ2,.., ϕv} for all topics. The LDA model also originates word assignments in addition to displaying these two levels, that is, the occurrence of words will be considered related to the topic. The dispensation of topics in the collection of all D documents can be calculated from the LDA model, ϕ D = (ϑD, 1, ϑD, 2,., .., ϑD, V). In the collection D. The most important participation of the LDA model is to display the topic using word dispensation and data representation using topic representation.

Showing topics indicating the words that are vital to the document & topic presentation. Identify which topics are important to the document. LDA can learn various topics from document collection and document digest by topic. There are various methods for new incoming documents to determine the content in the training topic. In this article, we use the topic format according to the format to display the document and offer the correct ranking method which determines the relevance of the new incoming document.

**II PATTERN ENHANCED LDA**

The pattern-based representation will overcome the constraints of word substitution, which provide the correct way to display documents. Moreover, as a representative based on the data structure model prepared by the alliance between words. In order to find a meaningful meaning from the document set to represent the topic and document, 2 steps are proposed:

(1) Process new transaction data set using the LDA results of document collection D

(2) Create a model that is represented by a set of transaction data to show the needs of users.

(3) Request the Equivalence Pattern class

**1) Create TD(Transactional Dataset)**

Let $Rdi, zj$ represent the word header definition for the $Zj$ topic in di $Rdi$ documents By user Each word under the topic that occurs in each document is called a document-specific transaction. A specific document transaction (TDT) is a unique set of words. For the assignment of all the words, $Rdi \mid 2j$; I'm $Mj\}$ where lij is a collection of words that originated in $Rdi$, $Zj$ lij, called transaction specific documents. For each topic in D, we can create a set of transaction V data ($\Gamma 1, \Gamma 2, ...., \Gamma v$).

**2) Generate Pattern based Representation**

In the form presented, the method according to the pattern frequently generated from each transaction set. $\Gamma j$ is used to represent the pattern $Zj$ as a collection of allied words for the minimum support limit that sets itemset items. X in $\Gamma j$ will occur only when supp (X)> = σ when supp (X) is the assist of X, which is the number of transactions in $\Gamma j$ with X.

Users specify the minimum support criteria. The frequency of the 'X' series is defined as a set of all common themes. The $Zj$ topic is shown as $Xzi = \{Xi1, X12,., Ximi\}$ where mi is the overall number of variations in $Xzi$ and v is overall number of topics

**3) PATTERN EQUIVALENCE CLASS**

The particular of patterns that are frequently received from the previous steps is very large and a lot of patterns do not need to be useful. There are many forms of strong presentation to present useful patterns instead of patterns that occur frequently from large data sets, such as the maximum format and closed patterns. For data sets, the particular concise forms is lesser than the particular of frequently created patterns.

Let EC1 and EC2 be two dissimilar uniform classes of the exactly indistinguishable transaction data set. EC1 ∩ EC2 = ϕ, this is in accordance to the equivalence class that is mutually exclusive. There are 2 related parts used in proposed format. Training used in creating user's patterns of interest from the assemblage of documents intended for training and the filtering section determines the pertinence of new incoming data. In the proposed model, check the meaning of the meaning model using the Stanford NLP library..

**III STANFORD NLP CLASSIFIER**

Stanford NER is implementation of the Named Entity Recognition. Named Entity Recognition (NER) a label for the order words in the text, which is the name of things such as the individual's name and company or gene name and protein.

Comes with a separate feature designed for Named Entity Recognition and many options for defining features. The download includes recognizers of entities with a good name for English, in particular for the 3 classes (LOCATION, ORGANIZATION,PERSON) and we also make available on this page several different models for vaious languages and circumstances, including trained models only in 2003 English CoNLL training data. Stanford Named Entity Recognition is recognized as CRFClassifier. The software provides a general implementation of conditional random field (CRF) sequence models (arbitrary order). That is, by training your own models on tagged data, you can use this code to build sequence models for NER or any other task.

**IV Word2Vec**

Word2vec is a group of related topics that are used in production called Word Weddings. These models are shallow neural networks that are trained to create a new vocabulary context. After training, the word2vec model can be used to assign each word to the vector of hundreds of elements, in general, that represent the relation of that word to other words. This vector is the hidden layer of the neural network.

Word2vec is based on cross bag or continuous word bag (CBOW) to create a wedding word. It was created by a team of researchers led by Tomas Mikolov on Google. The algorithm was later analyzed and explained by other researchers.

