



Linked in

python

Octopaste



Results

Data Source
& Collection

Algorithm

Mission
&
Design

Further
Ideas

Visualization

ISGB 7978 Web Analytics

Group 7
Choudhury, Sumi
Gu, Jiahua
Huang, Feifei
Ji, Minxia



Mission

- Helping company to identify the import parts in Job Description to attract more candidates
- Prediction Model based on Skills & Qualification

Design

- Job Title
- Job Location
- # of Views
- # of Applicants
- Seniority Level
- Industry
- Job Description
- Skills & Qualification

Software Engineer Intern
6sense · San Francisco, CA, US
Posted 3 weeks ago · 253 views
2 alumni work here
Save Apply on company website

Job description
PLEASE NOTE: We are currently recruiting for Spring 2018 & Summer 2018 interns ONLY - Fall 2017 has been filled. Thank you for your interest! Role: Software Engineering Intern
6sense is a Predictive Intelligence Engine that is reimagining how B2B companies do sales and marketing. We work with big data at scale, advanced machine learning and predictive modeling to find buyers and predict what they will purchase and when. Our customers include fortune 500 companies like Cisco, Dell and Oracle.
Role Description
This is a role for budding software engineers actually pursuing a degree in

Seniority Level
Internship

Industry
Information Technology and Services, Computer Software, Internet

Employment Type
Other

Job Functions
Engineering, Information Technology

181 Applicants
45 Applicants in the past day
Seniority level
61% have a Master's Degree (similar to you)
21% have a Doctor of Philosophy
12% have a Bachelor's Degree
6% have other degrees
Location
See more

The image shows a LinkedIn job application interface with a Chrome DevTools overlay on the right. The LinkedIn page displays job statistics and applicant details.

LinkedIn Page Content:

- Applicants for this job:** 181 Applicants (highlighted with a red box), 45 Applicants in the past day. A reminder to update the profile is shown.
- Top skills:** You have 8 out of 10 top skills among all other applicants. Skills listed include Python, Java, R, C++, Microsoft Office, Microsoft Excel, Data Analysis, Machine Learning, and SQL.
- Seniority level:** 63 Entry level applicants, 14 Senior level applicants, 1 Director level applicant.
- Education:** 61% have a Master's Degree (Similar to you), 21% have a Doctor of Philosophy, 12% have a Bachelor's Degree, 6% have other degrees.
- Location:** Dallas/Fort Worth Area.

Chrome DevTools Overlay:

- Elements Panel:** Shows the DOM tree with a selected element: `div#ember6504job-view-layoutjobs-details.ember-view`.
- Styles Panel:** Displays the default style for the selected element: `element.style { }`.
- Console:** Shows a message: "Competitive intelligence about other applicants".
- Bottom Bar:** Displays "Highlights from the Chrome 60 update", including "New Audits panel, powered by Lighthouse" and "Third-party badges".



Linked in

python

Octopaste



Further
Ideas

ISGB 7978 Web Analytics

Group 7
Choudhury, Sumi
Gu, Jiahua
Huang, Feifei
Ji, Minxia

Data Source & Collection



Happy Data Hunting

www.octoparse.com



Octoparse Version 5.4.2

My Task: Home | Task: New Task | Task: web | Task: web

1 Set up Basic Information 2 Design Workflow 3 Extraction Options 4 Done

Task completed! Choose the next step.

Run Task: Local Extraction, Cloud Extraction, Create an API, Return to Task Status

Task Information Preview: Name: web, Availability: From 12/1/2017 To 5/31/2018, Description: Extraction Flow Chart: Go To Web Page, Click Cycle Pages, Click Loop Item, Click Item, Extract Data, Click to Pageinate

Extraction Options: Display error messages during the extraction process, Disable image loading (Speed up the extraction), Use Web-Proxy (HTTP) Proxy Settings

Data Scientist, Oculus VR, Menlo Park, CA, US

Job title, keywords, or company name Location Find jobs

Data Scientist, Oculus VR, Menlo Park, CA, US

Posted 2 days ago 48 views

Be one of the first 50 applicants.

Data Extracted: 0 row Total Time Spent: 25sec Speed: less than 1 row/minute

Job Name	Company Name	Company Location	Posted Date	Num Views	Descr	Seniority Level	Industry	Employ Type	Filed	Current
1										

Welcome remembervoice, You are Professional Account. The expiration date is 2017-12-05.

Octoparse Version 5.4.2

My Task: Home | Task : New Task | Task : web | Task : web | Task : url

1 Set up Basic Information 2 Design Workflow 3 Extraction Options 4 Done

Task completed! Choose the next step:

- Run Task
 - Local Extraction: Run the task on your computer
 - Cloud Extraction: Run the task on Octoparse Cloud Platform
 - Schedule Cloud Extraction Settings: Schedule to run the task on Octoparse Cloud Platform
- Create an API: Link Octoparse directly to your system
- Return to Task Status: Check all the waiting tasks in Task Status

url (25 data records have been extracted)

Opening https://www.linkedin.com/jobs/view/470689693?refId=362603398151210576127969&efl_flag=1_search_wip_jobs

Sign in | Join now

This job is no longer accepting applications.

Job title, keywords, or company name | Location | Find jobs

Data Scientist
 Infospace
 Irving, TX, US
 Posted 4 days ago

People also viewed

Data Scientist NEW
 Premier Group Re...
 Jersey City, New... | 138

Data Extracted: 25 rows | Total Time Spent: 5min 14sec | Speed: 2 rows/minute

Field1	Field2	Field3	Field4	Field5
1	187 views	Leisure, Tra...	Mid-Senior L...	As the sole ...
2		Marketing a...	Associate	America Po...
3		Marketing a...	Entry level	Beyondsoft...
4	145 views	Information...	Not Applicable	Gartner has...
5	400 views	ElectricalE...	Mid-Senior L...	OSRAM's C...
6	1143 views	Oil & Energ...	Mid-Senior L...	Senior Data...
7		Health, Wel...	Not Applicable	Job Summa...
8		Banking, Re...	Associate	Senior Data...
9		Information...	Associate	Requiremen...
10		Marketing a...	Entry level	Job Descrip...
11	122 views	Marketing a...	Mid-Senior L...	Navigant D...
12		Information...	Entry level	About Us/...

Task Information Preview

Name: url
 Availability: From 12/1/2017 To 6/1/2018
 Description:

Extraction Flow Chart:

```

graph TD
    Start(( )) --> Loop[18 Loop Item]
    Loop --> GoTo[Go To Web Page]
    GoTo --> Extract[Extract Data]
    Extract --> End(( ))
  
```

Extraction Options

- ☐ Display error messages during the extraction process
- ☐ Disable image loading (Speed up the extraction)
- ☐ Use Web Proxy (HTTP) | Proxy Settings

Export to Database
 Export to Excel(2007)
 Export to Excel(2003)
 Export to CSV
 Export to TXT
 Export to HTML

