



DER TECHNISCHEN UNIVERSITÄT MÜNCHEN

Master's Thesis in Informatics

Utilizing Crowd Intelligence for Online Detection of Emotional Distress

Siddhant Goel





DER TECHNISCHEN UNIVERSITÄT MÜNCHEN

Master's Thesis in Informatics

Utilizing Crowd Intelligence for Online Detection of Emotional Distress

Insert thesis title in German here

Author: Siddhant Goel
Supervisor: Prof. Dr. Claudia Eckert
Advisor: Han Xiao, M.Sc.
Date: March 15, 2013



Ich versichere, dass ich diese Diplomarbeit selbständig verfasst und nur die angegebenen Quellen und Hilfsmittel verwendet habe.

I assure the single handed composition of this master's thesis only supported by declared resources.

München, den March 15, 2013

Siddhant Goel

Acknowledgments

If someone contributed to the thesis... might be good to thank them here.

Abstract

An abstracts abstracts the thesis!

Contents

Acknowledgements	vii
Abstract	ix
Outline of the Thesis	xiii
I. Introduction	1
1. Introduction	3
1.1. Latex Introduction	3
2. Related Work	5
3. Problem Definition	7
II. Methodology	9
4. Classification Methods	11
4.1. Introduction	11
4.2. Support Vector Machines	12
4.3. Multiple Kernel Learning	12
5. Text Representation	15
5.1. Introduction	15
5.2. Preprocessing	15
5.3. Representation and Vector Space Classification	16
6. Ensemble Learning	19
6.1. Bagging	19
6.2. Boosting	19
6.3. Stacking	20
III. Experimental Results	21
7. Experiments	23
7.1. Dataset	23
7.2. Approach and Setup	24

Contents

7.3. System Details	24
7.3.1. Ratings	24
8. Results	25
IV. Conclusion	27
9. Conclusion	29
Appendix	33
A. Appendix	33
Bibliography	35

Outline of the Thesis

Part I: Introduction

CHAPTER 1: INTRODUCTION This chapter presents an overview of the thesis and its purpose. Furthermore, it will discuss the sense of life in a very general approach.

CHAPTER 2: RELATED WORK Related Work

CHAPTER 3: PROBLEM DEFINITION Problem Definition

Part II: Theoretical Background

CHAPTER 1: CLASSIFICATION METHODS

CHAPTER 2: TEXT REPRESENTATION

CHAPTER 3: ENSEMBLE LEARNING

Part III: Experiments

CHAPTER 1: EXPERIMENTS

CHAPTER 2: APPLICATION

CHAPTER 3: RESULTS

Part IV: Conclusion

CHAPTER 1: CONCLUSION

Part I.

Introduction

1. Introduction

We start the thesis with an introduction. Do not spend time on formating your thesis, but on its content.

1.1. Latex Introduction

There is no need for a latex introduction since there is plenty of literature out there.

2. Related Work

Related Work goes here

3. Problem Definition

Define the problem here

Part II.

Methodology

4. Classification Methods

Machine learning is a branch of computer science that deals with building and analyzing systems that are able to learn from data. Algorithms based upon such analysis involve constructing a model from a given dataset, and then using this model to perform required tasks. Machine learning techniques can broadly be divided into two categories - supervised learning, and unsupervised learning.

- Supervised Learning

Methods falling in this category operate in two phases

- In the first step, the availability of training data is assumed, which is used to build a model so that it takes into account the structure of the given dataset
- In the second step, this model is used to make predictions on the testing data (the real world data). This is the data that the model has not seen yet, and is required to make predictions on.

- Unsupervised Learning

Methods falling under this category operate in a single phase. It starts with a model with zero knowledge about the structure of the given dataset. As data is fed into the model, it continuously learns the structure of the given dataset and calculates the predictions based on this knowledge. The main difference between this family of algorithms and supervised learning algorithms is the presence/absence of training labels.

The primary focus in this thesis remains on supervised learning methods.