### II Related Work

The most common methods for creating a domain-specific vocabulary are checked because it is based on content labeled with feelings or weakness in the domain. For example, researchers use Pointwise's data to learn vocabulary and emotions from tweets that have soft labels. Jacopo Staiano and Marco Guerini propose to take advantage of the news articles of the crowd. (www.rappler.com) .com) for vocabulary building, including frequency distribution of word documents and emotional distribution in documents.

Researchers have applied natural models such as the Dirichlet allocation (LDA) for vocabulary building. Yanghui Rao and colleagues have combined the emotional scores of users in the document. (http://news.sina.com.cn/society [In Chinese]), document frequency distribution and problem distribution of LDA documents to learn emotional vocabulary, vocabulary topics. M.Yang and his colleagues offer a semi-supervised LDA approach that uses a domain-independent emotional vocabulary set to guide the LDA process to learn topics related to emotions. However, the topics learned from this method are not always accurate because the coverage of vocabulary varies from domain to domain. However, supervised LDA (sLDA) offers a more accurate method of learning. Simulation of emotional topics for building vocabulary from emotional bodies.

A. Valitutti C presented SENTIWORDNET 3.0, which is a clearly designed vocabulary source to support the classification of feelings and the use of digging. SENTIWORDNET 3.0 reviews are updated versions of SENTIWORDNET 1.0, which is Data source The current public language vocabulary for research has been licensed for more than 300 research groups and used in various research projects. The SENTIWORDNET 1.0 and 3.0 worlds are the result of automatic WORDNET reviews based on the levels of positivity, negativity and SENTIWORDNET Neutrality 1.0 and 3.0 are different.

Sentinet is used for creating rich online resources, expanding opportunities for digging opinions and analyzing the great confidence Godbole and the faculty (2007) created a confidence analysis system. In the English vocabulary dictionary to assess the general reputation of the entity, Taboada and the faculty (2011) present a more complex model according to the model, including rejection and repetition using weight accepted Liu (2010) suggested methods Effective in modern times to analyze confidence and privacy in English.

Jijkoun and the Faculty (2010) focus on creating specific vocabulary topics. Li and faculty (2010) separate confidence and dependence on local topics as well as global ones. Gindl et al. (2010) conducted a confidence analysis across domain context with Pak and Paroubek (2010) focus on Twitter to conduct research on the spoken language of English.

Denecke (2008) conducted a multi-lingual sentiment analysis using SentiWordNet, Halalia and Faculty (2007) Multi-lingual personal language learning through cross-language Abbasi screening and faculty (2008). Separate the characteristics of Arabic that require specific language knowledge. Gˆınscˇa and faculty (2011) work to Improve confidence analysis system in Romania.

Banea and the faculty (2008) show that computer translation can be done quite well when extending the analysis of personal feelings to a multi-lingual environment which causes repetitive work motivation in the analysis.

Boiy and Moens (2009) analyzed the confidence in machine learning using multi-language text. A profound learning method will outline the built-in vocabulary that provides concise features that reflect the meaning of basic vocabulary.

Turian and the School (2010) Create powerful, integrated words when practicing true and damaged phrases. Zou and the faculty (2013) include automatic translation and word expressions to create bilingual resources. Socrates and the Faculty (2012) demonstrate an effective method for trusting in the English language by using embedded words that can be easily extended to other languages by training in the appropriate text of subject

Knautz et al., 2010 focuses on developing complex search algorithms that distinguish between the rental emotions associated with the product. For example, customers may search for banks, mutual funds or shares that people trust. Help organizations may search for activities and stories that create compassion and highlight them in a fundraising campaign. In addition, non-emotional intelligence systems may be victims of abuse. For example, recently found that online sellers intentionally hurt their customers because of negative online reviews that translate into higher rankings in Google search. In [5] Yang explains that emotions are seen to evolve through their adaptation values in dealing with basic tasks in life. Each mood has unique features: signs, physiology and previous events. Each emotion has the same characteristics as other emotions, such as fast attacks, short duration, unauthorized birth, automatic assessment and interaction between responses.

[3] Alessandro Valituttihas presents language resources to express words about emotions. This resource (called WORDNET) has been developed from WORDNET through selection and labeling a subset of coincidences that express emotions.

