



# **Analysis of skills affected due to Automation and AI**

Instructor: Rong Liu

Course: BIA-660-C

Team Members: Group 5

Rohan Shah

Parth Parab

Purva Khopkar

Bhagyashree Shende





# Introduction

Goal of the project: To cluster skills and activities from various job postings scraped from “careerbuilder.com”

- Cluster skills and activities from various job postings
- Analyze the activities in the cluster to predict which clusters may be impacted due to Automation and AI

**CAREERBUILDER<sup>®</sup>**  
WORK CAN WORK™

Jobs Upload/Build Resume Career Development & Learning ▾ Sign In | Sign Up | Pos

Customer Service in US  Job Type ▾ Date Posted ▾ Pay ▾ Easy Apply Only

**More Than 5,000 Jobs Found** [Save Search](#)

Sort by: Relevancy | Date

**TODAY Customer Service Representa...**   
Tempe, AZ | Full Time  
\$20 - \$23 / hour [Easy Apply](#)

**TODAY Retail Customer Service Representative**   
San Antonio, TX | Full Time  
\$10 - \$15 / hour [Easy Apply](#)

**TODAY CUSTOMER SERVICE REPRES...**   
OfficeTeam | Allentown, PA | Seasonal / Temp [Easy Apply](#)

**TODAY Customer Service Representa...** 

area. The Retail Customer Service Representative will interact with customers to answer general inquiries, resolves issues, and handle customer complaints. The ideal candidate will maintain a professional image both in our office and in our clients' atmospheres and uphold a can-do attitude at all times.

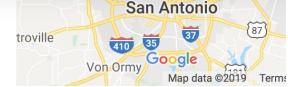
**Retail Customer Service Representative Responsibilities:**

- Provide excellent customer service to all clients/customers and maintain adequate appointment times
- Report any new customer acquisitions or account changes to the Director of Operations
- Resolve customer support related issues and provide customers with proper solutions to their concerns

- At least 1+ year experience in an office, retail, or other customer service environment
- Strong communication skills
- Basic office etiquette
- Ability to juggle multiple tasks simultaneously

**Recommended skills**

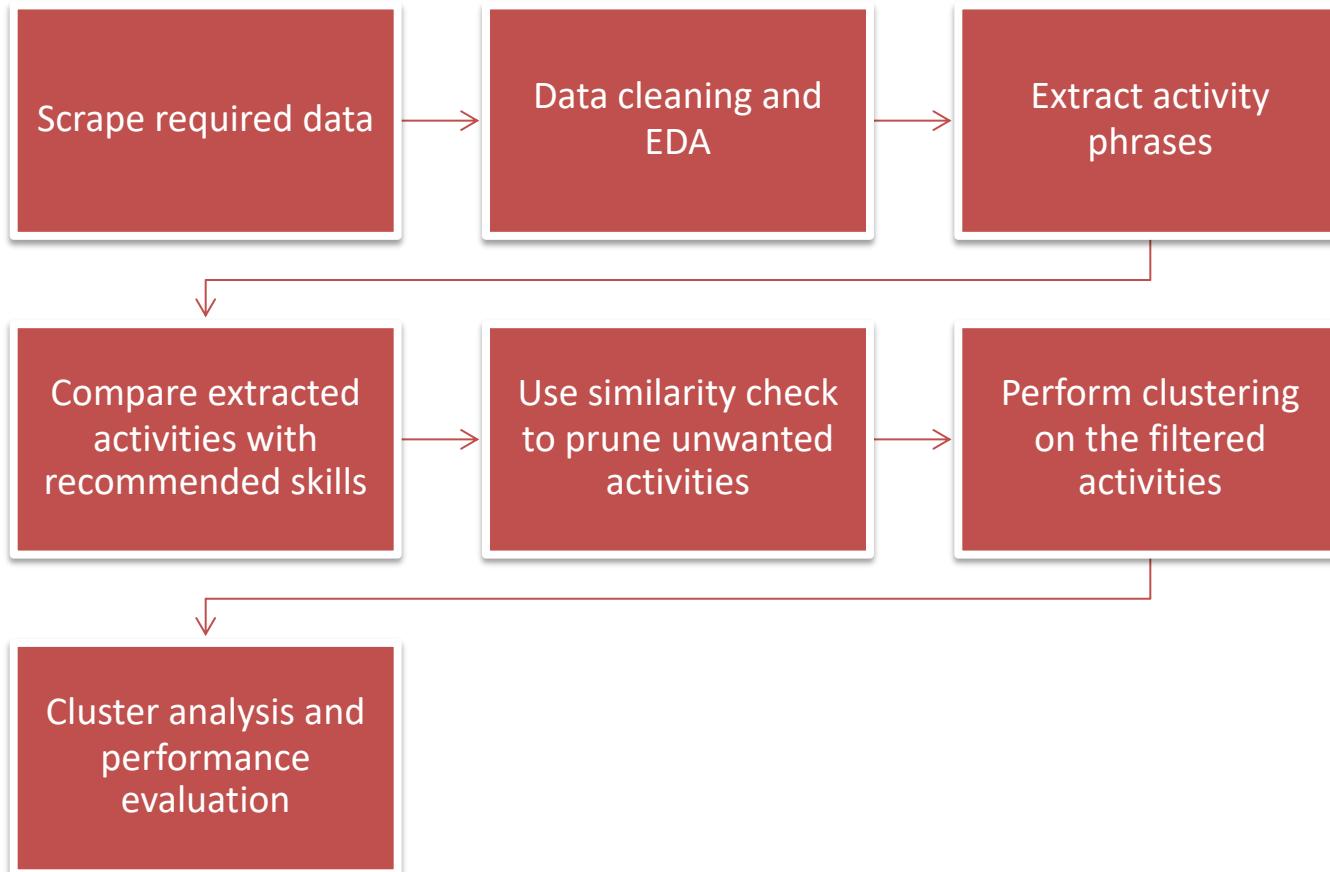
Professional Attitude Operations Customer Service Positive Attitude  
Setting Appointments Mergers And Acquisitions