4.1. Introduction

Classification is one of the fundamental problems in machine learning. Given a dataset \mathbf{X} , it is required to separate the samples contained within the dataset into two (or more than two, depending on the input) classes. Formally, given a dataset that contains N instances $(\mathbf{X}_n, \mathbf{Y}_n)_{n=1}^N$, where each instance $(\mathbf{x}_n, \mathbf{y}_n)$ is of the form $[(\mathbf{x}_{n,1}, \mathbf{x}_{n,2}, \dots, \mathbf{x}_{n,D}), \mathbf{y}_n]$, (each $\mathbf{x}_{n,d}$ being the value of the feature $d \in [1, D]$, and \mathbf{y}_n is the label of the sample which can take a limited number of possible values) the aim is to calculate the value of \mathbf{y}_n given the feature information. This separation can usually be done using a supervised learning method, in which case the training data (on which the model is built) is given and it is required to predict the labels of the test data, or using unsupervised methods, where the model is required to identify the categories of the samples without any information on \mathbf{y} .

The performance of a particular classifier varies with the type of the data to be classified. Not all classifiers are good for all classes of problems. Some classifiers suit a particular problem more than some others; choosing a classifier for a problem still remains a decision

which may or may not be completely scientific, even though there have been a number of tests been done to correlate classifier performance with data type [citation needed].

4.2. Support Vector Machines

Support Vector Machines (SVMs) form a fairly popular class of machine learning algorithms used mainly for binary classification and regression analysis. Given the training data, the goal of an SVM is to find a decision boundary (a hyperplane in a high or infinite dimensional space) that separates the two classes of data while maximizing the distance of the boundary from any data point. The resulting decision function is fully specified by a (usually small) subset of the data, and the points in this subset are referred to as support vectors.

All classifiers resort to a distance function, in some form or the other, that can provide a similarity measurement between two points. In the simplest form of an SVM, the distance function is simply the dot product between two points, and such SVMs are referred to as *Linear Support Vector Machines*. In the case that a simple linear SVM is not able to find a sufficiently accurate decision boundary that can separate the data points into two classes (simply because the input data is not linearly separable), the so-called *kernel trick* is used, which involves transforming the data from a low dimensional space to a much higher dimensional feature space (in which the input data may be separable) using an appropriate function $\phi(\mathbf{x}) : \mathbf{X} \in \mathbb{R}^L \rightarrow \mathbb{R}^H (L \ll H)$, and then using the kernel function $\mathbf{K}(\phi(\mathbf{x}), \phi(\mathbf{y}))$ to actually perform the separation. The trick involved is that even though the transformation $\phi(x)$ may be expensive, computing the final similarity value $\mathbf{K}(\phi(\mathbf{x}), \phi(\mathbf{y}))$ is not.

The most popular kernel functions include

- Linear (the simple SVM) - $\mathbf{K}(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})$
- Polynomial - $\mathbf{K}(\mathbf{x}, \mathbf{y}) = (\gamma \mathbf{x} \cdot \mathbf{y} + \mathbf{c})^d$
- Radial Basis - $\mathbf{K}(\mathbf{x}, \mathbf{y}) = \exp(-\gamma |\mathbf{x} - \mathbf{y}|^2)$
- Sigmoid - $\mathbf{K}(\mathbf{x}, \mathbf{y}) = \tanh(\gamma \mathbf{x} \cdot \mathbf{y} + \mathbf{c})$

4.3. Multiple Kernel Learning

Multiple Kernel Learning makes an improvement to standalone support vector machines by using multiple kernels instead of a single one. In contrast to standard SVMs which use a single kernel function to calculate the similarity score between two data points, MKL algorithms combine scores from multiple kernels to obtain one single score, which is then used as the final similarity score. As summarized by [1], multiple kernel learning algorithms can be put into 12 major categories, based on some of their key properties (learning method, functional form, target function, training method, base learner, and the computational complexity).

The two main uses of MKL algorithms are also discussed by [1]. The first one uses the fact that since different kernels may be used as different notions of similarity, the accuracy obtained from each kernel may be maximized, and then the learning algorithm could be

left to decide whether one kernel (which works better than the rest) or a combination of kernels is suitable for the task at hand. The second use is a more traditional use, where different kernels corresponding to different notions of similarity are combined, either in a linear or in a non-linear fashion to obtain a single kernel. As with all machine learning algorithms, multiple kernel learning does not suit all possible use cases, but in a lot of cases, combining multiple kernels may result in increased accuracy.

5. Text Representation

5.1. Introduction

Machine learning algorithms are designed to operate on real values. While classifying documents, a prerequisite is to convert the text in natural language to a format that such an algorithm can understand and work on. One such format is the vector space mode. Such processes include obtaining text from documents, building a vocabulary, text preprocessing (such as tokenization, stemming, stop words removal), and converting terms to actual usable feature values. All such issues are explored in this chapter, except the issues that arise due to usage of non-English languages (since this work focuses only on text in English).