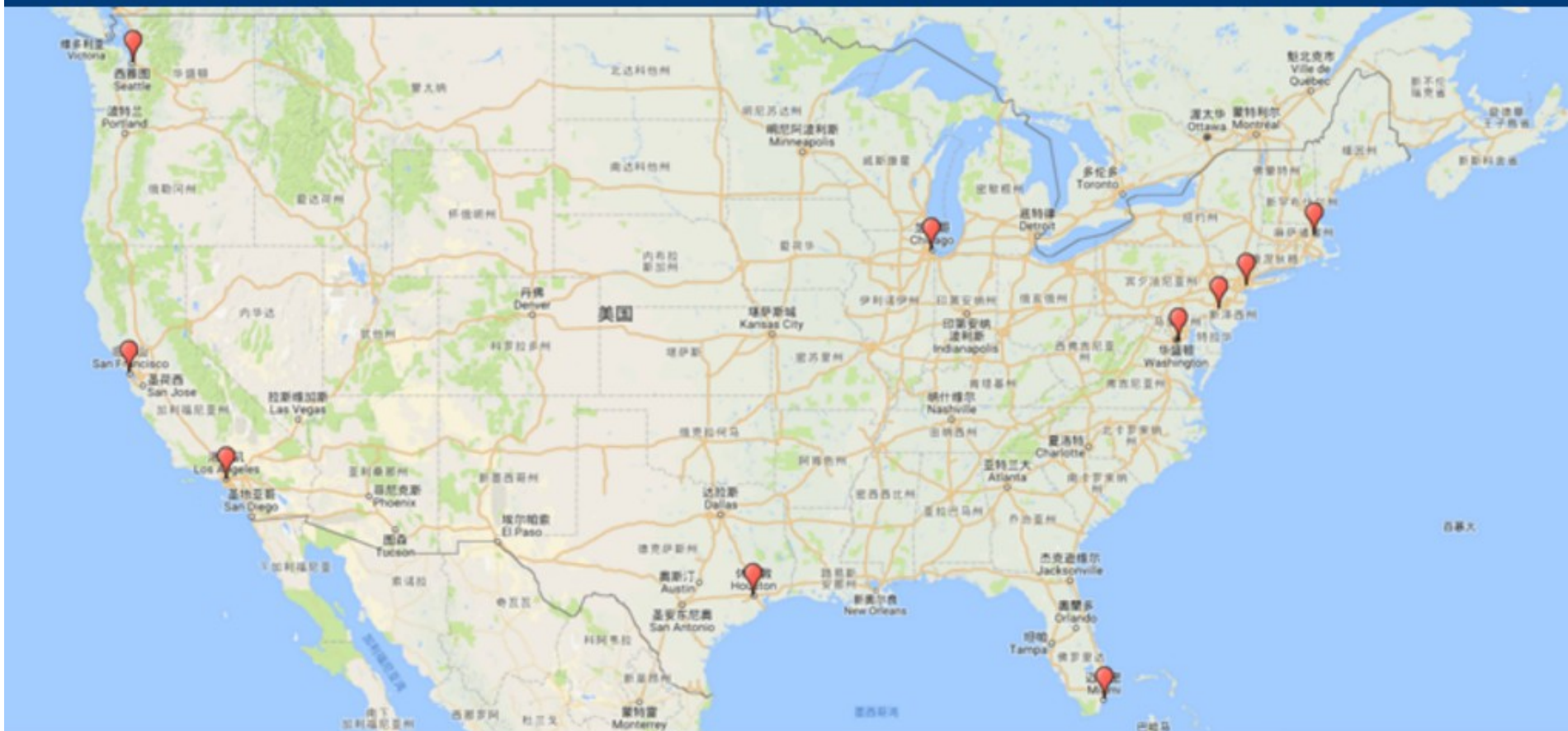
Quick Start

My Task

Task Status

Your Standard Account will expire in 1 days.

Renewals



		Renewals		Export to HTML											
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	
767	129	29	Not Applicable	Information Technology and Services	Lead analytics consultin	Consultant / Principal Consultant / Consulting Partner – Banking-Insurance / Retail / Telecom	Primary res								
768	1	362	Director	Management Consulting	Demonstrate technical b	As a Data Scientist, you will be: Providing essential insight and leverage on business problems through e									
769	383	174	Mid-Senior level	Investment Banking	Experience in building m	Machine Learning Software Engineer will be a member of a growing team with responsibilities for designi									
770	58	73	Not Applicable	Financial Services	2+ years working as a da	The Moody's Analytics Content team has a great opportunity for a talented, self-motivated and responsibl									
771	618	525	Executive	Computer Software	PhD or MS degree in Con	We are working with a multimedia content creator that is looking for a data scientist to join their research									
772	242	114	Not Applicable	Computer Software	If you have not written a	Work with internal client facing teams on broker/carrier analytics and reporting Perform data analysis b									
773	286	58	Mid-Senior level	Financial Services Information Services Information Tec	Master's degree or high	Senior Data Scientist - Machine Learning, Modeling, Ideation & Research	Equifax empowers businesses								
774	107	42	Entry level	Computer Software Information Technology and Services	Applying both technical i	Opentext - The Information CompanyAs the Information Company, our mission at OpenText is to create so									
775	319	78	Not Applicable	Construction Machinery	MS degree in Industrial E	Caterpillar Inc. is more than big, heavy equipment and much more than yellow iron. Our products have ev									
776	305	91	Entry level	Information Technology and Services	Advanced degree (PhD pr	OpportunitySessionM is looking to add a talented Data Scientist to support the continued growth of our D									
777	26	13	Entry level	Marketing and Advertising Information Technology and Serv	Mastery of either Python	We are looking to hire an entrepreneurial Data Scientist to join our team in Boston (Quincy). This person									
778	63	14	Not Applicable	Consumer Goods Food & Beverages Retail	Basic QualificationsEduc	This job contributes to Starbucks success by guiding business decisions through utilizing data analysis a									
779	131	39	Associate	Marketing and Advertising Information Technology and Services	You are awesome.You lov	Making projects happen: Coordinating and managing survey research projects on time and on budgets, we									
780	200	73	Mid-Senior level	Computer Software Internet Information Technology and Service	PhD in CS, Stats, Mathem	TalentReach is a boutique "upstart" search service obsessed with finding the most talented individuals in									
781	4	1	Associate	Real Estate	You enjoy exploring and	As a Principal Data Science Engineer, your job is to identify and solve challenging problems and productiz									
782	18	2	Director	Information Technology and Services Computer Software Mecha	Minimum 5-8 years of ex	The Go To Market Director will be responsible for developing both short and long-term strategies for how									
783	197	33	Not Applicable	Information Technology and Services Computer Software Compu	Bachelor's degree or equ	Work closely with engineers throughout the development cycle to reproduce the small and large customer									
784	82	33	Entry level	Medical Devices Medical Practice Hospital & Health Care	Bachelors, MS or PhD in	As a Data Scientist you will be working on consulting side of our business. You will be responsible for ana									
785	171	58	Entry level	Biotechnology Hospital & Health Care Pharmaceuticals	Level of education and ex	Be responsible for setting standards for core technical platforms and continual optimization and applica									
786	3	0	Mid-Senior level	Internet	3 to 5 years experience b	We are looking for an experienced Sr. Data Scientist whos passionate about machine learning, model buil									
787	70	7	Entry level	Electrical/Electronic Manufacturing Construction Automotive	Bachelor's degree in Con	We are looking for an experienced Data Engineer with an uncanny ability to integrate multiple heterogene									
788	63	13	Entry level	Marketing and Advertising Internet Retail	2+ years of industry expe	Build complex statistical models that predict behavior and improve over time with additional data and lea									
789	163	34	Mid-Senior level	Computer Software	At least 3 Years Of Experi	Passion for innovation & energetic company.Strong commitment & investment in our people's developmer									
790	136	45	Entry level	Information Technology and Services Computer Software Compu	Experience in creating ar	Working in a very dynamic and exciting environment with top engineering and product talent the data sci									
791	289	38	Not Applicable	Computer Software Information Technology and Services Intern	Master's degree in Math,	ABP is seeking an outstanding Data Scientist/Statistician who will help to build the next generation of crec									
792	8	1	Associate	Information Technology and Services Staffing and Recruiting Fir	PhD in Machine Learning	We are looking for an Applied Scientist to join the Devices Demand Planning team for the entire device fan									
793	351	80	Director	Internet	PhD in CS or related field	Provide technical leadership, identify key business challenges and opportunities, and develop end-to-end									
794	3	0	Associate	Computer Software Information Technology and Services	Building statistical mode	Work as a team to define and execute data science solutions that address client use cases and business r									