  
Map data ©2019 Google Terms

**Salary Details**  
This salary was provided in the Job Posting.

**\$10-\$15**  
Hourly Salary

# Project Workflow





**STEVENS**  
INSTITUTE of TECHNOLOGY  
THE INNOVATION UNIVERSITY®



# Web Scraping and Data Preprocessing





# Web Scraping

- Scrapped first 100 pages of Careerbuilder having 25 jobs each page to get details 10000 Jobs for 4 types of keywords namely Software Engineering, Manufacturing, Fabrication, Customer Service and store the data into file using BeautifulSoup

The image shows two side-by-side code editors, both titled "test.py". The left editor contains a script for extracting job data from careerbuilder.com. It uses the pandas library to manage data frames and the requests library to get the raw HTML content. The script iterates through 100 pages of jobs, extracts job titles, descriptions, locations, companies, and posting dates, and appends this data to a DataFrame. The right editor shows the continuation of this script, specifically handling the job description and salary information. It uses BeautifulSoup to parse the HTML and extract specific div elements. The job description is checked for length; if it's too short, it's set to "N/A". The salary information is also extracted and appended to the DataFrame. Both scripts use the same basic structure of importing libraries, defining URLs, and parsing the response content.

```
Editor - C:\Users\yohan\Downloads\test.py
Editor - C:\Users\yohan\Downloads\test.py

1 import pandas
2 import requests
3 import bs4
4
5 BASE_URL_careerbuilder = 'https://www.careerbuilder.com/jobs'
6
7 rawcode = requests.get(BASE_URL_careerbuilder)
8 soup = bs4.BeautifulSoup(rawcode.content, "html.parser")
9
10 num_pages_careerbuilder=100;
11 job_df_careerbuilder = pandas.DataFrame()
12
13 for i in range(1, num_pages_careerbuilder+1):
14
15     url = ''.join([BASE_URL_careerbuilder, '?page_number=', str(i)])
16     rawcode = requests.get(url)
17     soup = bs4.BeautifulSoup(rawcode.content, "html.parser")
18     job_divs = soup.select("div#jobs_collection div.data-results-content-parent")
19     temp_URL_careerbuilder = 'https://www.careerbuilder.com/'
20     for each in job_divs:
21
22         job_location_span = each.select("div.data-details span")
23         if(len(job_location_span) ==3):
24             job_location= job_location_span[1].get_text()
25
26         else:
27             job_location= "N/A"
28
29         job_company = job_location_span[0].get_text()
30
31         posted_date = each.select("div.data-results-publish-time")[0].get_text()

40
41         job_title = soup.select("div.data-display-header_info-content h1")[0].get
42
43         job_description_span= soup.select("div.col-2 div.col.big.col-mobile-full")
44
45         if(job_description_span != None and len(job_description_span) > 0):
46             job_description = job_description_span[0].get_text()
47             if(job_description==""):
48                 job_description = "N/A"
49
50         else:
51             job_description = "N/A"
52
53         salary_div=soup.select("div.data-snapshot div.block")
54         if(salary_div != None and len(salary_div) > 0 ):
55             salary = salary_div[0].get_text()
56         else:
57             salary = "Not Disclosed"
58         job_df_careerbuilder = job_df_careerbuilder.append({'job_title': job_titl
59
60
61
62
63
64
65
66 cols=['from','date','job_title','job_company','job_location','skills','salary']
67 job_df_careerbuilder = job_df_careerbuilder[cols]
68
69 path = 'C:\\\\Users\\\\pkhopkar\\\\Downloads\\\\BIA660\\\\BIA660\\\\' + 'job_careerbuilder'
70 job_df_careerbuilder.to_csv(path)
71
```



# Scrapped data from Careerbuilder

A	B	C	D	E	F	G	H	I
1	from	date	job_title	job_company	job_location	skills	salary	
2	0	Careerbui	Today	Branch Manager of DME Operations	Saint Paul	N/A	['Management', 'Recruitment', 'Accounting', 'Navigation', 'Medicaid', 'Training']	
3	1	Careerbui	1 day ago	Customer Service - Several Representatives	Riverside	N/A	['Sales Management', 'Office Management']	Starting \$900 - \$1200 Weekly (Based on Qualifications)
4	2	Careerbui	Today	Sales Support Specialist	Corps Team	Pearl River	['Payment Processing', 'Scheduling', 'Billing', 'Management', 'Accounts Receivable', 'Business Development']	
5	3	Careerbui	Today	Admin Clerk	Manpower	Charlotte	['Databases', 'Fax', 'Binding', 'Distribution']	\$3,000.00 / year
6	4	Careerbui	Today	Stocker	Manpower	Huntington Bay	['Unpacking', 'Loss Prevention', 'Merchandise Management']	\$15.00 / hour
7	5	Careerbui	Today	Accounting Assistant	Ocean Reef Community A	North Key Largo	['Accounts Receivable', 'Accounts Payable', 'Deposit Accounts', 'Accounting', 'Finance']	
8	6	Careerbui	Today	Warehouse Associate	Manpower	Louisville	['Solid Works (Cad)', 'Scheduling', 'Inventory Management']	\$12.10 - \$14.50 / hour
9	7	Careerbui	Today	Document Control Administrator II (Associate)	Yoh	Elkton	['Enterprise Document Management System', 'Filing', 'Data Entry', 'Document Management']	
10	8	Careerbui	Today	Wood Finisher - Pine Lawn - 1st shift	Manpower	Pine Lawn	['Sander (Metalworking Tools)', 'Maintenance']	
11	9	Careerbui	Today	Data Entry Clerks	Manpower	Mason	['Data Entry', 'Filing', 'Multitasking', 'Customer Service']	\$15.00 - \$15.50 / hour
12	10	Careerbui	Today	Freight Handler	FHI Group	Tipp City	['Palletizing', 'English Language', 'Mathematics']	
13	11	Careerbui	Today	Production Associate	Manpower	North Kansas City	['Drug Testing', 'Background Checks']	\$10.00 - \$10.50 / hour
14	12	Careerbui	Today	Selector	FHI Group	Collierville	['English Language', 'Mathematics']	\$14.00 / hour
15	13	Careerbui	Today	Dispatch Associates - Louisville, KY	Manpower	Louisville	['Scheduling', 'Clerical Works', 'Recruitment']	\$16.00 - \$17.00 / hour
16	14	Careerbui	Today	Maintenance	Manpower	Neosho	['Troubleshooting (Problem Solving)', 'Self Motivation', 'Poultry', 'Healthy', 'Preventive Maintenance']	
17	15	Careerbui	Today	Chocolate Candy Manufacturing - 3rd shift	Manpower	Scranton	['Quality Assurance', 'Corrective And Preventive Action']	\$12.00 / hour
18	16	Careerbui	Today	Systems Engineer /Support	ConsultNet	New York	['Python (Programming Language)', 'Confluence (Physical Geography)', 'Windows Servers', 'Linux']	
19	17	Careerbui	Today	Machine Operators - Hilliard, OH	Manpower	Hilliard	['Manufacturing', 'Machinery']	\$15,000.00 - \$18,000.00 / year
20	18	Careerbui	Today	Food Operator	Manpower	Luverne	['Cargos', 'Warehousing', 'Manifests', 'Warehouse Management']	\$8.00 - \$8.25 / hour
21	19	Careerbui	Today	General Labor - Hilliard, OH	Manpower	Hilliard	['Warehousing']	\$13.00 - \$15.00 / hour
22	20	Careerbui	Today	Extrusion Worker (Sanford, Maine)	Manpower	Sanford	['Standard Operating Procedure', 'Microscopy']	\$18.00 / hour