5.2. Preprocessing

Preprocessing documents typically involves all the steps that are needed to be performed before proceeding on to extracting meaningful information from text. Text documents in their original forms normally contain a lot of redundant information that usually does not affect how a classifier behaves. Preprocessing a document removes all the noise, and brings a document in a state where it can be used for scoring. The following are the various filters that are applied.

- Tokenization

Tokenization is the process of splitting a sequence of words into distinct pieces of alphanumeric characters, called tokens. Extra characters which are not needed, such as punctuations and white spaces are removed. For instance, tokenizing the text "*The quick, brown, fox jumps over the lazy dog;*" results in the following list of tokens -

the	quick	brown	fox	jumps	over	the	lazy	dog
-----	-------	-------	-----	-------	------	-----	------	-----

In most cases, tokens are split using whitespaces, line breaks, or punctuation characters. Splitting on white spaces may also result in some loss of information. For instance, if the string *San Francisco* appears in the text, then the tokens extracted will be *San* and *Francisco*, whereas the correct tokenization should treat both the terms in one token. Such problems are solved using n-grams, which are explained later in this section.

- Stop Words Removal

Some of the words in the English language such as *and*, *the*, and *a*, besides some others appear in almost every text. These words add very little meaning to the text

on the whole, and their frequency of occurrence is very high. Consequently, since these documents appear a lot in almost every document, it may lead to high, but inaccurate similarity scores between two documents; the best option is to remove such words from the text.

- Stemming

Stemming is the process of reducing words to their root form, usually by stripping some characters from the word endings. A strict requirement for stemming is that related words must stem to the same final word. For instance, the words *car*, *cars*, *car's* and *car's*, must all stem to the same word *car*. This helps in removing unnecessary information from the text which would otherwise inflate the vocabulary with low-signal information.

5.3. Representation and Vector Space Classification

After a document has been preprocessed, it needs to be represented in terms of the useful content that it has, such that the relative importance of each word has been taken into account. Since a document is just a collection of tokens, a document \mathbf{x}_n is represented as a vector $\mathbf{x}_n = (x_{n,1}, x_{n,2}, \dots, x_{n,D})$, where each dimension $x_{n,d}$ corresponds to a token, and the value depends on its occurrence, either only in the document, or in the document as well as in the entire corpus. This approach is called the bag-of-words model.

For instance, consider two simple text documents,

We are headquartered in Munich, Germany

We also have an office in Berlin, Germany

Based on these two documents, the token dictionary is built like the following -

"We": 1, "are": 2, "headquartered": 3, "in": 4, "Munich": 5, "Germany": 6, "also": 7, "have": 8, "an": 9, "office": 10, "Berlin": 11

and the documents are represented as the following vectors -

1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0
1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1

where each value represents the number of times the corresponding token appeared in the document. The order of words does not matter; the features simply indicate whether or not the word appears in the document.

In this approach, the values represent term frequencies. A slight improvement is to divide the term weights with the total number of terms in the document, which assigns the feature values while also taking into account the relative importance of the term in the document. But an even better improvement is to use the *tf-idf* value, which (for a term in the document) is defined as

$$\text{tfidf}_{t,d} = \text{tf}_{t,d} * \text{idf}_t$$

where $\text{tf}_{t,d}$ is the term frequency of the term t in document d , and idf_t is the inverse document frequency of the term t across the entire collection of documents, usually defined

as $\log(N/df_t)$, where N is the total number of documents, and df_t is the number of times this term appears in the entire collection of documents. The main idea here is to reduce the term frequency weight of a term by a factor that increases proportional to its frequency in the corpus. This tf-idf value is the highest when the term occurs many times within a small number of documents, and the lowest when the term occurs in almost all documents. Hence it provides a good representation of a document in the vector space format. This representation is now used to calculate the similarity between two documents. The standard way for this is to compute the dot product between the two vector representations

$$\text{sim}(\mathbf{x}, \mathbf{y}) = \frac{\vec{\mathbf{x}} \cdot \vec{\mathbf{y}}}{|\vec{\mathbf{x}}||\vec{\mathbf{y}}|}$$

Representing documents as vectors and calculating similarities based on their dot products leads to a view of a corpus as a *term-document* matrix, the rows for which represent the documents, and the columns for which represent the terms present in the corpus.