Linked in

python

Octopaste

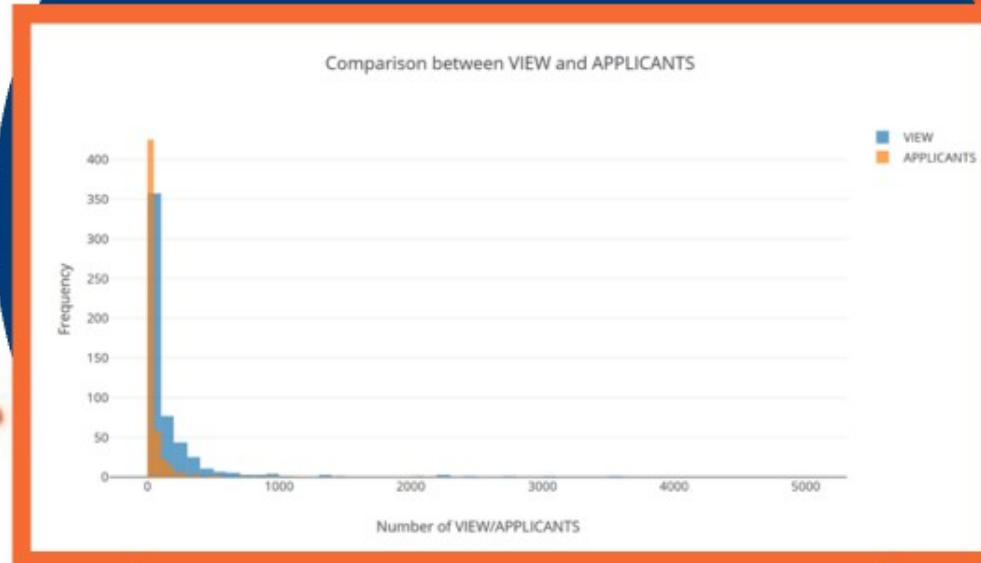


Further
Ideas

ISGB 7978 Web Analytics

Group 7
Choudhury, Sumi
Gu, Jiahua
Huang, Feifei
Ji, Minxia

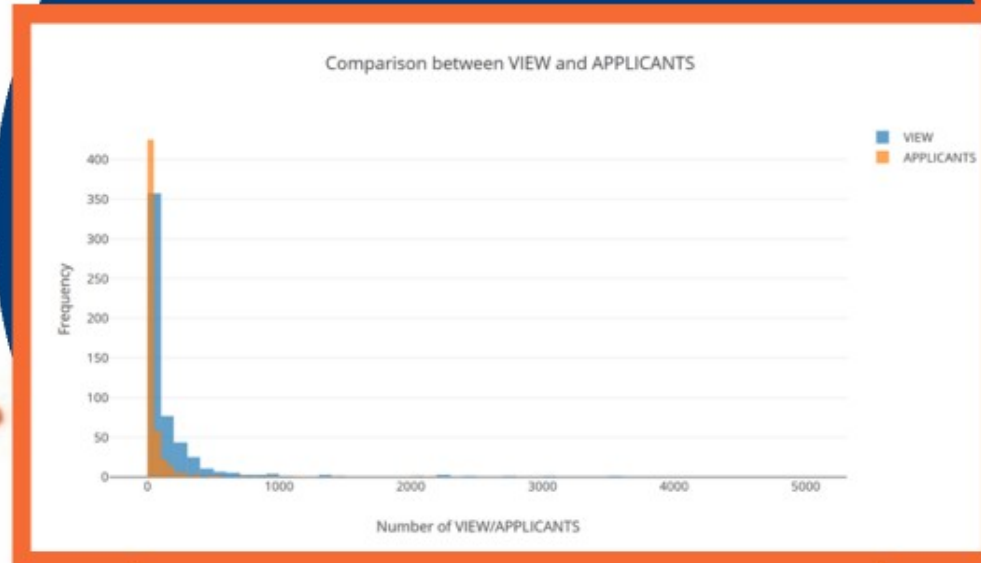
Visualization



Industry
Seniority Level

VS.
#Applicant
#Views

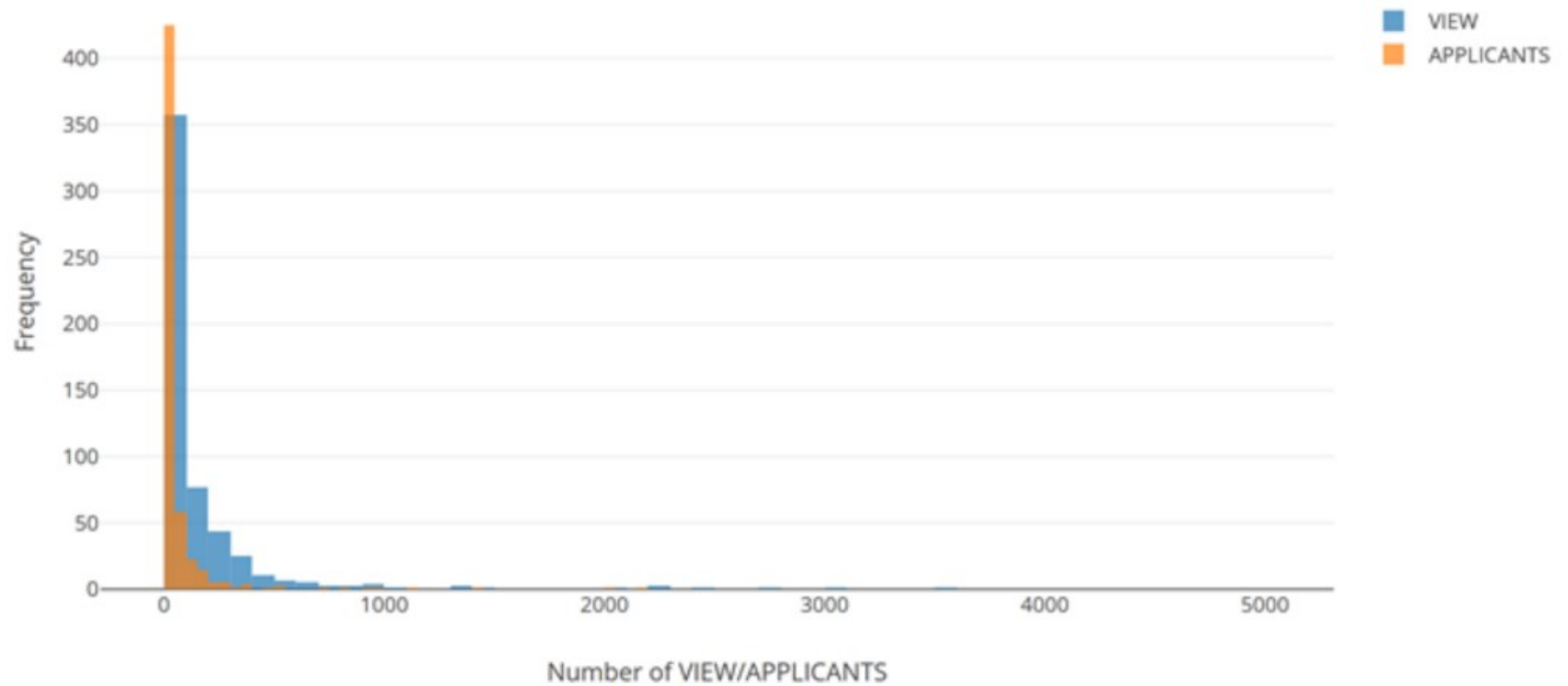
Visualization

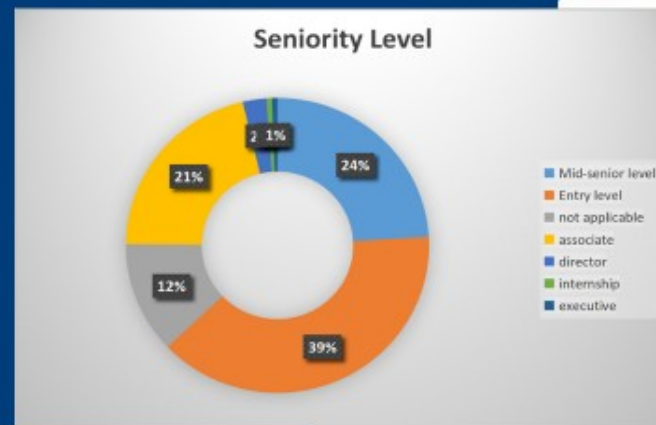
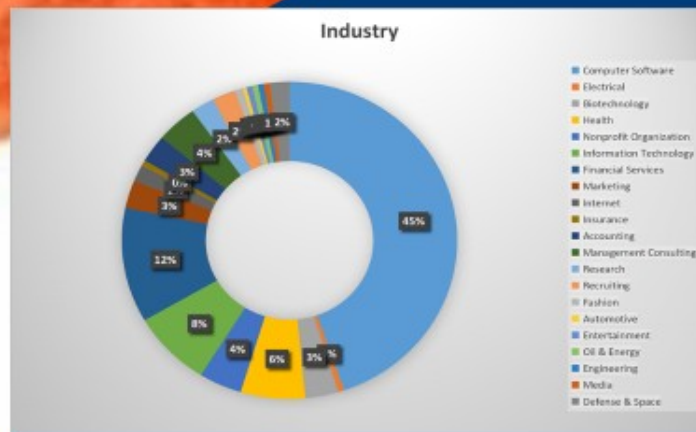