# Data Pre-processing

```
In [25]: 1 #removing columns Unnamed and From based on relevance
2 data = data.drop(columns=['Unnamed: 0'])
3 data = data.drop(columns=['from'])
4 data.shape

Out[25]: (6583, 7)
```

```
In [26]: 1 #replacing 'day' with 'days' in date for better processing
2 data1 = data["date"].replace(to_replace ="1 day ago",
3 value = "1 days ago")
4 data["date"] = data1
5 data.date[124]

Out[26]: '1 days ago'
```

```
In [27]: 1 #replacing date with numerical date values
2 for i in range(len(data)):
3     if(data.date[i] == "Today"):
4         data1 = data["date"].replace(to_replace ="Today",
5         value = datetime.datetime.now().strftime("%x"))
6         data["date"] = data1
7
8     elif((data.date[i][-1]) == 'o' and data.date[i] != "30+ days ago"):
9         today = datetime.datetime.now()
10        daysAgo = int((data.date[i]).split(' days ago')[0])
11        DD = datetime.timedelta(days= daysAgo)
12        earlier = today - DD
13        earlier_str = earlier.strftime("%x")
14        data1 = data["date"].replace(to_replace = data.date[i],
15        value = earlier_str)
16        data["date"] = data1
17
18 data.head()
```

Removing irrelevant columns and standardizing date format



# Data Pre-processing

```
In [43]: 1 #removing rows if column 'job_description' value is Null  
2 data = data[pd.notnull(data['job_description'])]  
3 #data.head()
```

```
In [29]: 1 #seperating list of skills into seperate rows for simplicity  
2 dataTest = data.copy()  
3 x = dataTest["skills"].apply(lambda s: list(ast.literal_eval(s)))  
4 dataTest["skills"] = x  
5 lst_col = 'skills'  
6 dataTest.skills.apply(pd.Series)  
7 dataTest = pd.DataFrame({  
8     col:np.repeat(dataTest[col].values, dataTest[lst_col].str.len())  
9     for col in dataTest.columns.difference([lst_col])  
10 }).assign(**{lst_col:np.concatenate(dataTest[lst_col].values)})(dataTest.columns.tolist())  
11 #data = dataTest[['date','job_title','job_company','skills','job_location']]
```

```
In [36]: 1 dataTest.head()
```

Out[36]:

	date	job_title	job_company	job_location	skills	salary	job_description
0	12/05/19	Senior Software Developer	Kforce Technology	Houston, TX	Advanced Audio Coding (Aac)	\n115,000.00–120,000.00 / year\n	Attending the agile ceremonies, performing so...
1	12/05/19	Senior Software Developer	Kforce Technology	Houston, TX	C Sharp (Programming Language)	\n115,000.00–120,000.00 / year\n	Attending the agile ceremonies, performing so...
2	12/05/19	Senior Software Developer	Kforce Technology	Houston, TX	C++ (Programming Language)	\n115,000.00–120,000.00 / year\n	Attending the agile ceremonies, performing so...
3	12/05/19	Senior Software Developer	Kforce Technology	Houston, TX	Software Development Life Cycle	\n115,000.00–120,000.00 / year\n	Attending the agile ceremonies, performing so...
4	12/05/19	Senior Software Developer	Kforce Technology	Houston, TX	.Net Framework	\n115,000.00–120,000.00 / year\n	Attending the agile ceremonies, performing so...

Dropping rows with null values based on column 'job\_description' and separating list of skills to separate rows

# Data Pre-processing

```
In [32]: 1 #Pre-processing for Graph for top 20 skills with maximum and minimum average salary
2 #getting salary values to float type
3 for i in range(len(eda_sal)):
4     if("Careerbuilder est." in eda_sal["salary"].iloc[i] and "$" in eda_sal["salary"].iloc[i]):
5         val = eda_sal["salary"].iloc[i].split('Careerbuilder est.'))
6         strip = val[0].strip().replace(',', '')
7         eda_sal["salary"].iloc[i] = float(strip[1::])
8         data = {'skills':eda_sal["skills"].iloc[i], 'salary':eda_sal["salary"].iloc[i]}
9         df.loc[len(df.index)] = list(data[0].values())
10
11 elif(" / year" in eda_sal["salary"].iloc[i] and "$" in eda_sal["salary"].iloc[i]):
12     x = eda_sal["salary"].iloc[i].strip().split(' / year')
13     val = x[0].replace('$', '').split(' - ')
14     eda_sal["salary"].iloc[i] = float(val[1].replace(',', ''))
15     data = {'skills':eda_sal["skills"].iloc[i], 'salary':eda_sal["salary"].iloc[i]}
16     df.loc[len(df.index)] = list(data[0].values())
17
18 elif(" / hour" in eda_sal["salary"].iloc[i] and "$" in eda_sal["salary"].iloc[i]):
19     x = eda_sal["salary"].iloc[i].strip().split(' / hour')
20     val = x[0].replace('$', '').split(' - ')
21     eda_sal["salary"].iloc[i] = (float(val[1].replace(',', '')) * 40 * 52)
22     data = {'skills':eda_sal["skills"].iloc[i], 'salary':eda_sal["salary"].iloc[i]}
23     df.loc[len(df.index)] = list(data[0].values())
```

```
In [33]: 1 df.head()
```