6. Ensemble Learning

Ensemble Learning is a class of supervised learning algorithms that use multiple models instead of a single one to obtain better predictions. A common problem with classifiers is that not all classifiers are suited for all kinds of problems. Ensembles combine multiple classifiers to construct a (hopefully) better classifier. The performance of an ensemble is not guaranteed to be better than the individual classifiers, but depending on the problem, an ensemble usually outperforms its constituent classifiers by a fair margin. Using ensemble learning requires performing more computation, but ensemble learning provides better results when the classifiers are not similar to each other. In this chapter, we explore the main ensemble learning techniques that we employed in our work.

6.1. Bagging

Bagging, also referred to as Bootstrap Aggregation is one of the most basic forms of ensemble methods. Given a training dataset D of n samples, each sample having f different features, the aim of bagging is to combine a certain number of classifiers, each trained using a (mostly different) subset of the main data, such that every classifier learns about a different structure of the input data, and when the outputs from all such classifiers is combined, the final output is expected to be better than the individual predictions. Bagging is known to improve the performance when the classifier is slightly better than a random classifier. On the other hand, in cases where the constituent models are already strong, the performance is known to often degrade.

To train different classifiers on a different subset of the data, the training data is either split on the number of samples, or on the number of features. In the first case, we pick $n' (< n)$ samples (with replacement) and each time train a new classifier using the new dataset. Since the data points are picked with replacement, some of the data points may be common to some classifiers. In the second case, we train every new classifier using $m' (< m)$ features (with replacement). Whenever the ensemble is required to make a prediction on a new sample, a majority vote from all the constituent classifiers is taken, which becomes the final output of the ensemble. Bagging is known to reduce variance and avoid the overfitting problem.

6.2. Boosting

Boosting aims to build an ensemble by iteratively training weak classifiers on a distribution, and then finally putting them all together to build a strong classifier. In the training phase, the classifiers are trained one by one. Each training instance is assigned a numeric weight, which is the same for all samples in the beginning. After the first classifier has been trained, the weights for the points which were misclassified are increased, and the

resulting input data (along with the modified weight values) is fed to the second classifier. This process continues until the last classifier has been trained. Each classifier has an incentive to focus more on classifying correctly those points which the previous classifier classified incorrectly.

At all times, each classifier's performance is maintained in a vector (usually named α). Whenever a new point has to be classified, the individual predictions are obtained from the constituent classifiers, and a weighted sum (from the α values) is taken across all the classifiers to obtain the final prediction. One of the very popular Boosting techniques is AdaBoost. The main aim of Boosting is to build a strong classifier from a set of weak classifiers. To that effect, even classifiers only slightly better than random are considered useful, because in the final predictions, they will still contribute positively to the aggregate prediction by behaving like their inverses because of having negative coefficients in the final linear combination of classifiers.

6.3. Stacking

Stacking combines classifiers not by analyzing their performances on the training data, but by treating the individual predictions (on the training data) as a second level of training data, and then using this piece of data to build a final model that is ultimately used for making the final prediction. Different from the previous ensemble learning methods, stacking operates by feeding the initial training data to the set of base classifiers chosen. After these base classifiers have been well calibrated, the predictions on the training data are obtained from the base classifiers, forming a new dataset. This dataset is now fed into a second level classifier, which is responsible for making the final prediction. Often, the base classifiers (at the first level) are trained using different features of the input data to increase classifier diversity.

For making a prediction on a new point, predictions are made using all the base classifiers, forming an equivalent data point in which the number of features is the same as the number of base classifiers. This point is then run through the second level classifier to obtain the final prediction. In contrast to other ensemble learning methods, no voting or aggregation takes place in prediction. This class of methods has proven to be very successful, one recent example being that of the winning team in the Netflix prize competition.

Part III.

Experimental Results

7. Experiments

7.1. Dataset

This work focuses on detecting emotional distress from publicly available tweets/blogs. To predict the label for a particular tweet, a similar dataset is needed that contains text (where the size of one piece of text is comparable to the size of a single tweet), each having a label corresponding to whether or not the text indicates depression. This dataset comprises the training dataset.