Industry
Seniority Level

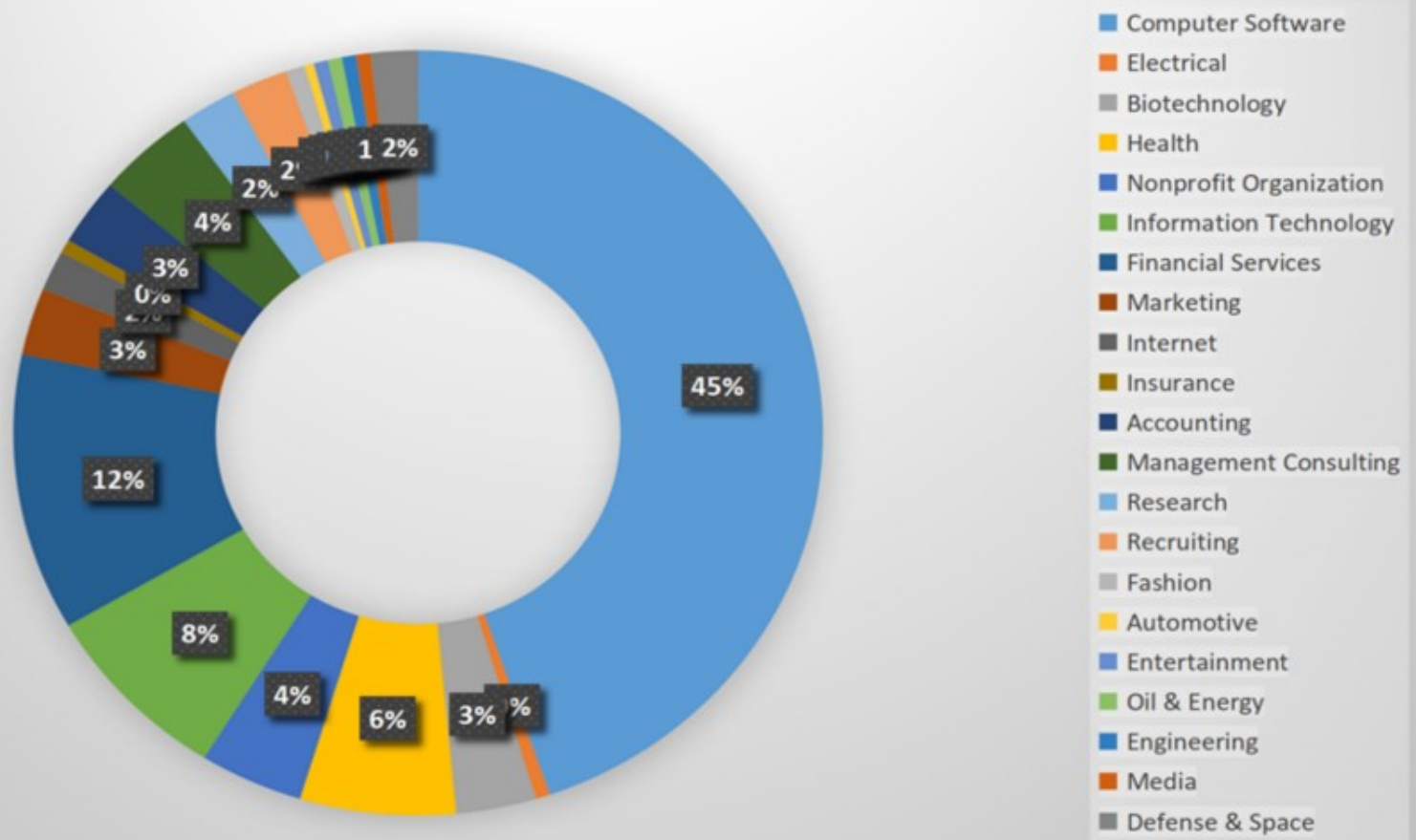
VS.
#Applicant
#Views

Comparison between VIEW and APPLICANTS

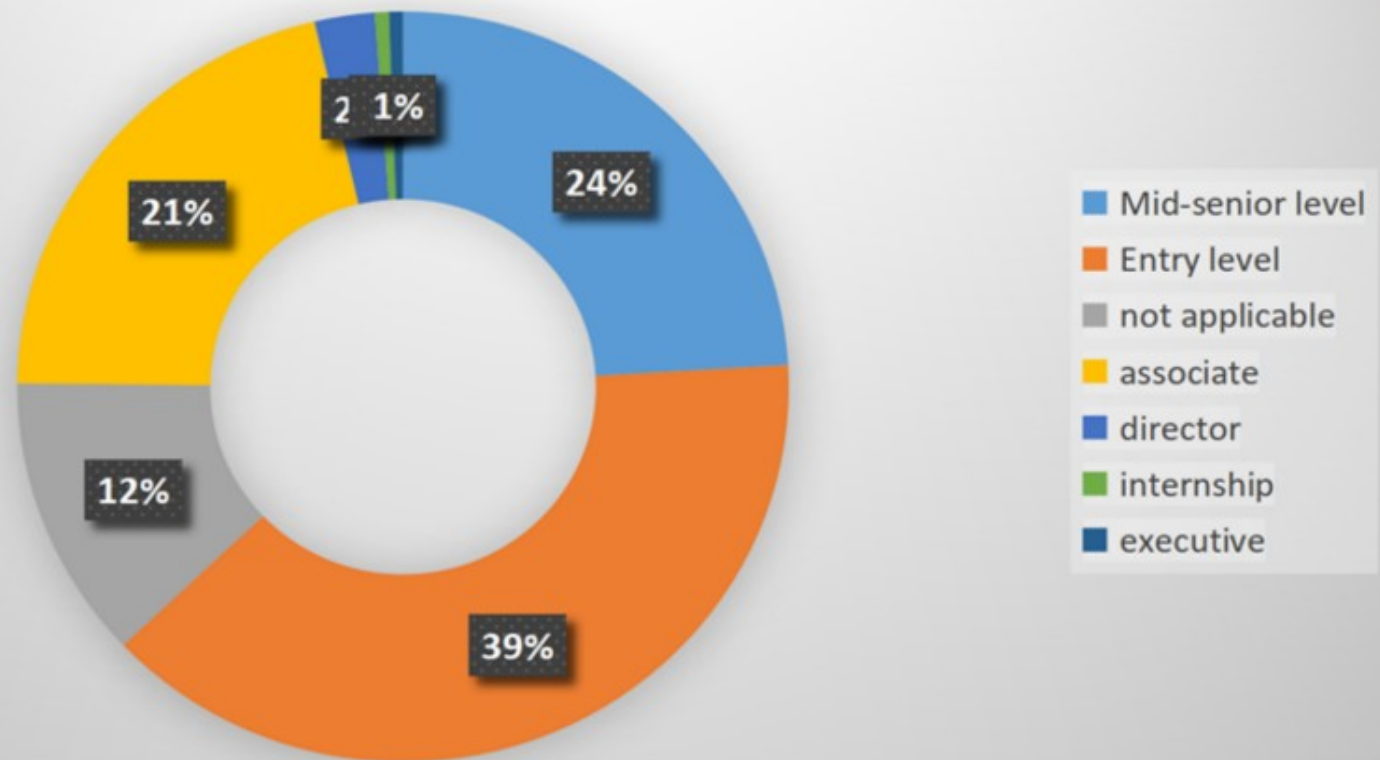


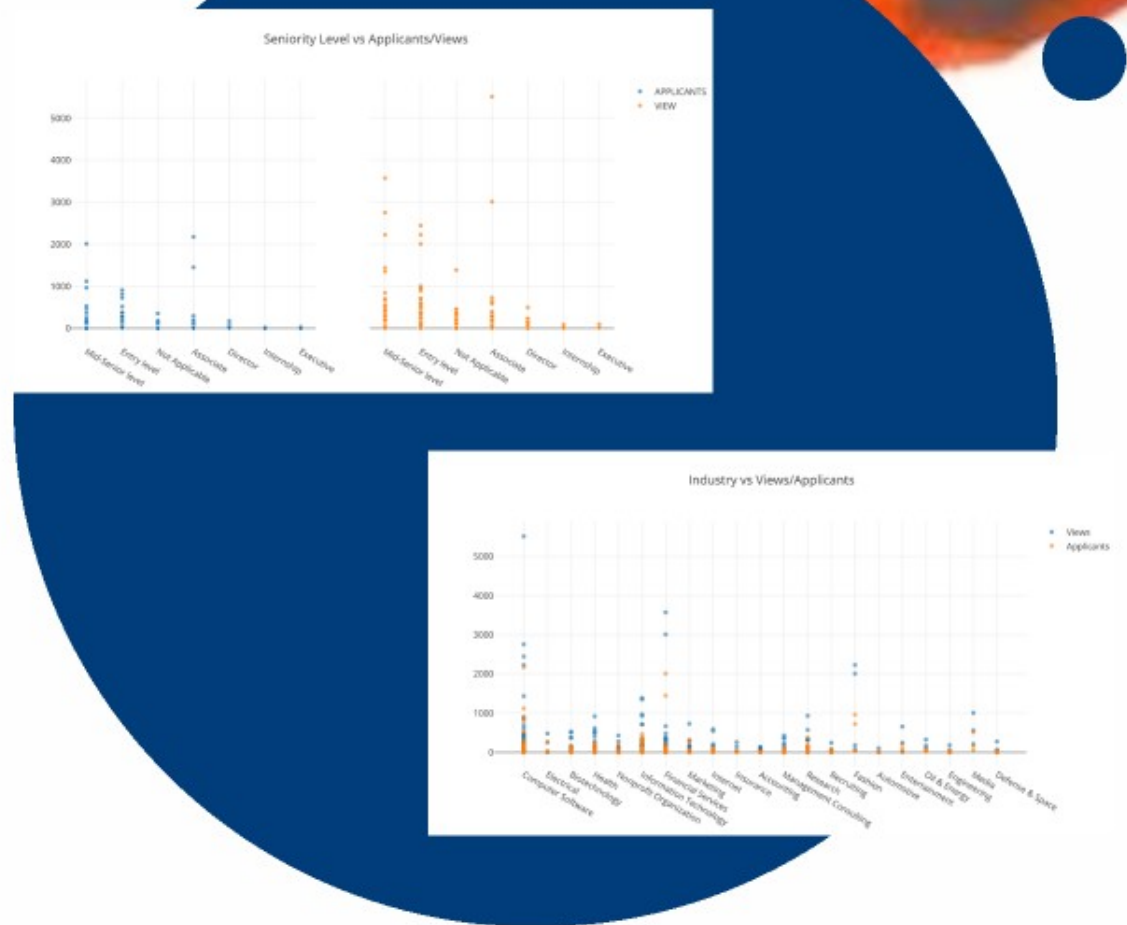


Industry

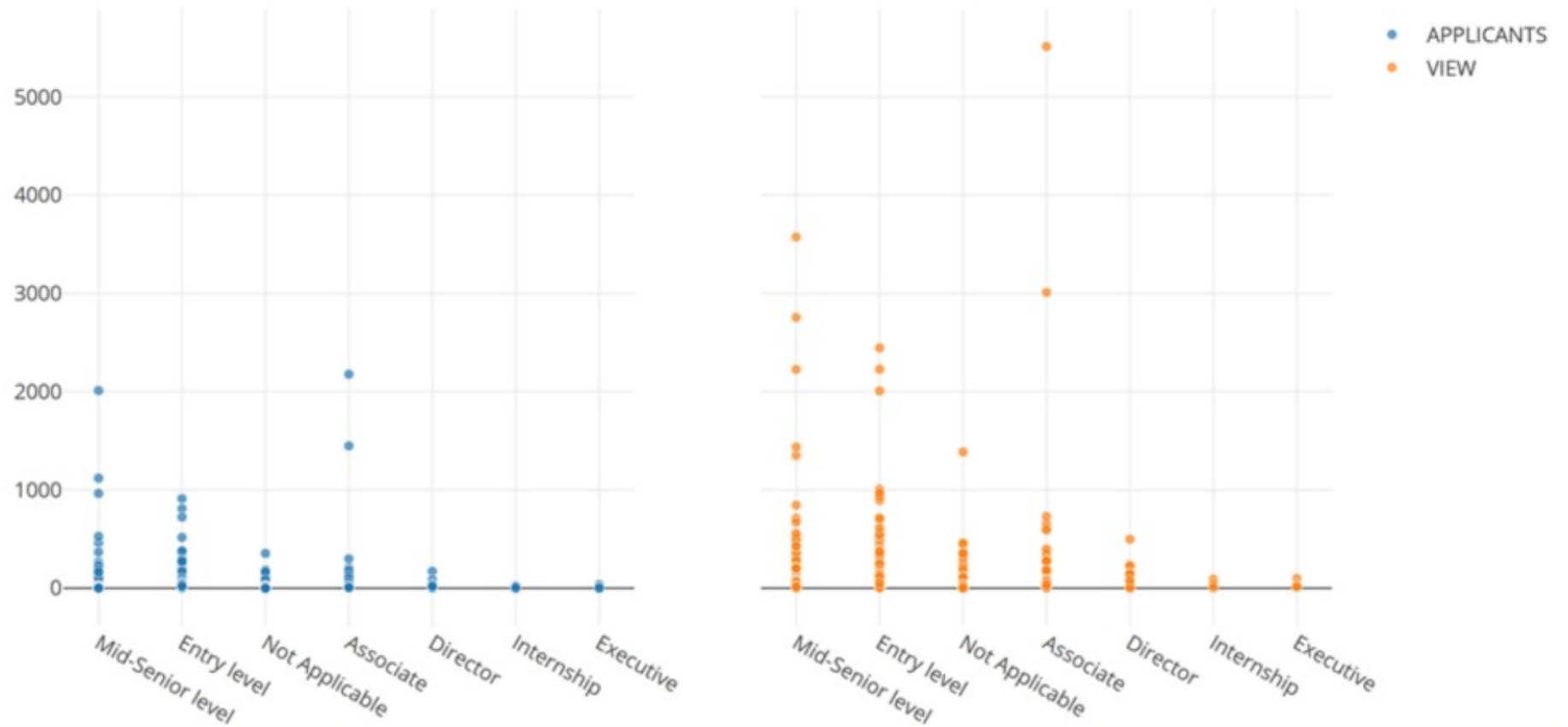


Seniority Level

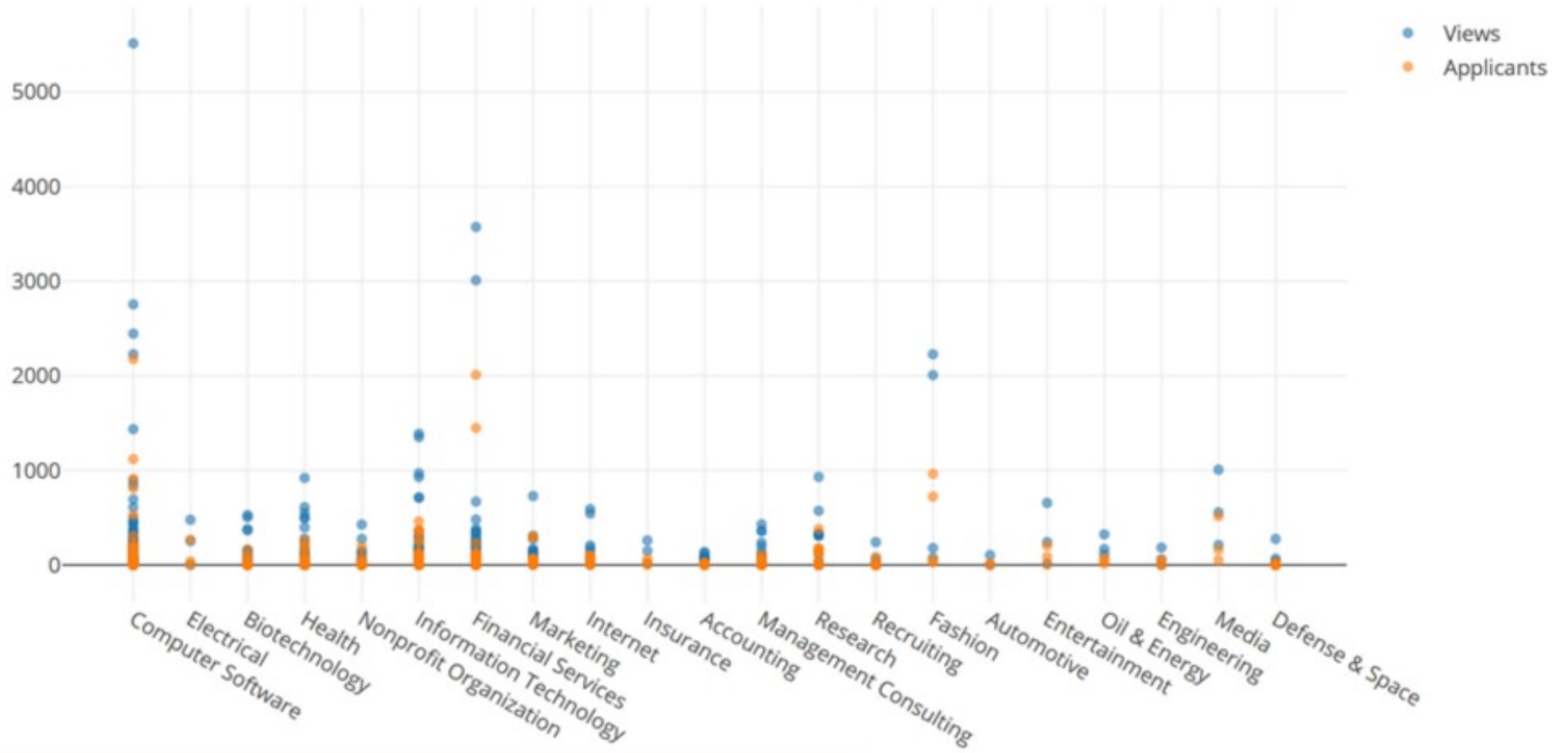


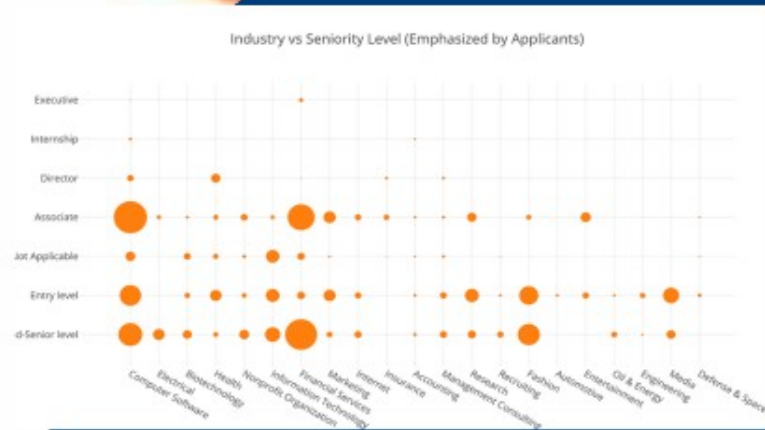


Seniority Level vs Applicants/Views

