

# Aspect-Sentiment Embeddings for Company Profiling and Employee Opinion Mining

Devamanyu Hazarika<sup>2</sup>, Rajiv Bajpai<sup>1</sup>, Kunal Singh<sup>4</sup>, Sruthi Gorantla<sup>5</sup>, Erik Cambria<sup>3</sup>, and Roger Zimmermann<sup>2</sup>

<sup>1</sup> Accenture, India

<sup>2</sup> National University of Singapore, Singapore

<sup>3</sup> Nanyang Technological University, Singapore

<sup>4</sup> Indian Institute of Technology Kharagpur, India

<sup>5</sup> Indian Institute of Science, Bangalore, India

`dr.rajivbajpai@gmail.com, {devamanyu, rogerz}@comp.nus.edu.sg,  
kunal.singh@iitkgp.ac.in, gorantlas@iisc.ac.in,  
cambria@ntu.edu.sg`

**Abstract.** With the multitude of companies and organizations abound today, ranking them and choosing one out of the many is a difficult and cumbersome task. Although there are many available metrics that rank companies, there is an inherent need for a generalized metric that takes into account the different aspects that constitute employee opinions of the companies. In this work, we aim to overcome the aforementioned problem by generating aspect-sentiment based embedding for the companies by looking into reliable employee reviews of them. We created a comprehensive dataset of company reviews from the famous website **Glassdoor.com** and employed a novel ensemble approach to perform aspect-level sentiment analysis. Although a relevant amount of work has been done on reviews centered on subjects like movies, music, etc., this work is the first of its kind. We also provide several insights from the collated embeddings, thus helping users gain a better understanding of their options as well as select companies using customized preferences.

## 1 Introduction

Emotions, sentiment, and judgments on the scale of good—bad, desirable—undesirable, approval—disapproval are essential for human-to-human communication. Understanding human emotions, deciphering humans’ emotional reasoning and how humans express them in their language is key to enhancing human-machine interaction. In this era of social media, the WWW provides new tools that create and share ideas and opinions with everyone efficiently. Capturing public opinion on social media about events, political movements or any other topics bears a potential for interest amongst the scientific community. That is mainly because of two reasons - firstly, these opinions can help individuals in their decision making process. Secondly, organizations will utilize such

data in order to glean public opinion regarding services and products so as to fine-tune/develop their business strategies.

Sentiment analysis is the study of opinions and sentiments from content that spans from the unimodal to the multimodal [27–29]. Recent approaches to sentiment analysis have focused on the use of linguistic patterns and deep neural networks. The ability to identify aspects within texts is also equally important. An aspect is defined as the product feature; for instance, in the sentence “the battery lasts long”, *battery* is the aspect for which positive sentiment is expressed. The main challenge of sentiment analysis is identifying aspects and their corresponding sentiment. In our case, we do so by merging linguistic patterns and an ELM classifier.

Employees are organizational assets and their opinion plays a vital role in any organization’s growth. Job Search Engines and review websites have evolved to become an ocean of employee reviews. Employee reviews play a vital role for a company’s growth as it improves the relationship between management and the employees via improving staff welfare and morale. These reviews also help prospective employees in selecting a company that meets their criterion. Despite such reviews being important data sources for sentiment mining, they have failed to draw the attention of the scientific community. To the best of our knowledge, only the work of [20] utilized company reviews from Glassdoor. However, even their work was limited to extracting only topics and sentiments from the reviews. In this work, we built a large dataset of employee reviews of companies in Singapore sourced from Glassdoor. Different types of analysis, e.g., aspect extraction, aspect based sentiment analysis were then carried out on this dataset by blending ELM with sentic patterns. To this end, we developed representational embeddings of the companies based on the sentiment score of various different aspects of the companies. In particular, each company is represented in a 30 dimensional space where each dimension corresponds to the average sentiment score of an aspect. Some of these aspects are ‘company culture’, ‘salary’, ‘location’, etc.

In this paper, we used Glassdoor as our source for the preparation of the dataset comprising of 40k reviews. The volume of reviews span a diverse range of aspects (positive and negative) that describe the company based on personal opinions of the reviewers. There are two main contributions of this paper:

- Creation of a large dataset derived from Glassdoor for aspect level sentiment analysis. This dataset contains the reviews of employees working or who previously worked at the corresponding companies.
- Introducing aspect-sentiment embeddings of the companies in order to find similarities between companies in the similar or differing sectors. Aspect-sentiment embeddings project each company onto an n-dimensional space where each dimension corresponds with the overall aspect-sentiment strength of the employees. Aspects are different features of a company, e.g., *salary*, *location*, *work-life balance*, etc. This is particularly useful for job seekers who are looking to find companies that suit their preferences.

The rest of the paper is organized as follows: Section 1 discusses sentiment analysis literature; collection and preparation of the dataset are featured in Sec-

tion 2; we discuss the algorithm details in Section 4; experimental results of this study are presented in Section 5; finally, Section 6 concludes the paper.

## Related works

Identifying emotions associated with employers is just one of the many possible applications of sentiment analysis. We could also analyze industries or professions as a whole, or consider the relationship between the emotional content of reviews with the corresponding salaries of employees. One would expect higher salaries to correlate with more positive emotions, but we might also see an inverse correlation in some cases, perhaps indicating the use of ‘golden handcuffs’.

Sentiment analysis systems can be broadly categorized into knowledge-based [6] or statistics-based systems [7]. Initially, knowledge bases were more commonly used for the identification of emotions and polarity in text. However, at present, sentiment analysis researchers more commonly use statistics-based approaches, with a specific focus on supervised statistical methods. For example, Pang et al. [23] compared the performance of different machine learning algorithms on a movie review dataset: using a large number of textual features, they obtained 82.90% accuracy.

Other unsupervised or knowledge-based approaches to sentiment analysis include Turney et al. [30], which used seed words to calculate the polarity and semantic orientation of phrases, as well as Melville et al. [12] which proposed a mathematical model to extract emotional clues from blogs and then used the information for sentiment detection.