Out[33]:

	skills	salary
0	Advanced Audio Coding (Aac)	120000.0
1	C Sharp (Programming Language)	120000.0
2	C++ (Programming Language)	120000.0
3	Software Development Life Cycle	120000.0
4	.Net Framework	120000.0

Getting salary values to float type



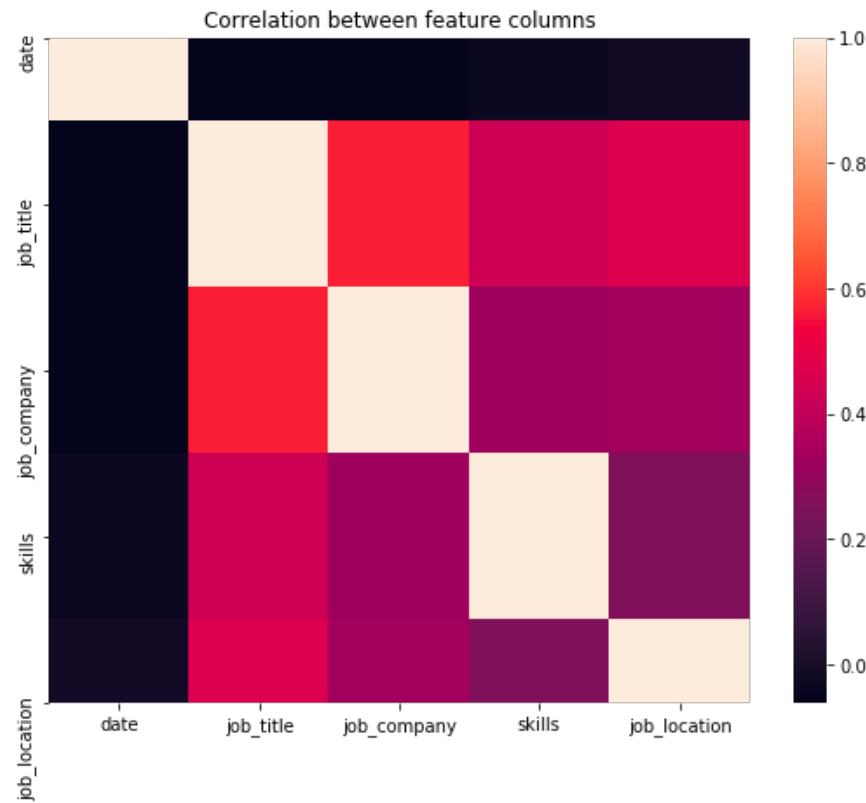
**STEVENS**  
INSTITUTE of TECHNOLOGY  
THE INNOVATION UNIVERSITY®



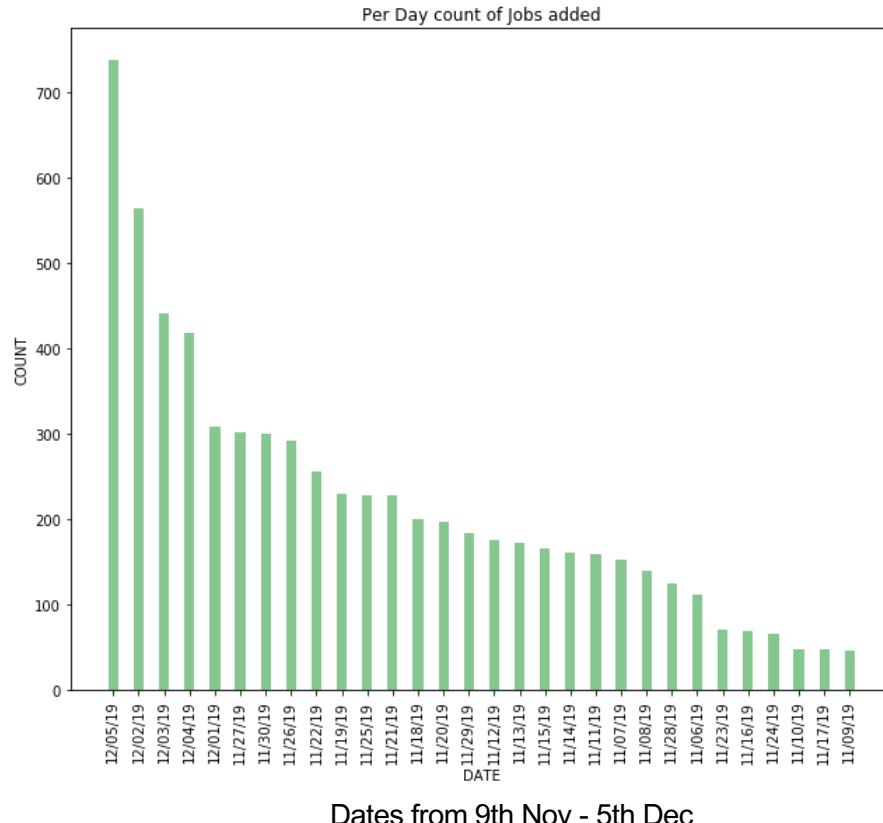
# Exploratory Data Analysis (EDA)



