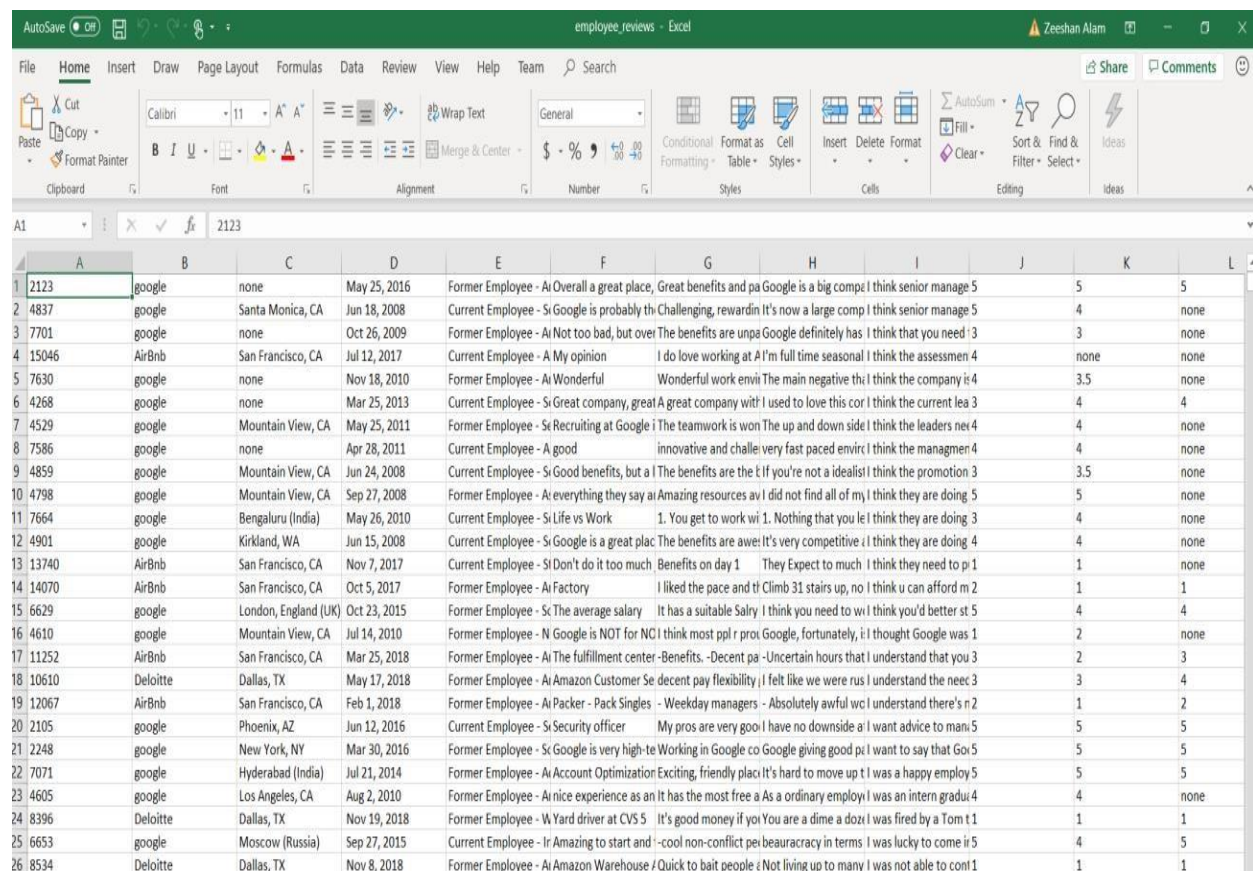


Zeeshan Alam  
1001577559

## Development Phase I: Search

<https://dmsearch.herokuapp.com/>

### Dataset:



ID	Company	Location	Date	Review Text	Rating
2123	google	none	May 25, 2016	Former Employee - Ai Overall a great place, Great benefits and pa Google is a big compe I think senior manage	5
4837	google	Santa Monica, CA	Jun 18, 2008	Current Employee - Si Google is probably th Challenging, rewardin It's now a large comp I think senior manage	4
7701	google	none	Oct 26, 2009	Former Employee - Ai Not too bad, but over The benefits are unpa Google definitely has I think that you need	3
15046	AirBnb	San Francisco, CA	Jul 12, 2017	Current Employee - A My opinion I do love working at A I'm full time seasonal I think the assessmen	4
7630	google	none	Nov 18, 2010	Former Employee - Ai Wonderful Wonderful work envii The main negative th: I think the company is	3.5
4268	google	none	Mar 25, 2013	Current Employee - Si Great company, great A great company with I used to love this cor I think the current lea	4
4529	google	Mountain View, CA	May 25, 2011	Former Employee - Si Recruiting at Google i The teamwork is won The up and down side I think the leaders nei	4
7586	google	none	Apr 28, 2011	Current Employee - A good innovative and challei very fast paced envirc I think the managem	4
4859	google	Mountain View, CA	Jun 24, 2008	Current Employee - Si Good benefits, but a l The benefits are the t If you're not a idealis I think the promotion	3.5
4798	google	Mountain View, CA	Sep 27, 2008	Former Employee - Ai everything they say ai Amazing resources av I did not find all of mj I think they are doing	5
7664	google	Bengaluru (India)	May 26, 2010	Current Employee - Si Life vs Work 1. You get to work wi 1. Nothing that you le I think they are doing	3
4901	google	Kirkland, WA	Jun 15, 2008	Current Employee - Si Google is a great plac The benefits are awe: It's very competitive i I think they are doing	4
13740	AirBnb	San Francisco, CA	Nov 7, 2017	Current Employee - Si Don't do it too much Benefits on day 1 They Expect to much I think they need to p	1
14070	AirBnb	San Francisco, CA	Oct 5, 2017	Former Employee - Ai Factory I liked the pace and ti Climb 31 stairs up, no I think u can afford m	2
6629	google	London, England (UK)	Oct 23, 2015	Former Employee - Si The average salary It has a suitable Salry I think you need to wi I think you'd better st	5
4610	google	Mountain View, CA	Jul 14, 2010	Former Employee - N Google is NOT for NC I think most ppl r pro Google, fortunately, i I thought Google was	1
11252	AirBnb	San Francisco, CA	Mar 25, 2018	Former Employee - Ai The fulfillment center -Benefits -Decent pa -Uncertain hours that I understand that you	2
10610	Deloitte	Dallas, TX	May 17, 2018	Former Employee - Ai Amazon Customer Se decent pay flexibility I felt like we were rus I understand the neec	3
12067	AirBnb	San Francisco, CA	Feb 1, 2018	Former Employee - Ai Packer - Pack Singles - Weekday managers - Absolutely awful wc I understand there's n	2
2105	google	Phoenix, AZ	Jun 12, 2016	Current Employee - Si Security officer My pros are very goo I have no downside a I want advice to mani	5
2248	google	New York, NY	Mar 30, 2016	Former Employee - Si Google is very high-te Working in Google co Google giving good pi I want to say that Goi	5
7071	google	Hyderabad (India)	Jul 21, 2014	Former Employee - Ai Account Optimization Exciting, friendly plac It's hard to move up t I was a happy employ	5
4605	google	Los Angeles, CA	Aug 2, 2010	Former Employee - Ai nice experience as an It has the most free a As a ordinary employ I was an intern gradi	4
8396	Deloitte	Dallas, TX	Nov 19, 2018	Former Employee - W Yard driver at CVS 5 It's good money if you You are a dime a doz I was fired by a Tom t	1
6653	google	Moscow (Russia)	Sep 27, 2015	Current Employee - Ir Amazing to start and i -cool non-conflict pei beauracracy in terms I was lucky to come ir	5