Sentiment analysis research can also be categorized as single-domain [23] versus cross-domain [2]. The work presented in [22] discusses the use of spectral feature alignment to: 1) group domain-specific words from different domains into clusters, and 2) reduce the gap between domain-specific words of two domains using domain independent words. Bollegala et al. [3] developed a sentiment-sensitive distributional thesaurus by using labeled training data from source domain and unlabeled training data from both source and target domains. Some recent approaches [8, 21] used SentiWordNet [1], a very large sentiment lexicon developed by automatically assigning polarity value to WordNet [19] synsets. In SentiWordNet, each synset has three sentiment scores along three sentiment dimensions: positivity, negativity, and objectivity.

As discussed in the introduction, there are hardly any works on mining opinions from company reviews written by employees. Moniz et al. [20] proposed an aspect-sentiment model based on the Latent Dirichlet Allocation (LDA). According to their study, the results of the articulate aspect-polarity model showed that it might be advantageous for investors to combine an appraisal of employee satisfaction with other existing methods for forecasting firm earnings. The research explained and analyzed the sentiments of a stakeholder group which is possibly neglected: the firm’s employees. The researchers initially used online employee reviews in order to capture employee satisfaction and utilized LDA to consider salient aspects in employees’ reviews. From that, they manually derived a latent topic that appeared to be associated with the firm’s outlook. Secondly,

they created an entire document by grouping employee reviews for each firm, and using the General Inquirer dictionary to count positive and negative terms, they measured sentiment as the polarity of the composite document. Their model suggested that employee satisfaction could be formulated as a function of the firm’s outlook and employee sentiment.

## 2 Dataset collection

In our research, we used the popular job recruiting site *Glassdoor.com* as our source to prepare the dataset. The website provides tons of reliable reviews for many companies written mostly by employees, ex-employees or directly associated clients. The content of the reviews possess different aspects (positive and negative) that represent the company from the individual perspectives of the writers. The anonymity of writers enhances the authenticity of the review, thus, this site was our primary source for the dataset collection. The language used in the reviews was found to be highly formal with very minimal usage of slang, decreasing the effort required for data cleaning, normalization, and tokenization. These ‘clean’ words had a high match rate with the dictionary we used [10] for creating the word embeddings, thus improving the performance of the ELM model used in our ensemble network.

The review structure in *Glassdoor.com* consists of a general description followed by both *pros* and *cons* to be written and listed by the writer. This rigid structure aids us in building a balanced dataset comprising of both positive and negative reviews associated with the companies. However, one must be wary of comments like “I really don’t have anything to say/complain about here” in the pros/cons section, which we believe to represent false positives or false negatives respectively. We decided to include these comments in our dataset to represent real-world instances where such false comments are prevalent.

The dataset has been collected by using the official API<sup>6</sup> of Glassdoor. We created a diverse list of 60 well-known companies representing various domains such as technology, finance, energy, hospitality, etc. Finally, we collected a total of 20,000 reviews for all companies. As mentioned earlier, the reviews listed both the pros and cons about the company that is reviewed. Given this sophistication, it was easier for us to split each review into two sub-reviews containing positive and negative opinions respectively. This in turn enabled the automatic labeling of said reviews. Despite doing so, we were careful to manually check the labeling afterwards in order to filter out wrong labels.

Here is an excerpt from one of the positive reviews written for *Accenture*—“*They have great career opportunities, a never ending supply of interesting work, competitive compensation, wonderful benefits, great people, wonderful training programs, a tremendous number of brilliant professionals in their fields ready to help, and great core values*”. This review, like others, is full of important aspects such as *opportunities, compensation, benefits, etc.*, which provides exten-

---

<sup>6</sup> <https://www.glassdoor.com/developer/index.htm>

sive opinions on different facets/characteristics of the company and thus allow our team to prepare a comprehensive review dataset.

Company	Reviews	Company	Reviews
Accenture	1000	HP	150
Adobe	998	HSBC Holdings	1850
Aerostale	874	IBM	150
Aflac	368	Intel Corporation	998
Autodesk	752	Intuit	972
Bank of China	212	Marriot International	980
Booz Allen Hamilton	976	Microsoft	1000
Broadcom	151	Mosanto	396
Brocade	990	Morningstar	733
Camden Property	203	National Instruments	712
Capital One	997	NetApp	800
CarMax	959	Nordstorm	839
Chesapeake Energy	725	OCBC	268
Cisco	980	Paychex	916
Citibank	1852	Qualcomm	919
Colgate-Palmolive	776	Quest Global	150
Creative Technology	150	Rackspace Hosting	732
Darden Restaurants	994	Samsung	150
DBS	288	SCB	942
Devon Energy	336	Singtel	480
DreamWorks Animation	978	Starbucks	828
EOG Resources	127	Starhub	150
FactSet	976	Stryker	904
FedEx Corporation	850	SVB Financial Software	277
Flextronics	150	J.M. Smucker Company	318
General Mills	1036	Ultimate Software	524
Goldman Sachs	970	Umpqua Bank	242
Google	756	Union Overseas Bank	265
Hasbro	458	World Foods Market	757
Herman Miller	540	Yes Bank	176

Table 1: Number of reviews per company

In Table 1, we show the number of reviews per company. The aim was to collect 500-1000 most helpful reviews<sup>7</sup>. However, for some companies the number of reviews available on Glassdoor numbered less than 500.

## 2.1 Preprocessing

As mentioned earlier, the reviews follow a rigid outline structure regardless of writers. Thus, we shuffled the dataset using a pseudo-random generator so as to break any kind of patterns embedded in the dataset. For the purposes of processing the text, we removed any urls, links and hashtags with the use of regular expressions. However, we retained the smileys and emoticons used in the reviews and included them in our vocabulary so as to exploit the emotional and sentimental hints present in them. As we used a context-based algorithm

<sup>7</sup> Reviews for each company in the dataset were selected based on them bearing the most ‘helpful’ review tag provided by [Glassdoor.com](https://www.glassdoor.com)

to create our aspect-sentiment embeddings, retaining these special, non-verbal ‘words’ held a key importance in the performance of the ELM classification. Following this, we used the *NLTK Tokenize Package* to tokenize the reviews into sentences and, finally, into words so as to build up our model’s vocabulary.