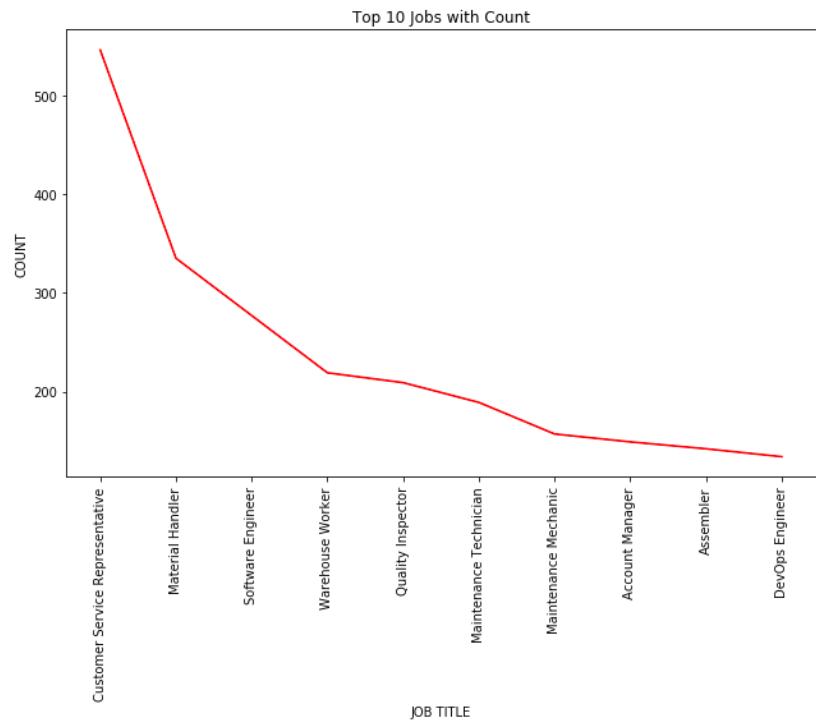
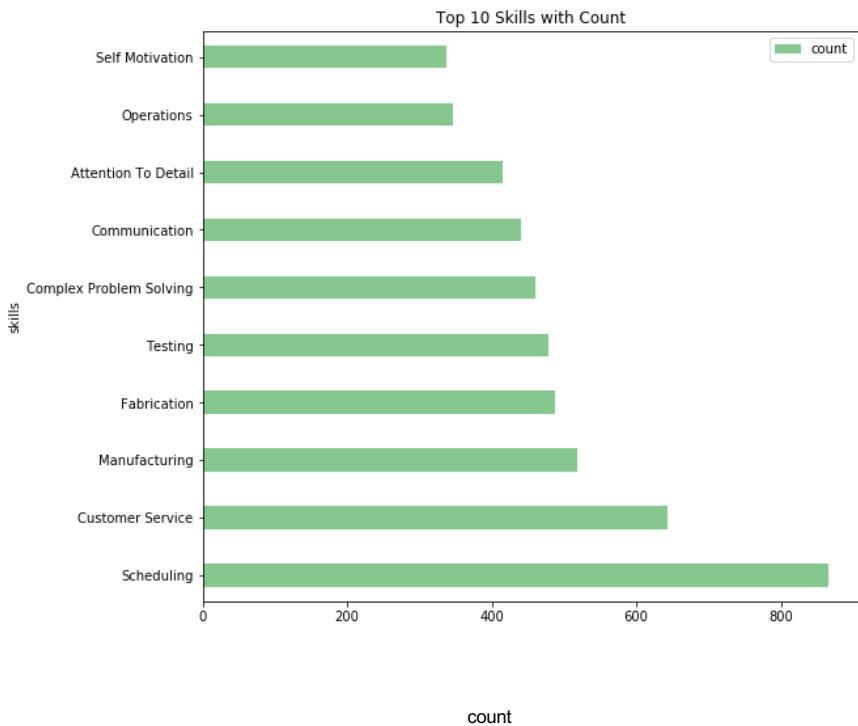
# Exploratory Data Analysis (EDA)



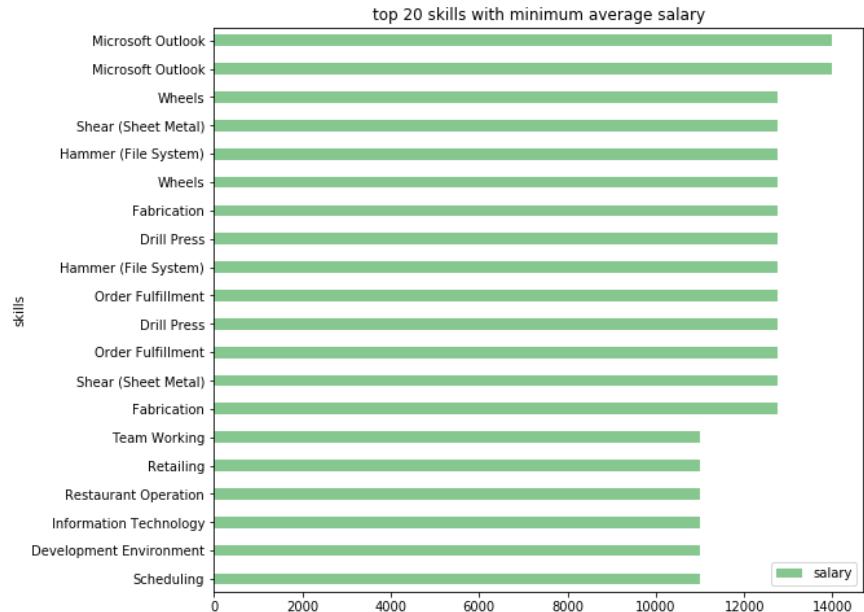
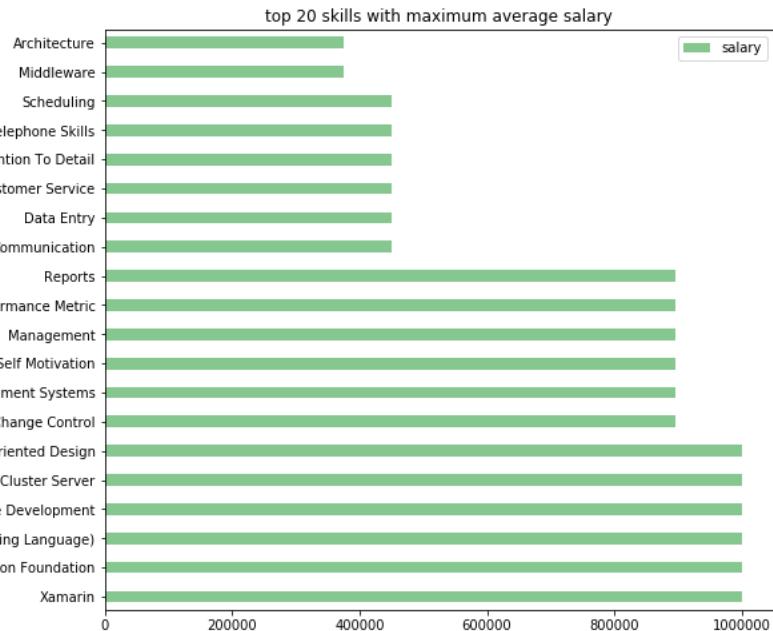
# Exploratory Data Analysis (EDA)