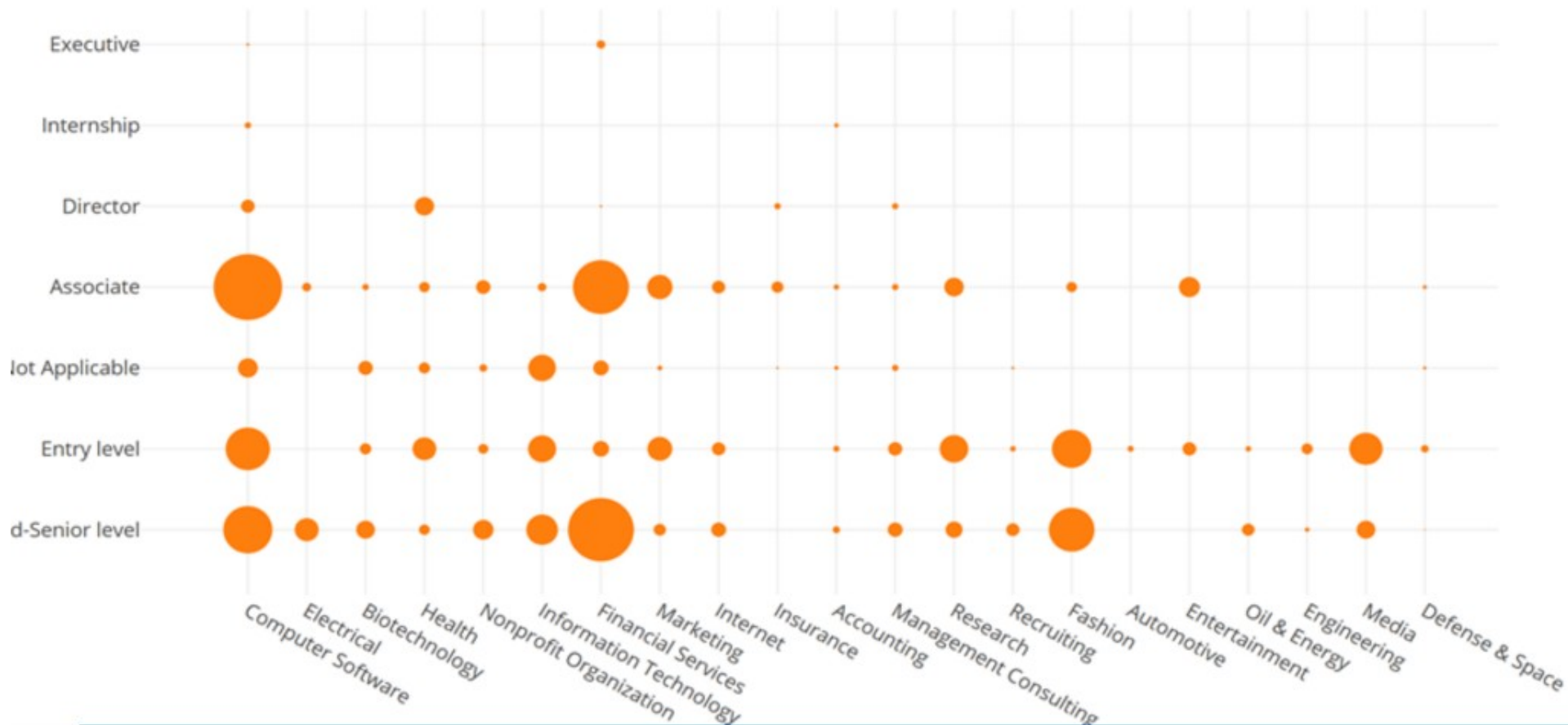


Industry vs Views/Applicants

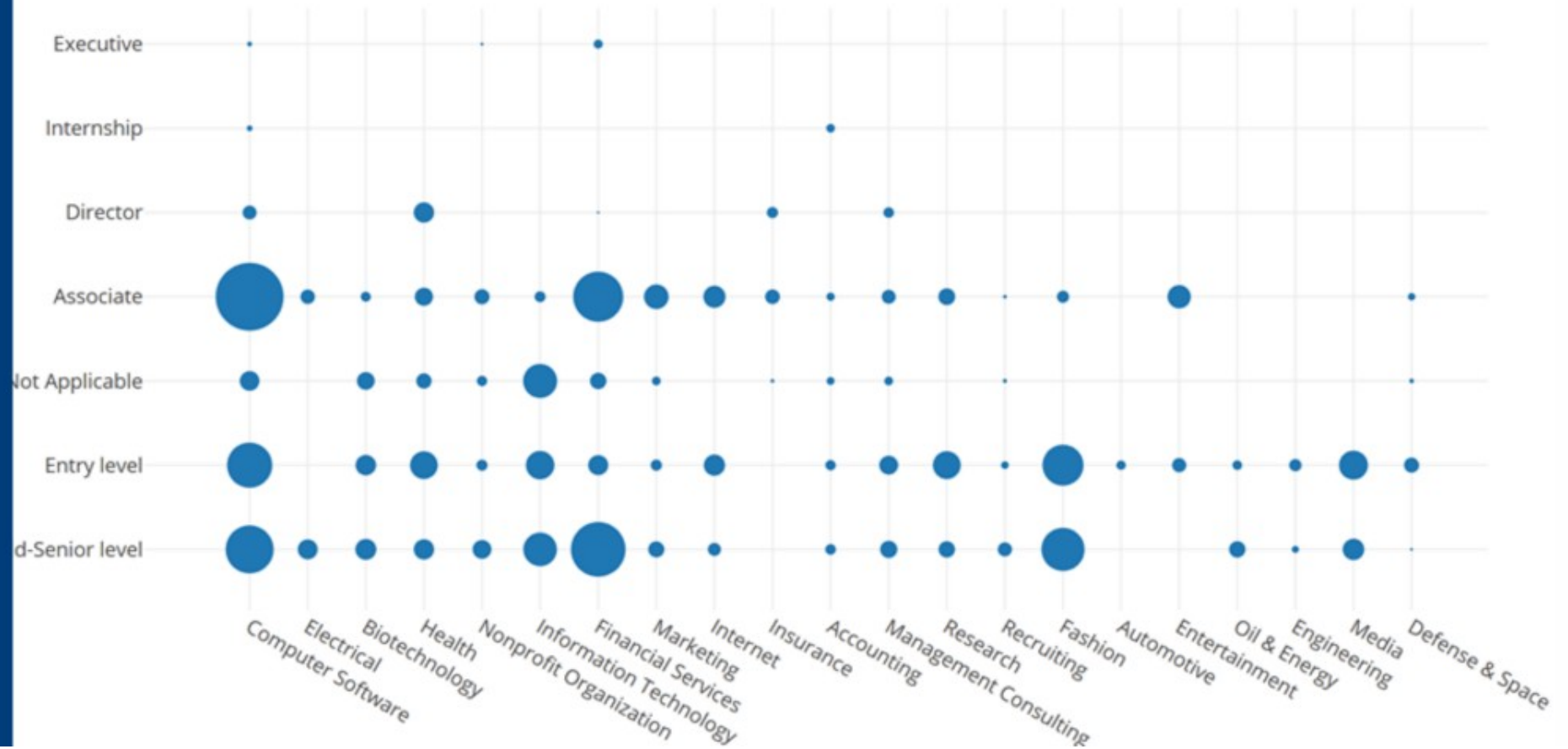




Industry vs Seniority Level (Emphasized by Applicants)



Industry vs Seniority Level (Emphasized by Views)





Linked in

python

Octopaste



Results

Data Source
& Collection

Algorithm

Mission
&
Design

Further
Ideas

Visualization

ISGB 7978 Web Analytics

Group 7
Choudhury, Sumi