### 3 Backgrounds

#### 3.1 Aspect-sentiment Embeddings

In this paper, we introduce Aspect-sentiment Embeddings, which projects companies onto an  $n$ -dimensional space. In particular, each company is given a sentiment strength for each aspect (e.g., salary, work life, location) based on the opinions mined from employee reviews of the company. In mathematical notation we can say,  $(s_1, s_2, \dots, s_n)$  is a vector where  $s_i$  is the sentiment score for aspect  $i$  and there are a total of  $n$  aspects. We constructed such vectors for every company in our dataset, which gave us aspect-sentiment embeddings of the companies.

#### 3.2 Doc2vec for review level embeddings

As we use ensemble architecture to prepare our aspect-sentiment embeddings, the ELM module plays a crucial role as one of the dual paths in the architectural model. To use the ELM module, we need to convert the raw text into review-level summarized embeddings. In our work, we use *Doc2vec* [17] to achieve this task. Doc2vec, also known as paragraph2vec, is a modification of the word2vec algorithm. Word2vec itself is a famous algorithm for word embeddings provided by [18] which trains a neural network to extract contextual-based word embeddings based on the CBOW architecture. Such training has been done on a 100 billion words corpus from Google News and the vectors formed are of 300 dimensionality. In contrast, Doc2vec is an unsupervised learning of continuous representations for larger blocks of text, such as sentences, paragraphs or entire documents. As the reviews in our dataset are very detailed, employing the Doc2vec algorithm is justified. We used the python implementation provided by *gensim*<sup>8</sup> to extract 300 dimensional embeddings to be fed into the ELM model for sentimental analysis and classification.

#### 3.3 Sentiment Dictionary/Lexicons

In order to assign aspect polarity score, we created a dictionary of terms along with their polarities to be used by our ensemble algorithm. We used two primary resources, SentiWordNet [10] and SenticNet [5] to create this dictionary. To filter out irrelevant words that act as noise and would thus fail to provide good polarity to the aspects, we kept an absolute threshold of 0.25 in the polarity scores of the terms in both resources. In order to avoid redundancy, when a particular term is

<sup>8</sup> <https://radimrehurek.com/gensim/>

present in both SentiWordNet and SenticNet, we gave priority to SentiWordNet and chose the term polarity pair from there. Once the dictionary was created, we used it as a reference lookup table in our algorithm mentioned below.

### 3.4 Aspect Level Sentiment Analysis

In opinion mining, different levels of analysis granularity have been proposed, with each having its own advantages and drawbacks [4]. Aspect-based opinion mining [11] and [9] focuses on the relations between aspects and document polarity. An aspect, also known as an opinion target, is a concept in which the opinion is expressed in the given document. For example, the sentence, “The screen of my phone is really nice and its resolution is superb” contains positive polarity for a phone review, i.e., the author likes the phone. However, more specifically, the positive opinion is about its screen and resolution; these concepts are called opinion targets, or, aspects of this opinion. The task of identifying the aspects in a given opinionated text is called aspect extraction.

There are two types of aspects defined in aspect-based opinion mining: explicit aspects and implicit aspects. Explicit aspects are words in the opinionated document that explicitly denote the opinion target. For instance, in the aforementioned example, the opinion targets ‘screen’ and ‘resolution’ are explicitly mentioned in the text. In contrast, an implicit aspect is a concept that represents the opinion target of an opinionated document, but which is not specified explicitly in the text. One can infer that the sentence, “This camera is sleek and very affordable” implicitly contains a positive opinion of the aspects ‘appearance’ and ‘price’ of the entity camera. These same aspects would be explicit in an equivalent sentence: “The appearance of this camera is sleek and its price is very affordable”.

### 3.5 Extreme learning machine (ELM)

For classification, we used ELM as a supervised classifier. The ELM approach [15] was introduced to overcome some issues in back-propagation network [26] training, specifically, potentially slow convergence rates, the critical tuning of optimization parameters [31], and the presence of local minima that call for multi-start and re-training strategies. The ELM learning problem settings require a training set,  $X$ , of  $N$  labeled pairs, using the equation  $(\mathbf{x}_i, y_i)$ , where  $\mathbf{x}_i \in \mathcal{R}^m$  is the  $i$ -th input vector and  $y_i \in \mathcal{R}$  is the associate expected ‘target’ value; using a scalar output implies that the network has one output unit, without loss of generality.

The input layer has  $m$  neurons and connects to the ‘hidden’ layer (having  $N_h$  neurons) through a set of weights  $\{\hat{\mathbf{w}}_j \in \mathcal{R}^m; j = 1, \dots, N_h\}$ . The  $j$ -th hidden neuron embeds a bias term,  $\hat{b}_j$ , and a nonlinear ‘activation’ function,  $\varphi(\cdot)$ ; thus the neuron’s response to an input stimulus,  $\mathbf{x}$ , is:

$$a_j(\mathbf{x}) = \varphi(\hat{\mathbf{w}}_j \cdot \mathbf{x} + \hat{b}_j) \quad (1)$$