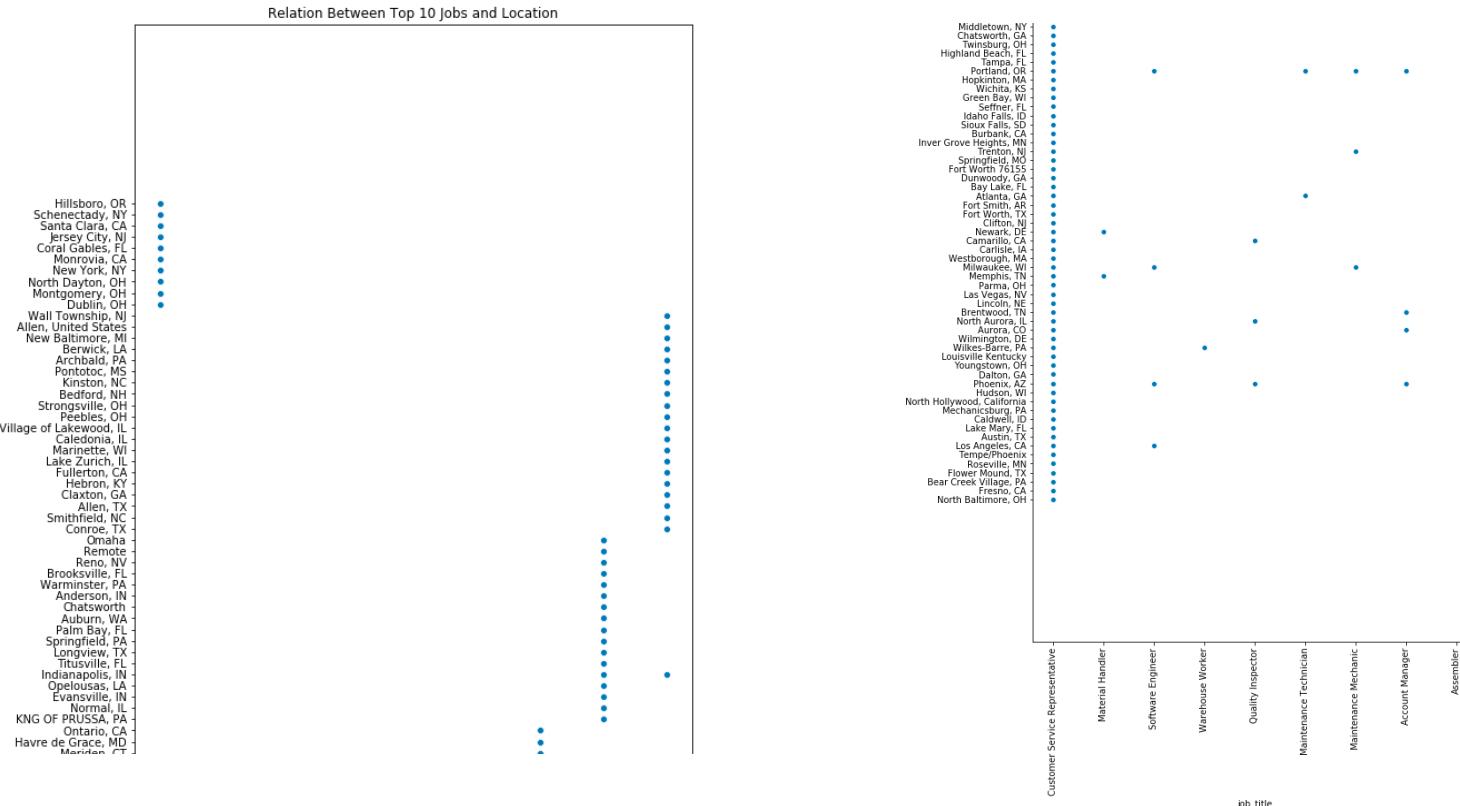
# Exploratory Data Analysis (EDA)



# Exploratory Data Analysis (EDA)



# Exploratory Data Analysis (EDA)





# Extract Activity Phrases

Why extract activity phrases??

- Recommended skills listed are abstract
- Give no idea about the context

Example:

Recommended Skill is “Analysis”

Doesn’t give clear picture of what is actually expected

Job Details

Company Overview

Our client, a fortune 500 company located in the Twin Cities, is looking to add a Customer Service Representative to their fast growing team. This is a golden opportunity for a career minded individual to get their foot in the door with a fortune 500 company. Our client, an industry leading company located in the Twin Cities, is looking to add a Customer Service in the door with a company invested in your success. Job Duties:

- Works with existing and potential clients in answering product and service questions
- Quickly and accurately enter customer information
- Maintains customer records by updating account information
- Resolves product or service problems
- Recommends potential products or services to management by collecting customer information and analyzing customer needs
- Contributes to team effort by accomplishing related results as needed
- Skills/Qualifications: Detail oriented
- Customer service mindset with a passion for helping others
- Previous experience in a customer service or sales environment preferred

If this sounds like the type of opportunity you are looking for, please apply today. We are an Equal Employment Opportunity employer committed to excellence, diversity and inclusion.

## Recommended skills

Analysis   Inclusion   Hardworking And Dedicated   Management  
Customer Service   Passionate



# Extract Activity Phrases

## Implementation Idea

### Part of Speech Tagging

- Using Bag of words approach doesn't capture structure of sentence
- It is also important to understand context
- Assigns part of speech to each word according to input sentence

### Chunking

- Chunking works on top of POS tagging, it uses pos-tags as input and provides chunks as output.



# Phrase Generator

Concept Used – Part of Speech Tagging

1. Part of Speech Tagging – Reads text and assign Part of Speech to each word
2. Examples of POS Tags – (NN, Noun) (VB, Verb) (JJ, Adjective) (RB, Adverb) etc.

Python library used: NLTK

## Steps Involved:

- Input : Job Description as text
- Tokenize text using word\_tokenize()
- Remove stop words
- apply pos\_tag to above step that is nltk.pos\_tag(tokenize\_text)



# Phrase Generator – Contd.

## Concept Used - Chunking

1. Chunking - add more structure to the sentence by following parts of speech (POS) tagging
2. primary usage of chunking is to make a group of "noun phrases."