Gu, Jiahua
Huang, Feifei
Ji, Minxia

Algorithm

- **Feature: Skills & Qualifications**
- **Label: Conversion rate**

Conversion Rate =
Number of Applicants/Number of Views

The median of conversion rate is **0.2**, which means every **5** views has **1** applicants

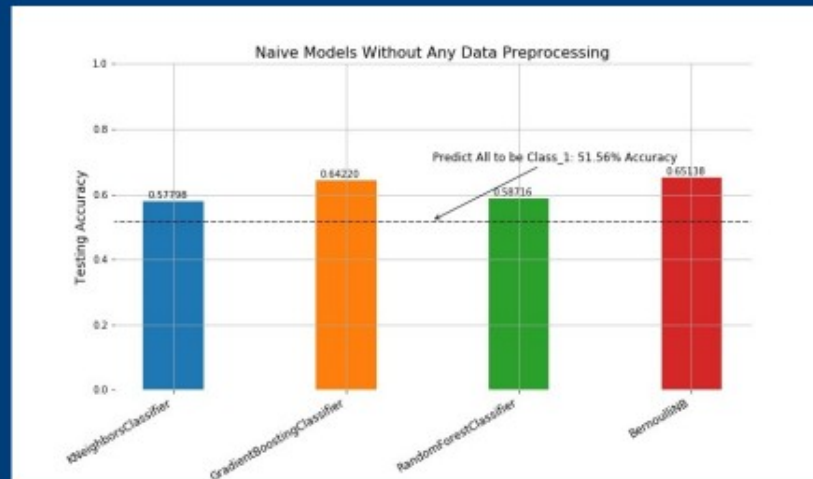
- label '1': conversion rate > 0.2 (good)
- label '0': conversion rate ≤ 0.2 (limited)