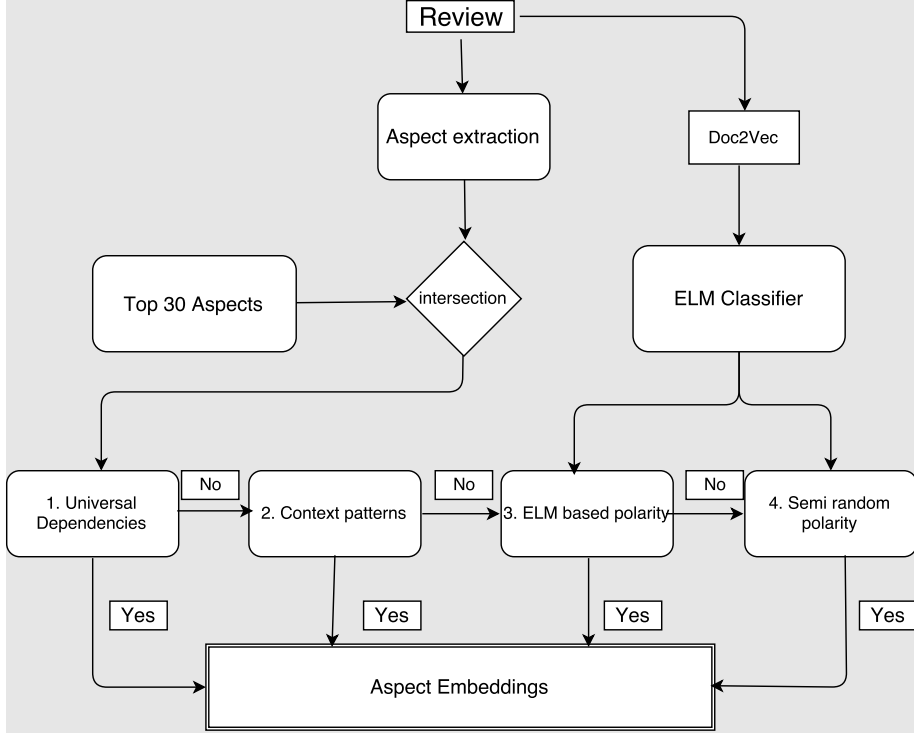


Fig. 1: The flowchart of the algorithm.

Note that (1) can be further generalized to a wider class of functions [14] but for the subsequent analysis this aspect is not relevant. A vector of weighted links,  $\bar{\mathbf{w}}_j \in \mathcal{R}^{N_h}$ , connects hidden neurons to the output neuron without any bias [13]. The overall output function,  $f(\mathbf{x})$ , of the network is:

$$f(\mathbf{x}) = \sum_{j=1}^{N_h} \bar{\mathbf{w}}_j a_j(\mathbf{x}) \quad (2)$$

It is convenient to define an ‘activation matrix’,  $\mathbf{H}$ , such that the entry  $\{h_{ij} \in \mathbf{H}; i = 1, \dots, N; j = 1, \dots, N_h\}$  is the activation value of the  $j$ -th hidden neuron for the  $i$ -th input pattern. The  $\mathbf{H}$  matrix is:

$$\mathbf{H} \equiv \begin{bmatrix} \varphi(\hat{\mathbf{w}}_1 \cdot \mathbf{x}_1 + \hat{b}_1) & \cdots & \varphi(\hat{\mathbf{w}}_{N_h} \cdot \mathbf{x}_1 + \hat{b}_{N_h}) \\ \vdots & \ddots & \vdots \\ \varphi(\hat{\mathbf{w}}_1 \cdot \mathbf{x}_N + \hat{b}_1) & \cdots & \varphi(\hat{\mathbf{w}}_{N_h} \cdot \mathbf{x}_N + \hat{b}_{N_h}) \end{bmatrix} \quad (3)$$

In the ELM model, the quantities  $\{\hat{\mathbf{w}}_j, \hat{b}_j\}$  in (1) are set randomly and are not subject to any adjustment, and the quantities  $\{\bar{\mathbf{w}}_j, \bar{b}\}$  in (2) are the only



degrees of freedom. The training problem reduces to the minimization of the convex cost:

$$\min_{\{\bar{\mathbf{w}}, \bar{b}\}} \|\mathbf{H}\bar{\mathbf{w}} - \mathbf{y}\|^2 \quad (4)$$

A matrix pseudo-inversion yields the unique  $L_2$  solution, as proven in [15]:

$$\bar{\mathbf{w}} = \mathbf{H}^+ \mathbf{y} \quad (5)$$

The simple, efficient procedure to train an ELM therefore involves the following steps:

1. Randomly set the input weights  $\hat{\mathbf{w}}_i$  and bias  $\hat{b}_i$  for each hidden neuron;
2. Compute the activation matrix,  $\mathbf{H}$ , as per (3);
3. Compute the output weights by solving a pseudo-inverse problem as per (5).

Despite the apparent simplicity of the ELM approach, the crucial result is that even random weights in the hidden layer endow a network with a notable representation ability [15]. Moreover, the theory derived in [16] proves that regularization strategies can further improve its generalization performance. As a result, the cost function (4) is augmented by an  $L_2$  regularization factor as follows:

$$\min_{\bar{\mathbf{w}}} \{ \|\mathbf{H}\bar{\mathbf{w}} - \mathbf{y}\|^2 + \lambda \|\bar{\mathbf{w}}\|^2 \} \quad (6)$$

## 4 Detailed Algorithm: Ensemble Architecture

We propose a hybrid algorithm, which works as an ensemble of Unsupervised and Machine Learning approaches, for assigning sentiment labels to the reviews. First, based on the dependency structure of a review we assign polarity to the aspects present in the reviews. This process assumes that the aspect word is connected to a word that is polar and present in the sentiment lexicon. If there is no polar word found connected to the aspect word, according to the sentiment dictionary, we resort to the use of supervised classifier, i.e., ELM. The usage of ELM has multiple benefits like comparable or better performance than other machine learning models like SVMs, and most importantly boasts a significant reduction in model building time. Training time is an important aspect in our work given its high probability to be adapted into an online and real-time application. The flowchart of the proposed algorithm is shown in Figure 1.