8534	Deloitte	Dallas, TX	Nov 8, 2018	Former Employee - Ai Amazon Warehouse i Quick to bait people t Not living up to many I was not able to cont	1

### Implementation:

#### TF-IDF:

TF-IDF is a weighting factor for feature, the more weight the more that type of term occurs in the document which is offset by the number of times the words appear in the entire document which helps removing really common words (stop-words) in the language.

$$\text{tf-idf}(t, D) = \text{tf}(t, d) * \text{idf}(t, D)$$

$$\text{tf}(t, d) = f(t|d)$$

# the number of time single words appear in a given specific documents

$\text{idf}(t,D) = \log(N / |\{d \in D : t \in d\}|)$

the log is used to dampen the effects of IDF function

Smoothing for IDF:  $\log(1 + N / |\{d \in D : t \in d\}|)$

The smoothing factor is used for introducing the lower bound of log(2) so that nothing will ever be multiplied by 0 by the IDF

```
12
13 def index(Query):
14     data = pd.read_csv('employee_reviews.csv')
15     text = data['Pros']
16     stop_words = stopwords.words('english')
17
18     def process_text(text):
19         text = str(text)
20         text = re.sub('[^a-z\s]', '', text)
21         text = [w for w in text.split() if w not in set(stop_words)]
22         return ' '.join(text)
23
24     data['Pros'] = data['Pros'].apply(process_text)
25
26     english_stemmer = SnowballStemmer('english')
27     analyzer = CountVecorizer().build_analyzer()
28
29     def stemming(text):
30         return (english_stemmer.stem(w) for w in analyzer(text))
31
32     count = CountVectorizer(analyzer=stemming)
33
34     count_matrix = count.fit_transform(data['Pros'])
35
36     tfidf_transformer = TfidfTransformer()
37     train_tfidf = tfidf_transformer.fit_transform(count_matrix)
38
39     def get_search_results(query):
40         query = process_text(query)
41         query_matrix = CountVecorizer().build_analyzer().transform([query])
42         query_tfidf = tfidf_transformer.transform(query_matrix)
43         sim_score = cosine_similarity(query_tfidf, train_tfidf)
44         sorted_indexes = np.argsort(sim_score).tolist()
45         return data['Company'].iloc[sorted_indexes[0] [-10:]]
46
47     working = get_search_results(Query)
```

## Challenges:

There were many challenges in my path to complete the search feature. First and foremost was to learn and understand properly about the search algorithm used. Initially my idea was to go for a mobile application but later I found that Windows 10 Home does not support Hyper-V and I was not able to use emulator for testing and building purposes and on the top of that not having an android phone made it even tougher to work on the mobile application Therefore; I switched to web application using python Django framework.

I had to learn the framework by following different tutorials and the structure for the Django framework was pretty unique to all the other frameworks I have used in the past. After learning and understanding about the Django framework I was done with my web application front-end and I used libraries (python) and followed tutorials on the internet to finish my search feature.

Final step in the process was to find the public server to deploy my application for that I choose Heroku as it is one of the well-known free option available in the market. There were many issues when deploying the web application to the server and I had to check the log files every time to check errors in my web application.

## References:

<https://medium.freecodecamp.org/how-to-process-textual-data-using-tf-idf-in-python-cd2bbc0a94a3>

<https://towardsdatascience.com/tfidf-for-piece-of-text-in-python-43feccaa74f8>

<https://github.com/Heetmadhu/Movie-Recommendation/blob/master/MovieSearch.ipynb>

[https://scikit-learn.org/stable/modules/generated/sklearn.feature\\_extraction.text.TfidfVectorizer.html](https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html)