[11]Turney, Saif Mohammad and Peter David show that the power and wisdom of the crowd can be used to create quality and large words. {Mood and vocabulary {Terminal terminology, terminal linking, fast and affordable. We identify challenges in increasing emotional feedback in fundraising situations and offering solutions to those problems. In addition to questions about emotions related to vocabulary, they also talk about how to combine questions, word choices that can discourage harmful input, help identify instances in which annotations may not be familiar. With words that are annotated with that goal) level of feeling (Instead of word level) We carried out demonstrations on how to ask questions, emotional annotations and show that asking which words are related to emotions leads to higher agreements between critics. Clearer than what was asked by the words, which caused the emotional management of the relationship with the customer by taking appropriate action based on the emotional status of the customer (For instance, disgust, fear, anger, happy, surprised, sad).

[9]Vel_asquez, Ravaja et al., Create a conversation system that responds divinely to the various emotional conditions of users For instance, in a game that recognizes emotions.

[10] Bellegarda, 2010, works to identify the emotions that headlines in newspapers are trying to make. R-ranking and categorizing data / answers in online questions {Answer to forums (Adamic et al., 2008). For example, high emotional responses may be ranked lower. Perceiving how individuals utilize words like emotions and metaphors to convince and oppress others. (For instance, in advocacy) (K_ovecses, 2003).

[8]Erik Cambria, Soujanya Poria, Rajiv Bajpai explains the difference between traditional AI systems and human intelligence. They conclude that the ability to control general knowledge is compiled from a lifetime of learning and a wise decision-making experience. This allows humans to adapt to the new situation in which the AI has been disastrous due to lack of specific rules and general implications. General knowledge also provides basic information that helps humans to succeed in operations in social situations where such knowledge is often assumed. Because general information contains information that humans accept, collecting is a very difficult task.

Previous versions of SenticNet focused on gathering this type of knowledge for confidence analysis. But they were very limited by the inability to conclude.
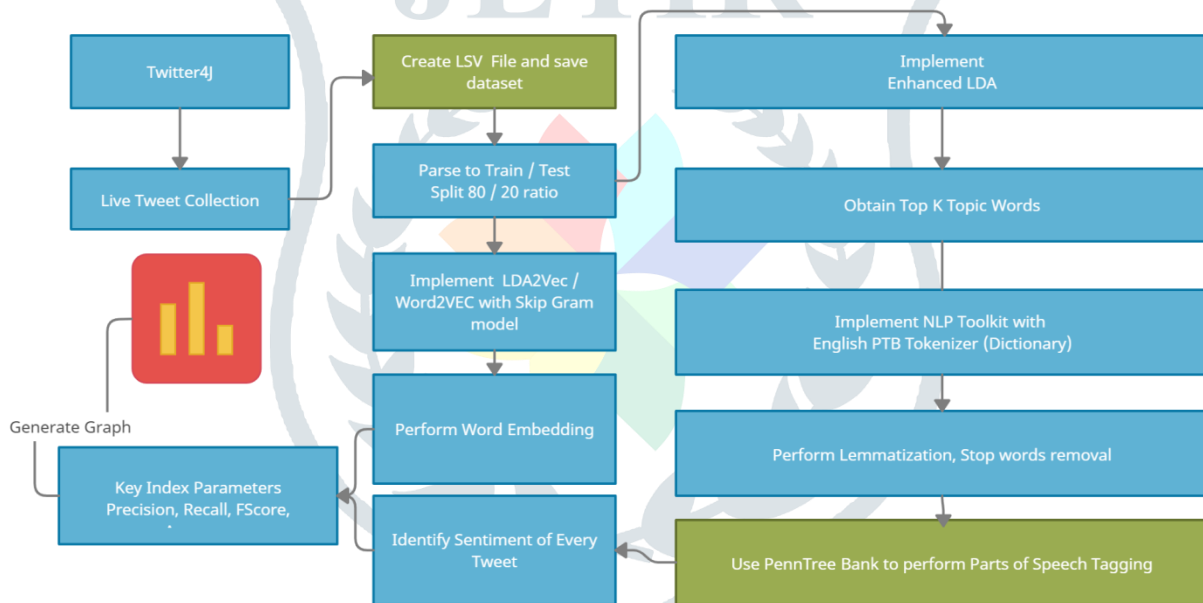
[5]Yun Ching and colleagues use SenticNet and Bay models for contextual thinking concepts. Mauro Dragoni and colleagues offer a vague framework that combines Word_Net, Concept_Net and Senti-Net to draw key concepts from sentence. iFeel is a system that allows users to create their own confidence analysis framework using Senti_Net, Senti WordNet and other confidence analysis methods. Jose Chenlo and David Losada use SenticNet to separate bags of concept and terminal properties for Personal detection and other confidence analysis work.