### 4.1 Aspect Extraction

We used the aspect extraction method by Poria et al. [25], who proposed a hybrid classifier that uses a Convolutional Neural Network for aspect extraction, as well as linguistic patterns for the purposes of pruning the aspect extraction process. The extracted aspects are given below:

- **Job**, is an umbrella aspect that represents the overall characteristic and nature of the job at that particular company, e.g., *Stable and secure job, good people*
- **Employees/Co-workers**, represents the quality of employees, co workers and the company's relationship with them, e.g., *Very healthy organization with a high-performance culture and very talented employees.*
- **Working time**, clubs aspects of diverse meanings ranging from 'extra-time' to 'time-off' which signify work hours and trends, e.g., *Great people, very generous vacation and time off including sabbaticals every 5 years.*
- **Management**, explores the managerial aspects of the company, e.g., *Great place to work. Inclusive management process. Great products.*
- **Office culture**, summarizes the office environment and the working style, e.g., *Strong focus on procedures, policies, culture, and people. Great benefits.*
- **Location**, represents the comments on location of the company, e.g., *Competitive salary, Nice location, Full freedom.*
- **Work life**, speaks in particular about the type and quality of work, along with the work-life balance present in the company, e.g., *Great office space and location, interesting products to work on.*
- **Salary**, provides information on the salary margins of the company, e.g., *Salary is OK Bonus is good unless there is a food fight. Too laid back (leads to no innovation).*
- **Perks/Benefits**, compensations, bonuses and miscellaneous benefits are included in this aspect, e.g., *Good perks for this type of job - and vary even across levels of employment. Fitness reimbursement, stock options, sabbaticals, etc.*
- **Job opportunities**, scope of the job in the company, e.g., *Strong culture, good reputation, interesting opportunities, management cares about the careers of the employees they are managing.*
- **Employee experience**, lists the sentiments of employees with different degrees of experience and also the quality of experience to be acquired in the company, e.g., *The size of the org can make it difficult for individuals to have their voices heard, especially for new hires, regardless of their experience.*
- **Official staff**, talks about the strength, quality and hospitality of the staff in the company, e.g., *Hours, upper management, lack of staff.*
- **Job training**, expresses the training frameworks and opportunities provided by the company, e.g., *Inconsistent hours, sometimes no hours. No proper training.*
- **Personal growth**, the possibility of technological and experiential growth for the employee, e.g., *Need patience to sense growth since it is a challenging business and changing company.*
- **Leadership**, discusses the role and efficacy of the leadership/senior officials in the areas of motivation and leadership skills. *Leadership does not know how to utilize experienced, professional talent*
- **Politics**, represents the interaction between people for power, e.g., *Politics drive people for more power.*

- **Company business**, explores the business aspects such as performance, trade, etc. of the company, e.g., *Business of the company is booming.*
- **Career development**, forecasts the future of the individuals and company, e.g., *International job within 2 years.*
- **Vacation**, represents the number of holidays the company provides its staff with, e.g., *Paid vacation for the summer. Free travel within the same country.*
- **Company support**, summarizes the quality of interaction between employees, e.g., *Supervisors take care of their employees.*
- **Flexibility**, represents how flexible the company’s environment is, e.g., *Full freedom, work from home allowed.*
- **Performance**, speaks about the type and quality of work done by the employees, e.g., *Extraordinary skills shown by the employees.*
- **Job respect**, shows how employees admire their peers and supervisors, e.g., *Mutual respect amongst the employees.*
- **Work projects**, companies projects, products and plans are included in this aspect, e.g., *Diverse products,numerous projects are provided by the company.*
- **Market viability**, provides information about the type and quality of the market, the company battles, e.g., *Competitive, changing market.*
- **Technology**, explores the technological aspects of the company, e.g., *The machinery used in this company is built on ancient technology.*
- **Work issues**, highlights the operational, legal and internal issues of the company, e.g., *Most of the machines are not working.*
- **Knowledge scope**, lists skills required by the company and knowledge acquired by the employees, e.g., *Technical knowledge in required for this task.*
- **Employee communication**, discusses the interaction between employees and the companies, e.g., *People are very interactive in this company.*
- **Stress**, stress and pressure the employees and companies feel, e.g., *Getting underpaid, stress good for working in a competitive environment.*

We also show the corpus frequency of these aspects in Table 2.

Aspect	Frequency	Aspect	Frequency
Employees/Co-workers	7659	Employee experience	1288
Work Life	7305	Location	627
Perks/Benefits	4565	Leadership	599
Office culture	4192	Technology	490
Working time	3658	Politics	451
Salary	3323	Flexibility	340
Management	2654	Company business	116
Job opportunities	2129		

Table 2: Extracted aspects with their corpus frequency.

## 4.2 Assigning Polarity to the Aspects

In this section, we describe the process of assigning polarity to the aspects.

**Universal Dependent Modifiers** Keeping in mind the goal of finding the polarity score of each aspect (out of the top 30 extracted aspects) present in a review, we start with finding the universal dependencies<sup>9</sup> of aspects in the review. We focus primarily on three dependencies, namely, *adjectival modifier* (*amod*), *adverbial modifier* (*advmod*) and *nominal subject* (*nsubj*).

For better understanding, we provide some examples from reviews in our dataset, with which we demonstrate the aforementioned dependencies associated with the aspects present.

- *Great opportunities for career growth.* *amod*(opportunities, Great)
- *Very political and conservative company. Old school, stodgy.* *advmod*(political, Very)
- *Great people to work with, perks of business traveling.* *nsubj*(travelling, perks)

We use *StanfordCoreNLP Parser* as the tool to extract these universal dependencies. After we find the dependencies, we use these ‘trigger’ words to determine the sentimental polarity of the corresponding aspect. As the aspect sentiment is determined by these modifiers, we lookup their polarities in the prepared sentiment dictionary. These polarities serve as the corresponding aspect score for that particular aspect. This process is repeated for the top 30 aspects that are found in the review.

**Context patterns** In the event that the trigger word is not in the sentiment dictionary, we move on to the second step in the ensemble. Here, we look at the context (window size - 5 words used, including the aspect). We try to find dependency patterns mentioned by [24] and use them to determine the overall polarity score for the aspect. Our assumption is that the presence of highly polar words in the context will contribute to the overall polarity of the aspect. Negations have been appropriately handled as they flip the polarity.

