

# CS 565: Project Proposal

Rakesh Gupta - 170101050  
Rohith Chodavarapu - 170101017  
Sri Harsha Nakka - 170101040  
Sahithya Vemuri - 170101077

## Problem Statement

To refine content in an individual's social media accounts by predicting their impact on job prospects.

## Problem Description

Our key idea is to refine content by marking posts and comments containing hate or abusive speech and stress on skills needed or associated with a job to increase his chance of getting hired.

For example, when an individual is applying for a job at a company as a Lead designer, our aim is to refine public content on his/her social media to highlight artworks made, seminars and sessions conducted/ lead, UI/ UX developed etc and also to remove repulsive content such as toxic/ negative/ bullying content.

## Proposed Direction

We achieve this by performing sentimental analysis and/or by tagging on every post on social media of the individual that is publicly available to determine its impact on his/her hiring. Tagging is of great help as if the individual applies at some other company for other role then we can use the existing model instead of training the model again.

- **Stage 1** - Our initial step and main goal will be to deal with tagging toxic/ negative/ bullying contents posted by an individual.
- **Stage 2** - Then we will be predicting the impact of such comments on job prospects and advise them to remove or edit those posts. Enterprises can add tags for a role to filter the candidates tailored for the job.

Above step will be performed after checking the feasibility based on the results obtained from Stage 1 and we are currently focused on developing stage 1 in the view of the course.

# Major Challenges

- As we are dealing with fairly new ideas involving tagging, There are no appropriate dataset that best Fits our Needs. So we need to proceed with the existing generic taggers and filter the results according to our needs.
- As there can be many languages on Social Media we need to filter them out.(At Least initially we will be working on English)
- Filtering out unwanted content such as images,videos,audio files and other attachments from text.
- Giving intended meanings to slangs. e.x.  
E.x when someone says "Party last night was lit!" here lit doesn't mean something was set on fire but means very good/ enjoyable.
- We also tend to use slang from different languages like Amigo,bon appetit,Hola!, etc... Ignoring the whole post if it just contains them is not good.So we need to come up with a way to avoid it.
- The metrics used by different employers varies for a role. For example, a startup needs someone who can work in a broader area of work but it may not be necessary for a global conglomerate which has several employees dedicated for a specific area of work.
- The evaluation of an individual is very personalized and is hard to generalize.
- Most of the Sentimental Analysis tools use a very narrow range for classification(+ve, neutral, -ve). But we need a wide scope as our range of predictions are many.

## Brief Review of existing models

- The model in [1] is one of the founding papers in text classification using ml. In this paper, the focus is only to classify the documents into several particular categories. The proposed approach consists of preprocessing the text through stemming, removing stop words and building a vector representation of text.  
Later on, feature selection and feature transformation is done and a machine learning algorithm composed of multiple text categorization classifiers is used to determine the category.
- The model in [2] implements an extension of the problem in 2 parts which are to filter the texts(messages/posts etc) which are offensive/vulgar and to black list the user based on a user estimate.  
For the filter part big data is collected from an online social network and a simple bag of words and naive Bayesian methods are used to classify the texts into positive, negative or neutral and filters all the negative texts.

- In the model by Murat Demirbas[3], the focus was mainly on microblogging sites(such as twitter). As short texts do not provide sufficient word occurrences to use a bag of words. They proposed an approach which uses a small set of domain-specific features to effectively classify texts into predefined sets of generic classes. Experiments are conducted with the available implementation of Naïve Bayes classifier in WEKA3 using 5-fold cross validation which yielded better results than the bag of words approach.
- In the paper[4], the analysis was to classify tweets into data and sentiments(medical field specific) from Twitter. For this they proposed a system to process short text and filter them precisely based on the semantics. The information from tweets are extracted using keyword based knowledge extraction and it is further enhanced using domain specific seed based enrichment technique. For classifying tweets the seed list benefited in such a way that tweets with keyword “Blood glucose” are also classified with keyword “diabetes”

## References

- [1] [Ikonomakis, Emmanouil & Kotsiantis, Sotiris & Tampakas, V.. \(2005\). Text Classification Using Machine Learning Techniques.](#)
- [2] [V, Subramaniaswamy & R., Logesh & Varadarajan, Vijayakumar & Indragandhi, V.. \(2015\). Automated Message Filtering System in Online Social Network.](#)
- [3] [Sriram, Bharath & Fuhry, Dave & Demir, Engin & Ferhatosmanoglu, Hakan & Demirbas, Murat. \(2010\). Short text classification in twitter to improve information filtering.](#)
- [4] [Batool, Rabia & Khattak, Asad & Hashmi, Jahanzeb & Lee, Sungyoung. \(2013\). Precise tweet classification and sentiment analysis.](#)
- [5] [Randa Benkhelifa and Fatima Zohra Laalla. \(2016\) Facebook Posts Text Classification to Improve Information Filtering](#)
- [6] [Mr.Vignesh Kumar , Ms.Dhivya N, Mr.Abishek R K & Mr. Dharun Briram. \(2018\). Post summarization and text classification in social networking sites.](#)