[6]Jay Kuan-Chieh Chung and colleagues use the SenticNet concept as a seed and propose a random method of walking in ConceptNet to draw additional ideas along with other work terminal scores. Offer sharing of cognition of information and machine learning. For analyzing sentiment on Twitter, categorizing short messages 40 and digging frame comments They conclude that the blended approach is intended to understand the rules of thought that control the feelings and hints that can bear these concepts from perception to the use of words in the mind.

**Proposed System**

For developing domain specific terms, most occur under supervision as they are dependent on emotionally labeled or weak content in the domain. For example, researchers use Pointwise's data to learn vocabulary-emotions from tweets that are softly labeled with emotional hashtags and by taking advantage of the crowd-like emotional news articles. (www.rappler.com) for creating dictionaries by combining documents and emotional distribution through documents.

Our system will consists of phrases or sentences as well as labels of emotions. It will work with emotion dataset and training dataset and to obtain valency in the form of emotional and neutral concept that covers many sentences. Training the system will handle eliminating stop words and neutral words to get the emotions of the electrodes, then distinguish in their class.



**Figure 1.0  Proposed Architecture**

The chief characteristics of the proposed model are as follows:

(1) Every subject is entitled by patterns

(2) Information is filtered using Stanford NLP framework.

(3) Provide a more precise document modeling method for classification.

In (form) pattern based topic model, which has been employed in Information Filtering, can be acknowledged as a "Post-LDA" model based on the patterns that are produced from the topic representations of the LDA model.

Patterns can represent more specific meanings than single words. By comparing the word-based topic model with pattern-based topic models, the pattern based model can be used to represent the semantic content of the user's

documents more accurately than word based document. However, many times the number of patterns in few of the topics can be huge and many of the patterns are not enough to represent specific topics.

We propose to overcome the limitation of existing system by using Natural Language Processing Natural language processing (NLP), i.e., the Stanford NLP library used in enhanced LDA algorithm for filtering semantic meanings of patterns from the collections of topics. The specificity (precision) of the emotion class can be affected and we can have two closely related emotion classes, say, *ecstatic* and *excited* as two separate classes, or *afraid* and *scared* as two separate classes, instead of one class with label *excited* and *afraid*, respectively.

Feature representation is the next stage in the process that includes representations and the use of skip-gram and n-gram, characters instead of words in a sentence, inclusion of a part-of-speech tag, or phrase structure tree.

Next process is to obtain aspect of data, using heuristic rules that we can define from our NLP framework and Penn Tree bank and obtain different aspects such as Nouns, Pronouns, Adjectives etc.

We need to calculate that the number of neurons and layers in a neural network has on an emotion classification task.

**Steps**

- Learn initial model from training data.
- Set mixture parameter $\lambda$ using Word2Vec representation.
- Set estimation of hidden variable $Zw$.
- Perform Maximization step (M-step) and obtain parameter Theta($e$ )
- Generate model for each document. (LDA, Theta Value for Dt)
- Set Burnout Parameter.
- Perform Gibbs Sampling
- Calculate Emotional Valence
- Calculate Neutral Valence
- Obtain vector value for words and generate lexicon.

**EnhancedLDA Psuedocode**

```
Input: user interest model UE = { E(Z1), . . ., E(ZV)}, a list of incoming
document Din
Output: rankE(d), d ∈ Din
   1: rank(d) = 0

   2: for each d ∈ Din do

   3: for each topic Zj ∈ [Z1, Zv] do

   4: for each equivalence class ECjk ∈ E(Zj) do
5: scan ECk,j and find maximum matched pattern which exists in d
   6: update rankE (d) using equation(1)
```

7: rank(d) := rank(d) + ||$^{0.5}$ * fjk * υD,j * uniform distribution * equivalent class frequency

```
   8: end for
9:  end for
10: end for
```

**Input 1:**

• (required): word list (key = "getVecFromWord")

**Output 1:**

• 300-dimensional vector representation of a given word

- Input 2:
    (Required): List of 300-dimensional vectors (key = "getWordFromVec")
- Output 2:
    The top 10 words that are most consistent with the vector defined in the vector space
- Input 3:
    (Required): Two words list (key = "similarity between the words")
- Output 3:
    Similarity scores between the two words received
- Input 4:
    (required): word list (key = "doesntMatch")
- Output 4:
    Return words that do not match the other words in the list.
- Input 5:
    (required): vector arithmetic using the algorithm proposed in the original word2vec paper (key = "vectorArithmetic")
    (only need one): List of words that will be positive in vector math (key = "positive")
    (only need one): List of words that will be negative in vector calculations (key = "negative")
    (Optional): The number of results I want to return. The default is 10 (key = "numResults")
- Results 5:
    N maximum (if specified, otherwise N = 10), words that are close to the product vector of mathematical operations