**ELM based polarity score** Should the two procedures mentioned above fail to assign a score to the aspect, we use the prediction made by our ELM model ( 1 - positive, 0 - negative). We directly lookup the polarity of the aspect word in the sentiment dictionary and adjust it based on the ELM output as per the following formula:

$$aspect_{score} = (e_{out}) * (lookup(aspect)) + (1 - e_{out}) * (-1 * lookup(aspect))$$

here,  $e_{out}$  is the output predicted by the ELM for the review,  $lookup(aspect)$  is the polarity score of the aspect word as obtained from the sentiment dictionary.

<sup>9</sup> <http://universaldependencies.org/>

**ELM based semi random score** If the aspect word itself is not present in the sentiment dictionary, we initialize a random polarity value based on the ELM output. The following formula is used to generate the random score:

$$aspect_{score} = \begin{cases} rand(0, 1) & , e_{out} \equiv 1 \\ rand(-1, 0) & , e_{out} \equiv 0 \end{cases}$$

Here,  $rand(a, b)$  is a random generator function which generates random real numbers within the range  $[a, b]$ .

We had to assign a score randomly to the aspects in only 2% of cases. This indicates that the semi-random polarity generation of the aspects did not impact the overall aspect-sentiment embeddings much. However, a fully automatic process is always desirable and as such, we plan on doing so in our future work.

## 5 Experimental Results

In this section, we describe the experimental results and insights drawn from the crawled datasets. Table 6 (Appendix A) shows the aspects and aspect terms that we extracted from the dataset.

We utilized both SVM and ELM models (Table 3) on the dataset in order to detect sentiment. In the experiments, the SVM method performed better than ELM in terms of accuracy. However, the difference between the performances of these classifiers on the given dataset was not statistically significant, as per the paired t-test ( $p > 0.05$ ). Also, in the case of training time, we observed that the ELM method was almost 30 times faster than that of the SVM method on this dataset.

Model	Accuracy	Macro F1-score
SVM	75.13%	74.85%
ELM	74.89%	75.09%

Table 3: Performance of SVM and ELM on the dataset.

### 5.1 Aspect-sentiment Embeddings

In Section 4.2, we described the process of assigning polarity to the aspects. We then calculated the score of an aspect belonging to a company by simply taking the average score of the polarities belonging to that aspect in all reviews of said company.

If we consider each aspect as a separate dimension, and the polarity value of a company for one aspect is the projection along that dimension, then we can

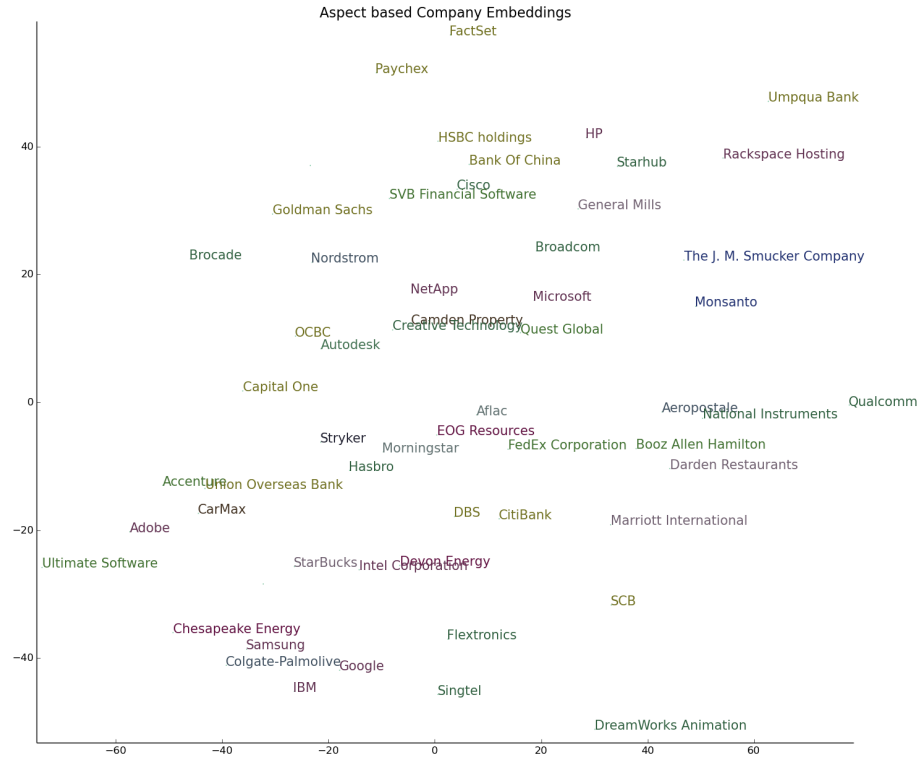


Fig. 2: Projection of the Aspect-sentiment Embeddings of the companies. Note: The same color represents companies from the same sector.

project each company in a  $n$ -dimensional space where  $n = \text{number of aspects}$ . In our case,  $n = 30$ . In Figure 2 we show projection of the companies using aspect-sentiment embeddings.

The motivation for constructing aspect-sentiment embeddings for the companies was to be able to calculate the similarities between companies based on the sentiments of the employees working at those companies.

In Table 4, we present the cosine similarity scores between companies from similar or differing sectors. We see that even though *Goldman Sachs* and *DBS* are in same sector, i.e., Banking and Finance, they have a lower similarity score. However *DBS* and *SCB* (Standard Chartered Bank) have a relatively higher similarity score.

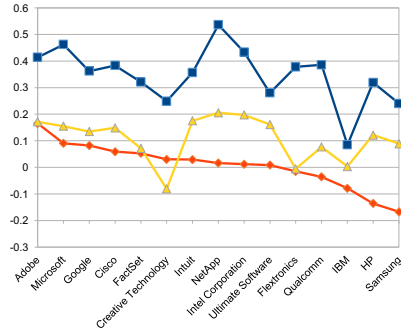
Table 5 presents the best and worst companies in the technological and finance sectors based on sentiment values of the salary, location and work life aspects. Analysis (Figure 3b and 3a) shows that employees are mostly happy with the salary they receive in both the finance and tech sectors. However, in relation to most banks, employees provide negative feedback on work culture.