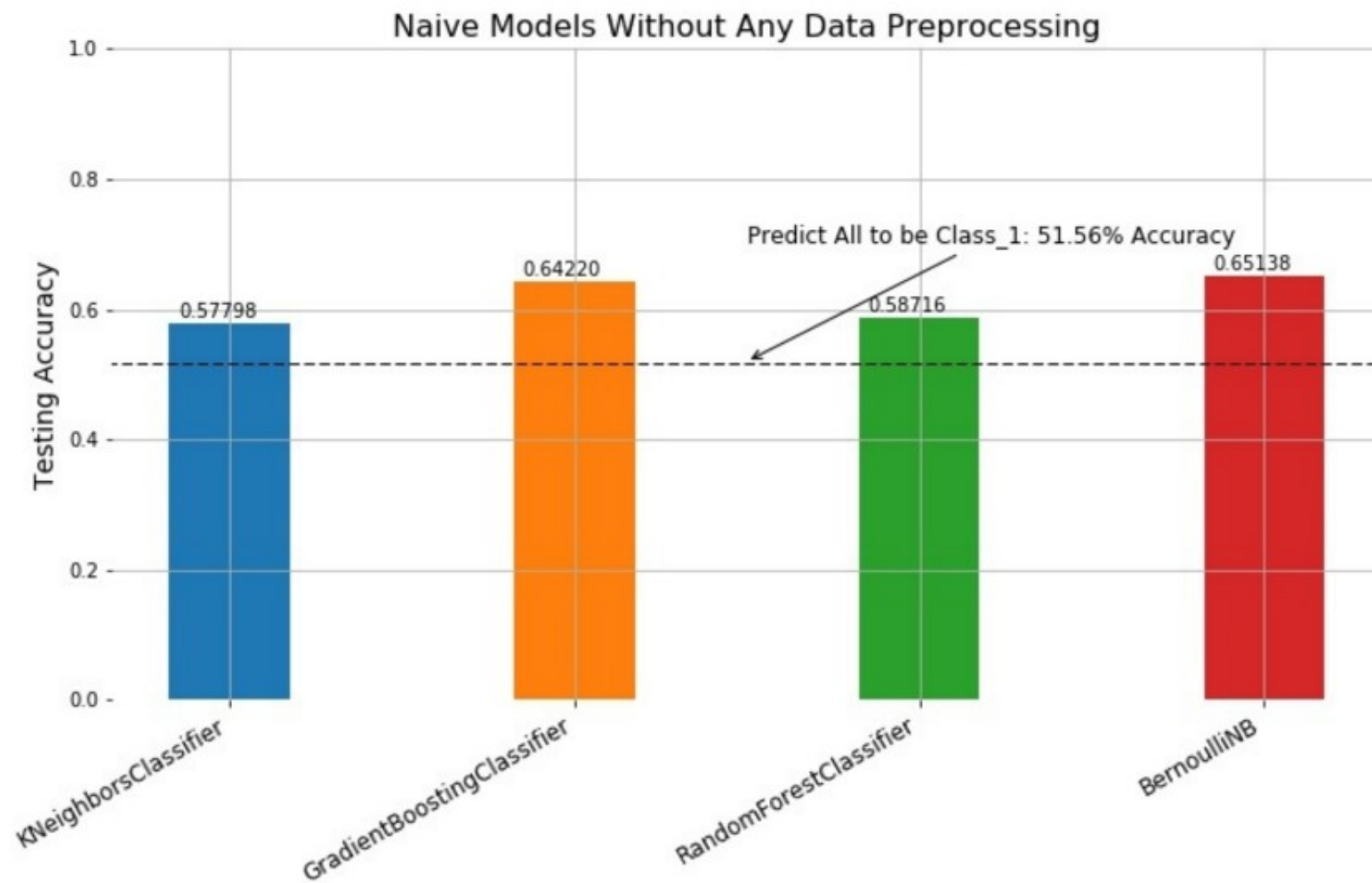
Naive
Model

Data Cleaning
&
Transforming

Reduce
Dimension
&
Classification

Naive Model





Data Cleaning & Transforming

- Removing Stopwords and Punctuations : remove commonly used words in the English language such as I, me, my, etc.
- Text Stemming: reduces words to their root form. For example, "moving", "moved" and "movement" to the root word, "move".
- Transformed each feature to TF*IDF

Reduce Dimension & Classification

Reduce Dimension:

- ICA: Independent Component Analysis
- PCA: Principal Component Analysis
- KernelPCA: Kernel Principal Component Analysis

Classification:

- KNeighbors Classifier
- GradientBoosting Classifier
- RandomForest Classifier
- BernoulliNaiveBayes Classifier

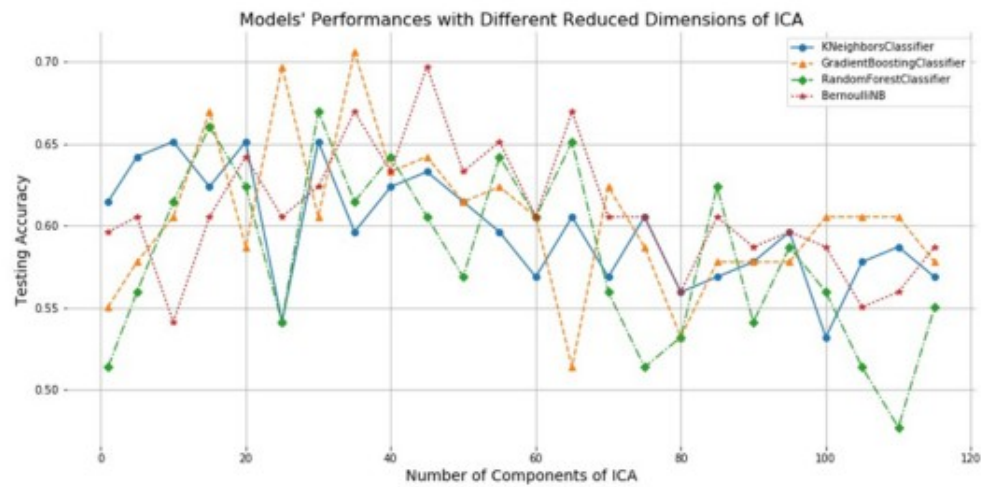

```

def train_after_DR(learners,dr,X_train_tfidf,X_test_tfidf,y_train,y_test):
    """
    function to trian machine learning algorithms after dimension reduction method
    'learners': list of machine learning algorithms
    'dr': dimension reduction method
    'X_train_tfidf,X_test_tfidf,y_train,y_test': training data and testing data
    """

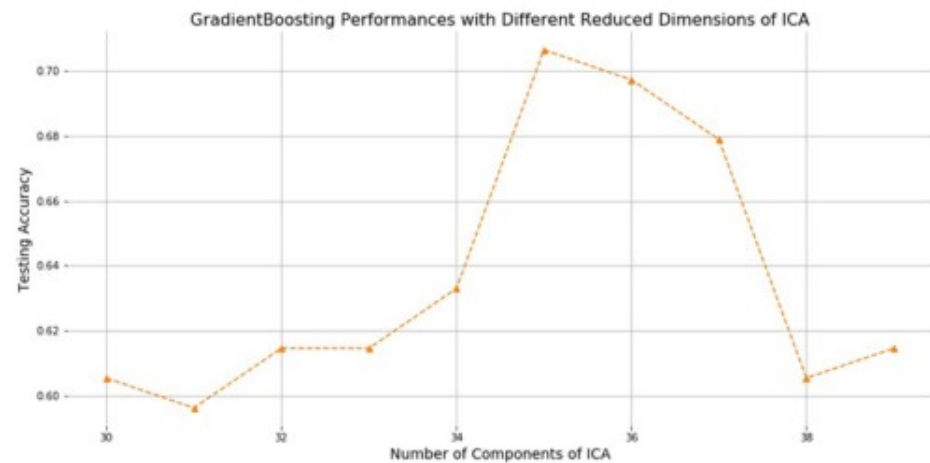
    result = {}
    # convert sparse matrix to ndarray
    X_train = X_train_tfidf.toarray()
    X_test = X_test_tfidf.toarray()
    # fit Dimension Reduction algorithm
    reduced_X_train = dr.fit_transform(X_train)
    reduced_X_test = dr.transform(X_test)
    # fit machine learning algorithms
    for learner in learners:
        learner.fit(reduced_X_train,y_train)
        pred = learner.predict(reduced_X_test)
        acc = np.mean(pred == y_test)
        result[learner.__class__.__name__] = acc

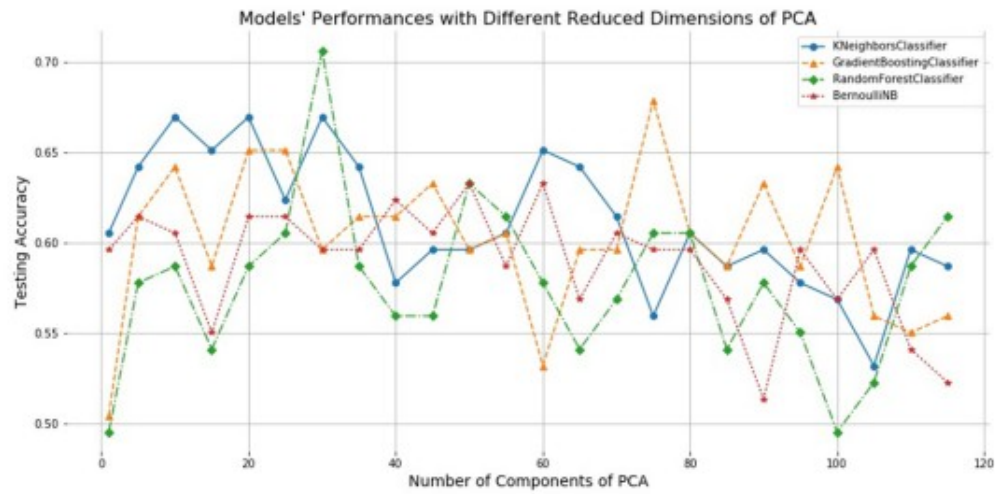
    return result