**Input 6:**

    (Required): Vector arithmetic that uses a different algorithm. (key = "vectorArithmeticCosmul")
    (Only one required): A list of words that will be positive in vector arithmetic. (key = "positive")
    (Only one required): A list of words that will be negative in vector arithmetic. (key = "negative")
    (Optional): Number of results I want to return. Default is 10. (key = "numResults")

**Output 6:**

- Top 10 words that are closest to the product vector of the arithmetic operation.

We evaluate a lexicon's ability to classify a collection of target words hand-labeled with emotions. More formally, given an arbitrary word w, the task is to predict an emotion label e (E) for w using the word-emotion lexicon. Because it quantifies the associations between words in a vocabulary V and a range of emotions in E, for any given arbitrary word w, the dominant emotion e being expressed is calculated using the lexicon.

**About Dataset**

The SemEval data set contains news headlines drawn from major newspapers such as the New York Times, CNN and BBC News, as well as from the Google News search engine. We decided to focus on the news topic for two reasons. For the first time, news often contains a lot of emotional content because they describe national or global milestones and write in a format that refers to attracting readers' attention. Secondly, the structure of news headlines is appropriate for the goal of making sentence annotations at the emotional level.

**Conclusion**

The main objective of using deep learning is that they aim to extract those features which disentangle the hidden factors of variations. This will help to perform the transfer across different domains. In this case, they were expecting the concept which characterized the review. They considered some of the factors like positive reviews to check the disentanglement of the dataset. We have considered unlabeled data from different and labels from a single domain and followed two-step procedure for sentiment analysis and the word2vec algorithm trains the linear classifier on transformed labeled data thereby assisting in sentiment analysis and detection.

## REFERENCES

[1] A Bandhakavi, N Wiratunga, and S Massie, Robert Gordon University Queen's University, Belfast

[2] A. Esuli, F. Sebastiani  and S. Baccianella, "Sentiwordnet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining," Proc. Int'l Conf. Language Resources and Evaluation, 2010, pp. 18–38.

[3] A. Valitutti C. Strapparava and , "Wordnet-Af fect: An Af fective Extension of Wordnet ," Proc . Int 'l Conf. Language Resources and Evaluation, 2004, pp. 1083–1086. 3. S.M. Mohammad, "#emotional tweets," Proc. 1st Jt. Conf. Lexical and Computational Semantics, 2012, pp. 246–255.

[4] Rao Y. et al., "Building Emotional Dictionary for Sentiment Analysis of Online News," World Wide Web, vol. 17, no. 4, 2014, pp. 723–742. 5.

[5] C. Yang, K.H.Y. Lin, and H.H. Chen, "Emotion Classification Using Web Blog Corpora," Proc. IEEE/WIC/ACM Int'l Conf. Web Intelligence, 2007, pp. 275–278.

[6]. C. Yang, K.H.Y. Lin, and H.H. Chen, "Emotion Classification Using Web Blog Corpora," Proc. IEEE/WIC/ACM Int'l Conf. Web Intelligence, 2007, pp. 275–278.

[7] M. Yang et al., "A Topic Model for Building Fine-Grained Domain- Specific Emotion Lexicon," Proc. 52nd Ann. Conf. Assoc. Computational Linguistics, 2014, pp. 421–426.

[8] D.M. Blei and J.D. McAuliffe, "Supervised Topic Models," Advances in Neural Information Processing Systems, 2007, pp. 121–128.

[9] Erik. Cambria, "Affective Computing and Sentiment Analysis," IEEE Intelligent Systems, vol. 31, no. 2, 2016, pp. 102–107.

[10] Callison-Burch, C. (2009). Fast, cheap and creative: Evaluating translation quality using amazon's mechanical turk. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP-2009).

[11] J. R. Bellegarda. 2010. Emotion analysis using latent affective folding and embedding. In Proceedings of the NAACL HLT 2010 workshop on computational approaches to analysis and generation of emotion in text, pages 1–9. Association for Computational Linguistics.

[12] S. M. Mohammad and P. D. Turney. 2013. Crowdsourcing a word–emotion association lexicon. Computational Intelligence, 29(3):436–465.