Company 1	Company 2	Cosine Similarity
Accenture	Booz Allen Hamilton	0.548
Accenture	FedEx	0.733
Google	Microsoft	0.549
Microsoft	Intel	0.486
Adobe	HP	0.370
Adobe	Google	0.276
Adobe	IBM	-0.214
OCBC	Goldman Sachs	0.422
Goldman Sachs	DBS	-0.098
DBS	SCB	0.313
Goldman Sachs	SCB	0.041
Singtel	Broadcom	0.424
Singtel	Starhub	-0.244
Microsoft	Stryker	0.687
National Instruments	Microsoft	-0.117
NetApp	Hasbro	0.655
Monsanto	Quest Global	-0.380

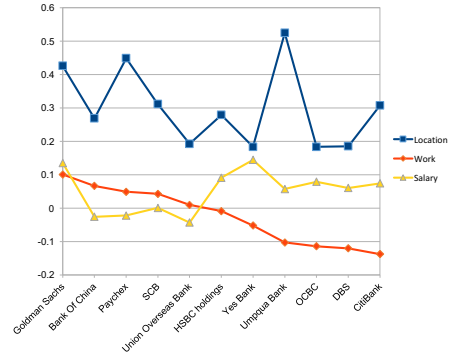
Table 4: Cosine Similarities between the Companies (calculated based on the aspect-sentiment embeddings)

	Location		Salary		Work Life	
	Tech	Finance	Tech	Finance	Tech	Finance
BEST	Microsoft	Umpqua Bank	Intel	Yes Bank	Adobe	Goldman S
	Intel	Goldman S	Adobe	Goldman S	Microsoft	Bank of China
	Adobe	SCB	Microsoft	HSBC	Google	SCB
	Google	Citibank	Cisco	OCBC	Cisco	UOB
	HP	HSBC	Google	Citibank	FactSet	HSBC
WORST	IBM	Yes Bank	Creative	UOB	Samsung	Citibank
	Creative	OCBC	Flextronics	Bank of China	HP	DBS
	NetApp	DBS	FactSet	SCB	NetApp	OCBC
	Cisco	UOB	Samsung	Umpqua Bank	Creative	Umpqua Bank
	FactSet	Bank of China	HP	DBS	Intuit	Yes Bank

Table 5: Companies with best/worst salary and work culture rating in tech and finance sectors.



(a) Tech sector.



(b) Finance sector.

Fig. 3: The plot of the polarity of the aspects of companies in tech and finance sector.

And, although tech companies receive positive feedback for their work culture, the intensity of such positivity is comparatively lower than salary satisfaction.

## 6 Conclusion

In this paper, we described the process of constructing aspect-sentiment embeddings of companies. In particular, we employed aspect-level sentiment analysis on the previously barely researched employee reviews of various companies available on the site Glassdoor. Several experimental insights of the data are given in this study. We addressed the overall employee sentiment on different aspect granularities, e.g., *salary*, *location*, *work life*, etc. This study presents a useful tool for companies to address employees' concerns and increase staff morale. On the other hand, job seekers will also be able to use this study to better find the best employers in their domain of interests.

Future work will mainly focus on considering the rating already given by the users in order to develop a user-product-sentiment model. A more comprehensive aspect-level sentiment analysis is also an important part of this future work.

## Acknowledgment

This research was supported in part by the National Natural Science Foundation of China under Grant no. 61472266 and by the National University of Singapore (Suzhou) Research Institute, 377 Lin Quan Street, Suzhou Industrial Park, Jiang Su, People's Republic of China, 215123.



## References

1. Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. Sentiwordnet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining. In *LREC 2010*, volume 10, pages 2200–2204, 2010.
2. John Blitzer, Mark Dredze, Fernando Pereira, et al. Biographies, bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification. In *ACL 2007*, volume 7, pages 440–447, 2007.
3. Danushka Bollegala, David Weir, and John Carroll. Cross-domain sentiment classification using a sentiment sensitive thesaurus. *IEEE Transactions on Knowledge and Data Engineering*, 25(8):1719–1731, 2013.
4. E. Cambria, B. Schuller, B. Liu, H. Wang, and C. Havasi. Statistical approaches to concept-level sentiment analysis. *IEEE Intelligent Systems*, 28(3):6–9, May 2013.
5. Erik Cambria, Soujanya Poria, Rajiv Bajpai, and Björn Schuller. SenticNet 4: A semantic resource for sentiment analysis based on conceptual primitives. In *COLING*, 2016.
6. Erik Cambria, Bjoern Schuller, Bing Liu, Haixun Wang, and Catherine Havasi. Knowledge-based approaches to concept-level sentiment analysis. *IEEE Intelligent Systems*, 28(2):12–14, 2013.
7. Erik Cambria, Bjoern Schuller, Bing Liu, Haixun Wang, and Catherine Havasi. Statistical approaches to concept-level sentiment analysis. *IEEE Intelligent Systems*, 28(3):6–9, 2013.
8. Kerstin Denecke. Using sentiwordnet for multilingual sentiment analysis. In *Data Engineering Workshop at ICDEW 2008*, pages 507–512. IEEE, 2008.
9. Xiaowen Ding, Bing Liu, and Philip S Yu. A holistic lexicon-based approach to opinion mining. In *2008 International Conference on Web Search and Data Mining*, pages 231–240. ACM, 2008.
10. Andrea Esuli and Fabrizio Sebastiani. Sentiwordnet: A publicly available lexical resource for opinion mining. In *Proceedings of LREC*, volume 6, pages 417–422, 2006.
11. Mingqing Hu and Bing Liu. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 168–177. ACM, 2004.
12. Xia Hu, Jiliang Tang, Huiji Gao, and Huan Liu. Unsupervised sentiment analysis with emotional signals. In *WWW 2013*, pages 607–618, 2013.
13. Guang-Bin Huang. An insight into extreme learning machines: Random neurons, random features and kernels. *Cognitive Computation*, doi: 10.1007/s12559-014-9255-2, 2014.
14. Guang-Bin Huang, Lei Chen, and Chee-Kheong Siew. Universal approximation using incremental constructive feedforward networks with random hidden nodes. *Neural Networks, IEEE Transactions on*, 17(4):879–892, 2006.
15. Guang-Bin Huang, Dian Hui Wang, and Yuan Lan. Extreme learning machines: a survey. *International Journal of Machine Learning and Cybernetics*, 2(2):107–122, 2011.
16. Guang-Bin Huang, Hongming Zhou, Xiaojian Ding, and Rui Zhang. Extreme learning machine for regression and multiclass classification. *Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on*, 42(2):513–529, 2012.
17. Quoc V Le and Tomas Mikolov. Distributed representations of sentences and documents. In *ICML*, volume 14, pages 1188–1196, 2014.

18. Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
19. George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
20. Andy Moniz and Franciska de Jong. Sentiment analysis and the impact of employee satisfaction on firm earnings. In *European Conference on Information Retrieval*, pages 519–527. Springer, 2014.
21. Bruno Ohana and Brendan Tierney. Sentiment classification of reviews using sentiwordnet. In *9th. IT & T Conference*, page 13, 2009.
22. Sinno Jialin Pan, Xiaochuan Ni, Jian-Tao Sun, Qiang Yang, and Zheng Chen. Cross-domain sentiment classification via spectral feature alignment. In *WWW 2010*, pages 751–760. ACM, 2010.
23. Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. Thumbs up?: Sentiment classification using machine learning techniques. In *EMNLP 2002, Volume 10*, pages 79–86. ACL, 2002.
24. Soujanya Poria, Basant Agarwal, Alexander Gelbukh, Amir Hussain, and Newton Howard. Dependency-based semantic parsing for concept-level text analysis. In *Computational Linguistics and Intelligent Text Processing*, pages 113–127. Springer, 2014.
25. Soujanya Poria, Erik Cambria, and Alexander Gelbukh. Aspect extraction for opinion mining with a deep convolutional neural network. *Knowledge-Based Systems*, 2016.
26. Sandro Ridella, Stefano Rovetta, and Rodolfo Zunino. Circular backpropagation networks for classification. *Neural Networks, IEEE Transactions on*, 8(1):84–97, 1997.
27. Rajiv Ratn Shah. Multimodal analysis of user-generated content in support of social media applications. In *Proceedings of the International Conference on Multimedia Retrieval (ICMR)*, pages 423–426. ACM, 2016.
28. Rajiv Ratn Shah, Yi Yu, Akshay Verma, Suhua Tang, Anwar Dilawar Shaikh, and Roger Zimmermann. Leveraging multimodal information for event summarization and concept-level sentiment analysis. *Proceedings of the Knowledge-Based Systems (KBS)*, pages 1–8, 2016.
29. Rajiv Ratn Shah, Yi Yu, and Roger Zimmermann. Advisor: Personalized video soundtrack recommendation by late fusion with heuristic rankings. In *Proceedings of the International Conference on Multimedia (MM)*, pages 607–616. ACM, 2014.
30. Peter D. Turney. Thumbs up or thumbs down?: Semantic orientation applied to unsupervised classification of reviews. In *ACL 2002*, pages 417–424. ACL, 2002.
31. Thomas P Vogl, JK Mangis, AK Rigler, WT Zink, and DL Alkon. Accelerating the convergence of the back-propagation method. *Biological cybernetics*, 59(4-5):257–263, 1988.

## A Top aspects and their aspect-terms

Aspects	Aspect Terms	Aspects	Aspect Terms
Company business	Professional, Booming, Structured Challenging, Competitive, Steady	Career development	Rewarding, International, Guided Challenging, Difficult, Solid
Employee communication	Strong, Transparent, Cryptic, Remote, Awful, Effective	Office culture	Cooperative, Balanced, Exciting, Suffered, Abysmal, Bias
Employees /Co-workers	Excellent, Cooperative, Competent, Stagnant, Unfriendly, Pretend	Employee experience	Diverse, Useful, Firsthand, Horrific, Odd, International
Flexibility	Strict, Dependent, Tremendous, Encourage, Minimal, Great	Personal growth	Exponential, Poor, Constant, Hierarchy, Potential, Constrain
Work issues	Legal, Serious, Inherent, Internal, Operational, Demographic	Overall job	Excellent, Temporary, Overnight, Changing, Tough, Secure
Knowledge scope	Immense, Required, Sharing, Technical, Limited, Vast	Leadership	Appreciate, Strong, Poor, Unwilling, Dedicated, Driving
Location	Strategic, Remote, Accessible, Attractive, Multiple, Uncertain	Management	Flexible, Fluctuate, Inexperienced, Mindful, Dishonest, Focus
Market viability	Changing, Shrinking, Impact, Competitive, Unknown, Successful	Job opportunities	Excellent, Driven, Mindset, International, Lacking, Unique
Perks/Benefits	Unique, Scares, Incredible, Illusion, Lousy, Incentives	Performance	Personal, Necessary, Measurable, Extraordinary, Encouraged, Technical
Politics	Dysfunctional, Drive, Dirty, Extreme, Everywhere, Internal	Work projects	Numerous, Diverse, Challenging, Pushed, Creative, Unbearable
Job respect	Professional, Mutual, Utmost, Solid, Diminishing, Great	Salary	Optimal, Advancement, Hikes, Midrange, Fantastic, Unattractive
Official Staff	Understanding, Excellent, Competent, Mean, Motivated, Dysfunctional	Stress	Underpaid, Excessive, Good, Constant, Additional, Incompetent
Company support	Excellent, Tedious, Benefits, Supervisors, On site, Rare	Technology	Excellent, Ancient, Global, Latest, Green, Innovative
Working time	Exciting, Peak, Tough, Stressful, Extra, Irregular	Job training	Prepares, Competent, Notch, Outstanding, Outdated, Tough
Vacation	Decent, Mandatory, Paid, Planned, Considering, Balance	Work Life	Competent, Mundane, Exciting, Friendly, Versatile, Stressful

Table 6: Top 30 aspects along with their respective aspect terms.