### Rules for Chunking:

- No pre-defined rules, we can combine them according to need and requirement.
- Our rule – {<JJ.?>+<NN|NNS|NNP.?>+<NN|NNS|NNP.?>\*}}
- Examples of chunks returned:  
[ ‘strong experience JAVA’, ‘verbal communication skills’, ‘punctual delivery’, ‘trade-off analysis’ , ‘operational support’ ]



# Similarity Check

## Concept Used –Word Vectors, Sum and Mean of Word Vectors

- Train word2vec model for skill phrases obtained from Phrase Generator
- Word2vec model builds vectors of individual words
- Find mean vector for skill phrases
  1. Use word vector of each word in skill phrase from the vocabulary
  2. Add all word vectors
  3. Divide it by the number of word in the phrase
- Use the mean word vector to compare skill phrases to recommended skills

- Examples of word vectors:

**Construction Management** [-2.53128121e-03 -1.77700538e-04 -2.65417900e-03 -3.87543434e-04  
-1.15432322e-03 1.62815582e-03 1.00338110e-03 -1.52967067e-03.....-3.40627413e-03 -3.35135451e-03  
-2.22289260e-03 3.08321090e-04 -3.39231896e-03 -2.66436837e-03  
2.76020076e-03 4.11090790e-04 2.04549381e-03 9.39245452e-04  
-1.60909526e-03 -1.57659827e-03 -1.14465761e-03 -3.12960986e-03  
2.36541731e-03 -1.92798441e-03 -6.19727012e-04 -3.52221658e-03  
-1.10748247e-03 -3.67690041e-03 3.10475053e-03 -1.39216939e-03]

**internal company business systems** [ 2.3599188e-03 -1.2206250e-03 7.9751259e-04 7.0600610e-05  
-1.8862416e-03 -1.2368266e-04 -9.1530092e-06 8.7684812e-05.....  
1.4695667e-03 -1.0500317e-03 -8.4408821e-04  
-9.4342022e-04 5.0632062e-04 -3.1003775e-04 -1.4311107e-04  
-4.1986880e-04 -1.4218048e-04 4.7879951e-04 -1.1398435e-03]



# Similarity Check

Concept Used – Word Vectors, Sum and Mean of Word Vectors

Similarity Calculator:

Measure used -> Cosine Similarity

Python library used -> scipy import spatial

Finds the angle between the two vectors  
and calculates cosine value of that angle

Similarity score of 0.25 used as threshold

All the phrases greater than threshold value  
are selected

2 new features	Integration	1	2	0.35856548
20 strong knowledge software engineering discipline	Java (Programming Language)	3	0	0.273530692
20 strong knowledge software engineering discipline	C++ (Programming Language)	3	1	0.297596008
20 strong knowledge software engineering discipline	C (Programming Language)	3	2	0.279707402
20 workflow management tools Azkaban	Open AI	5	3	0.252860934
7 equivalent expertise	Aws Best Practices	6	0	0.254608423
23 new automation technologies	Systems Design	7	5	0.291011453
0 wireless analytics space IEEE	Systems Architecture	10	4	0.259406716
15 new processes	Data Integration	11	5	0.273268759
12 negotiate solutions Hardware/Software	Software Development Life Cycle	12	4	0.276912421
14 system-level computational architectures	Software Development Life Cycle	12	4	0.273831517
6 new projects	Integration	14	2	0.341670662
18 continuous improvement system stability performance	Architecture	14	1	0.335764796
1 multi-threaded software development experience Applic	Databases	15	1	0.273800403
6 continuous integration Software development experienc	Software Engineering	15	2	0.291767538
12 different problems	Software Engineering	20	5	0.282554746
15 reusable components	Kubernetes	20	0	0.252689928
11 related experience Master	Spring Framework	27	5	0.277940065
19 WORKING CONDITIONS Typical office environment	Hibernate (Java)	27	4	0.26278162
26 technical staff Establishes	Enterprise Java Beans	27	0	0.256945223
28 interactive Web User Interfaces	Hibernate (Java)	27	4	0.325454235
7 common Java Design	Java (Programming Language)	29	2	0.344111264
6 steady upward career progression art Software Engineer	Software Engineering	30	4	0.486126363
© abandoned engine reconnection art Software Engineer Engineering				



# Word2Vec Model

- Pretrained Model with GloVe - Global Vectors for Word Representation
- GloVe is an unsupervised learning algorithm for obtaining vector representations for words
- Gigaword 5th Edition corpora (6B tokens, 400K vocab) with 100 dimensions



# Clustering

- No labels for the matched activities- Unsupervised learning
- Used two clustering methods
  - K means clustering – Centroid model
  - Hierarchical clustering – Connectivity model
- Clustering performance evaluation (Based on cohesion and separation)
  - Silhouette Score (-1 to 1)  
S is calculated as :  $s = b - a / \max(a, b)$  where,
    - b - the mean distance between a sample and all other points in the next nearest cluster
    - a - the mean distance between a sample and all other points in the same cluster
  - Calinski-Harabaz Index (Higher the better)  
S is calculated as :  $s = b/a$  where
    - b - mean between cluster separation
    - a - mean within cluster separation



# K-Means clustering

- It is a partitioning algorithm - partitions dataset into as number of clusters to minimize intra-partition distances
- Starts with a first group of randomly selected centroids, which are used as the beginning points for every cluster
- Performs iterative (repetitive) calculations to optimize the positions of the centroids

**Library – `sklearn.cluster import Kmeans`**

## **Parameters:**

- K – Number of clusters ( 7-11)
- Max\_iterations ( 20 ,25)