Non availability of such data in the beginning of this work led to explore a slightly different problem in the same domain that operates on a similar dataset. The experiments were initially performed on a dataset made available by a machine learning competition website ([2]), where the aim is to predict whether a certain piece of text (a comment from a conversation on the internet) can be insulting to a user or not. The dataset provided contained the list of comments, each with a binary label. The dataset was split into two parts, the first file containing 3947 comments, while the second one containing 2235 comments.

After conducting experiments on this dataset, the dataset which can finally help towards building a system that can identify emotional distress from the given text was identified and consolidated. Stories from the website Reddit ([3]) were downloaded, which is an online community for people to interact with one another, hosting two main subreddits of interest - the subreddit where people post if in case they are planning to end their lives ([5]), and the subreddit where people post if in case they want to share their happy moments with others ([4]). The process of building the dataset was integrated in the main web interface, and every time there is a request to increase the size of the dataset, 500 stories from each of the subreddits mentioned are downloaded and stored in the database. Since this system also aims to incorporate crowd intelligence into the system (letting people assign labels to the training data), 2000 stories are initially labelled manually to build a strong foundation for the system, after which stories are left to be labelled on demand.

To actually incorporate identifying emotional distress into our system, the main data is fetched from Twitter ([?]). Twitter's public streaming API ([6]) is used to fetch 100 tweets every 3 hours, hence fetching 800 tweets every single day (interruptions in this process were faced from February 5 until February 15). This gives an overall view of the general sentiment of the public, on which the analysis is performed.

To summarize, the first part of the training data (comments on the web) comes from a competition on Kaggle, the second part comes from Reddit (main title of the stories posted by users), and the actual data for prediction comes from Twitter.

7.2. Approach and Setup

Various machine learning techniques are experimented with and evaluated, including standalone support vector machines, multiple kernel learning algorithms, and finally ensemble learning methods. This work can broadly be seen as being done in two phases. In the first phase, the learning techniques are evaluated using cross validation. In the second phase, a system is built that monitors emotional distress on the internet.

In the first phase, the comments [2] dataset is used, constructing n-grams (of length 2) and then obtaining the vector space representation for each comment along with the label for each. This data is now split into training and test data as per holdout cross validation. The model is then trained on the training data, and the predictions are obtained for the testing data. Since the actual predictions of the testing data are already known, the accuracy of the classifier thus trained is obtained. A 70-30 split on the data is performed, i.e. 70% of the data is used to train the models, and the remaining 30% is used to calculate predictions with.

This procedure is repeated for standalone support vector machines (using linear, polynomial, gaussian, and RBF kernels), multiple kernel learning methods, and the ensemble learning methods which include bagging, boosting, and stacking. The calculated accuracy for each method is then reported.

In the second phase, the system that is able to monitor public content (with focus on fetching content from Twitter) and identify the tweets which may signal emotional distress is implemented.

7.3. System Details

The final output of this work is a web based system that allows users to -

- assign labels to stories fetched from Reddit, which helps in building up a target set of training data, as well as tapping into crowd intelligence
- monitor a general *level of distress* amongst people who are posting on Twitter, grouped by date
- keep a check on certain tweets that have been classified by the model as depressed

The system mainly comprises of two modules - *ratings*, and *monitor*.

7.3.1. Ratings

The *ratings* module is responsible for allowing users to help build the training data. As mentioned before, the main source of text is Reddit. Stories fetched are simply stored in the database. When a user chooses to visit the ratings module, he/she is presented with the next story that does not have a label. The user can then proceed to assign a positive (depressed) or a negative (not depressed) label to it, which is then stored in the database.

8. Results

Results

Part IV.

Conclusion

9. Conclusion

Conclude

Appendix

A. Appendix

Appendix

Bibliography

- [1] Mehmet Gönen and Ethem Alpaydin. Localized algorithms for multiple kernel learning. *Pattern Recognition*, 46(3):795–807, 2013.
- [2] Kaggle. Detecting Insults in Social Commentary. <http://www.kaggle.com/c/detecting-insults-in-social-commentary>.
- [3] Reddit. reddit: the front page of the internet. <http://www.reddit.com>.
- [4] Reddit. /r/happy. <http://www.reddit.com/r/happy>.
- [5] Reddit. /r/suicidewatch. <http://www.reddit.com/r/suicidewatch>.
- [6] Twitter. Twitter Streaming API. <https://dev.twitter.com/docs/streaming-apis/streams/public>.