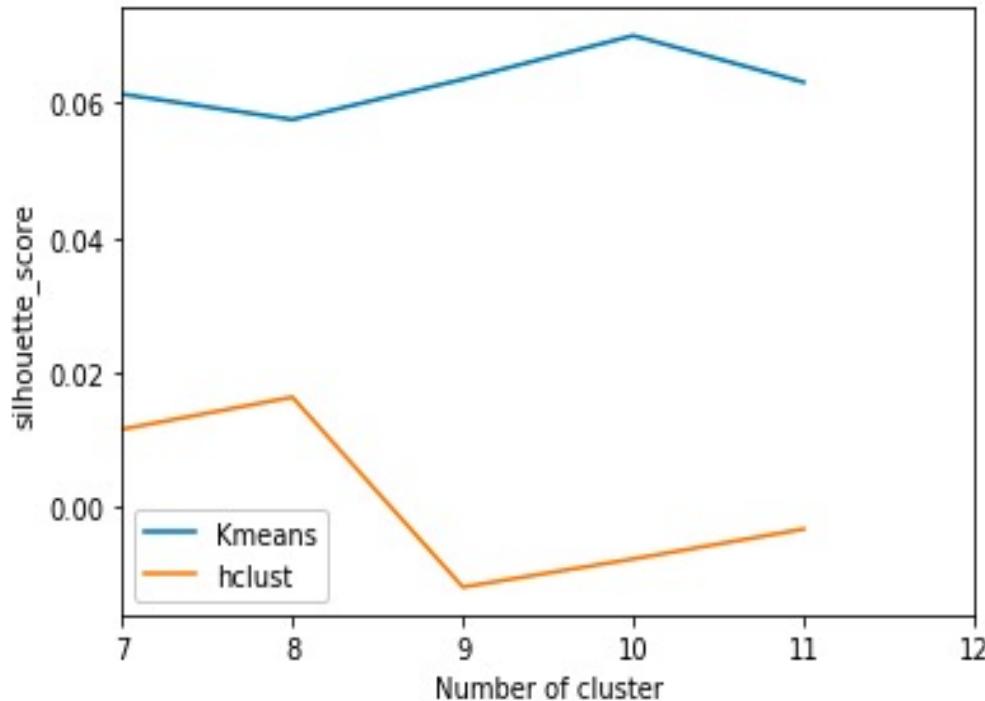
# Hierarchical Agglomerative Clustering

- In this technique, each data point is considered as an individual cluster.
- At each iteration, the similar clusters merge with other clusters until K clusters are formed.
- **Library - from sklearn.cluster import AgglomerativeClustering**

## Parameters:

- Linkage = ward
- Number of clusters k =range(7,11)

# Performance Analysis of Clustering Algorithms



For clusters: 10

The average silhouette\_score is : 0.06998661

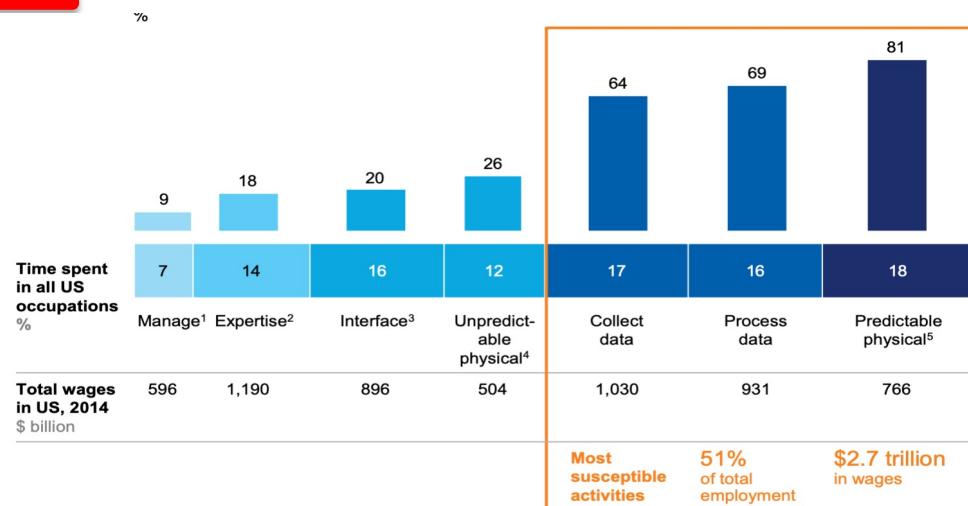
The value of Calinski\_value is : 233.82560805202831

# Dimensions of activities and impact of Automation

- According to studies, all the activities in any occupation can be categorized into the following segments:

- Managing and leading people
- Expertise to decision making, planning and creative tasks
- Interfacing with stakeholders
- Collecting and Processing data
- Physical work

**MOST IMPACTED BY AUTOMATION**



Source: US Bureau of Labor Statistics ; McKinsey Global Institute Analysis



# Cluster Analysis

## Physical work and Manual intervention

### Cluster 0:

[ 'able work', 'able lift', 'assigned work area', 'heavy work', 'assist help callers', 'regular basis help', 'preferred ability lift']

### Cluster 1:

[ 'on-site fabrication', 'hydraulic system', 'disassemble mechanical equipment', 'de-dusting equipment']

### Cluster 2:

[ 'professional customer service', Customer service experience Persuasive, 'extraordinary customer service experience customers']

### Cluster 6:

[ 'assist customer questions needs', ' assist customer questions needs', ' effective resolution issues', ' protect quality condition items']

### Outlier Cluster:

[ 'degree Information Systems', ' degree Marketing', ' degree Electrical Engineering']

## Expertise to decision making, planning and creative tasks

### Cluster 3:

[ 'service-oriented architecture', 'object oriented programming', 'third party auditors Judgment Decision Making', 'Complete video closeout audits']

### Cluster 4:

[ 'critical thinking problem-solving skills', ' deep consulting skills', ' technical leadership', ' Technical writing skills']

### Cluster 5:

[ 'big data analysis application Data integration', ' complex data models', ' in-depth knowledge products services clients']

### Cluster 7:

[ 'equivalent Docker Experience OpenShift Kubernetes', ' Solid Experience Azure DevOps Able create Jenkins pipelines', ' AWS Certified ']

### Cluster 8:

[ 'implement optimize Data Engineering ETL', ' stable Java applications', 'continuous integration server Fortran code', ' popular Big Data frameworks Hadoop']



# Conclusion

- Converting sentences to word vectors and using the word vectors for K-Means clustering gives the best results for this project
- Value of k=10 is the most optimal value of k as the cohesion and separation scores are maximized at k=10



# Future Work

- Phrase chunking using more patterns to improve the overall performance
- Check the clusters by using more clustering models(DBScan and BTM)
- Compare the activities based on the historical data of job postings and compare which all activities are automated
- Implementation of Doc2vec for activities and check cluster performance



**STEVENS**  
INSTITUTE OF TECHNOLOGY  
THE INNOVATION UNIVERSITY®

# THANK YOU!!

Questions???