```

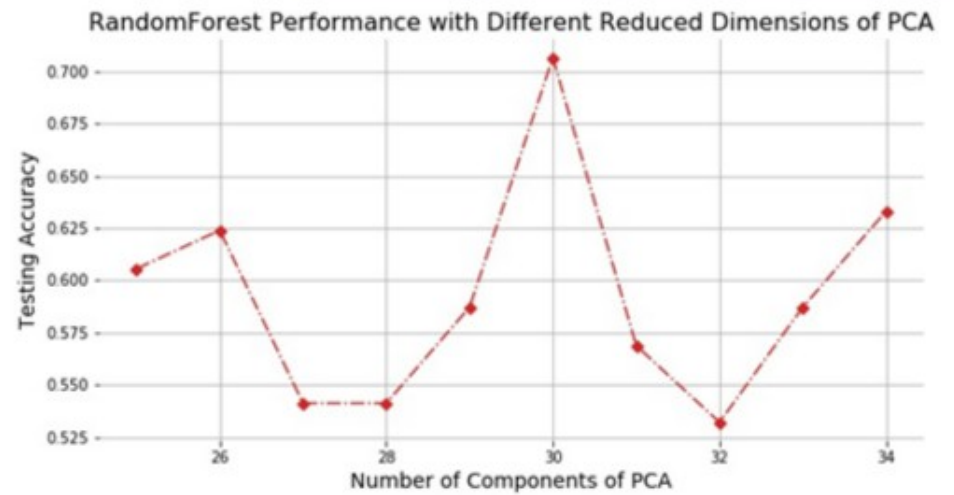


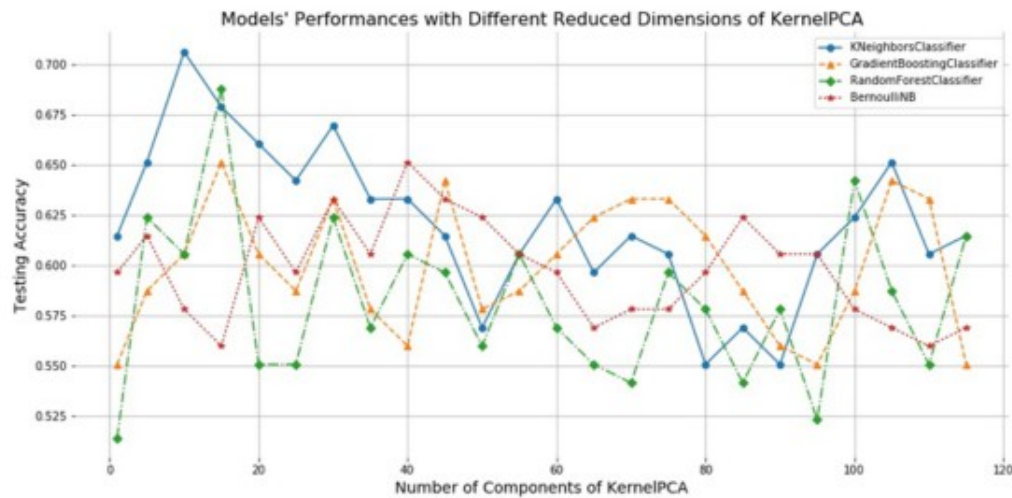
Gradient Boosting Classifier has the best performance based on first 35 ICs testing accuracy around 71%





Random Forest Classifier has the best performance based on first 30 PCs: testing accuracy around 71%



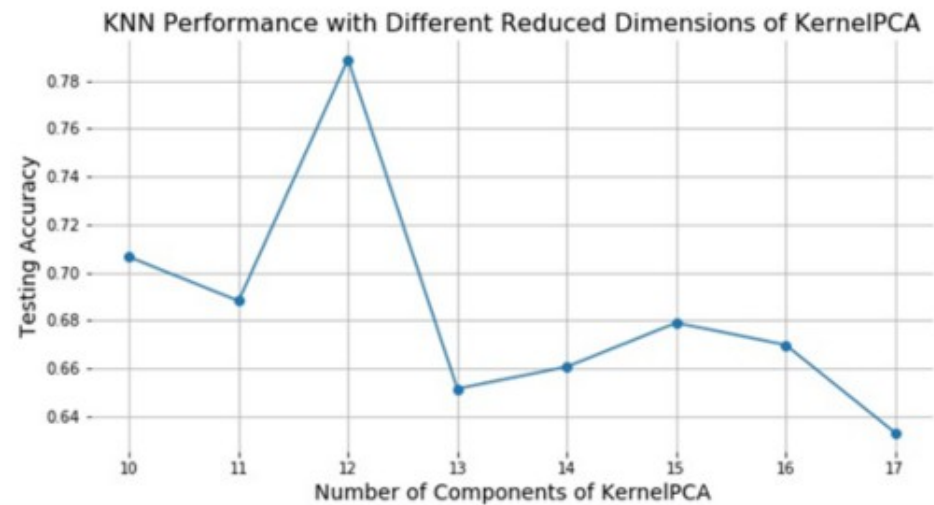


Best! 79%

Data: first 12 PCs with Kernel PCA

Model: KNeighbors Classifier

**KNeighborsClassifier has the best performance based on first 12 PCs:
testing accuracy around 79%**




```
{
  'PC0': [u'experi', u'learn', u'machin', u'statist', u'data'],
  'PC1': [u'data', u'abil', u'year', u'model', u'busi'],
  'PC10': [u'knowledg', u'least', u'minimum', u'advanc', u'year'],
  'PC2': [u'year', u'learn', u'degre', u'analyt', u'experi'],
  'PC3': [u'data', u'year', u'least', u'experi', u'strong'],
  'PC4': [u'data', u'experi', u'viewpoint', u'plu', u'mathemat'],
  'PC5': [u'least', u'busi', u'abil', u'statist', u'quantit'],
  'PC6': [u'requir', u'knowledg', u'minimum', u'busi', u'model'],
  'PC7': [u'data', u'plu', u'least', u'big', u'languag'],
  'PC8': [u'languag', u'common', u'data', u'experi', u'program'],
  'PC9': [u'experi', u'softwar', u'requir', u'least', u'analyt']}

```

	precision	recall	f1-score	support
Unpopular_Post	0.84	0.70	0.76	53
Popular_Post	0.75	0.88	0.81	56
avg / total	0.80	0.79	0.79	109



Linked in

python

Octopaste



Results

Data Source
& Collection

Algorithm

Mission
&
Design

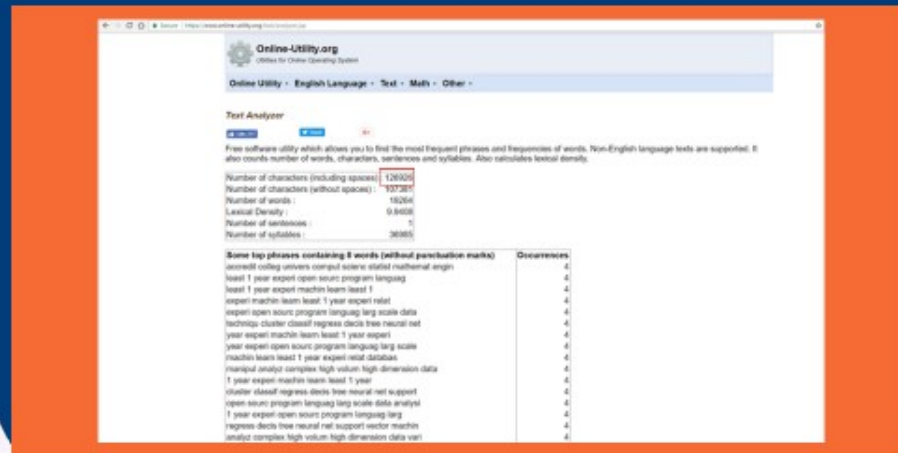
Further
Ideas

Visualization

ISGB 7978 Web Analytics


Group 7
Choudhury, Sumi
Gu, Jiahua
Huang, Feifei
Ji, Minxia

Results



The screenshot shows the Online-Utility.org Text Analyzer interface. It displays various statistics for a text sample, including character counts, word counts, and lexical density. A table at the bottom lists the top phrases containing 8 words, along with their occurrences.

Text Analyzer	
Free software utility which allows you to find the most frequent phrases and frequencies of words. Non-English language texts are supported. It also counts number of words, characters, sentences and syllables. Also calculates lexical density.	
Number of characters (including spaces) :	126909
Number of characters (without spaces) :	107281
Number of words :	18264
Lexical Density :	0.8668
Number of sentences :	1
Number of syllables :	36885
Some top phrases containing 8 words (without punctuation marks)	
	Occurrences
seawell collag univers comput scienc statid mathemat engin	4
least 1 year expert open sourc program languag	4
least 1 year expert machin learn least 1	4
expert machin learn least 1 year expert relat	4
expert open sourc program languag larg scale data	4
techniq cluster classif regress decisi tree neural net	4
year expert machin learn least 1 year expert	4
year expert open sourc program languag larg scale	4
machin learn least 1 year expert relat databas	4
manipul analysi complex high volum high-dimension data	4
1 year expert machin learn least 1 year	4
cluster classif regress decisi tree neural net support	4
open sourc program languag larg scale data analysi	4
1 year expert open sourc program languag larg	4
regress decisi tree neural net support vector machin	4
analysi complex high volum high-dimension data vari	4



Online-Utility.org
Utilities for Online Operating System

[Online Utility](#) - [English Language](#) - [Text](#) - [Math](#) - [Other](#) -

Text Analyzer

Like 217

Twitter

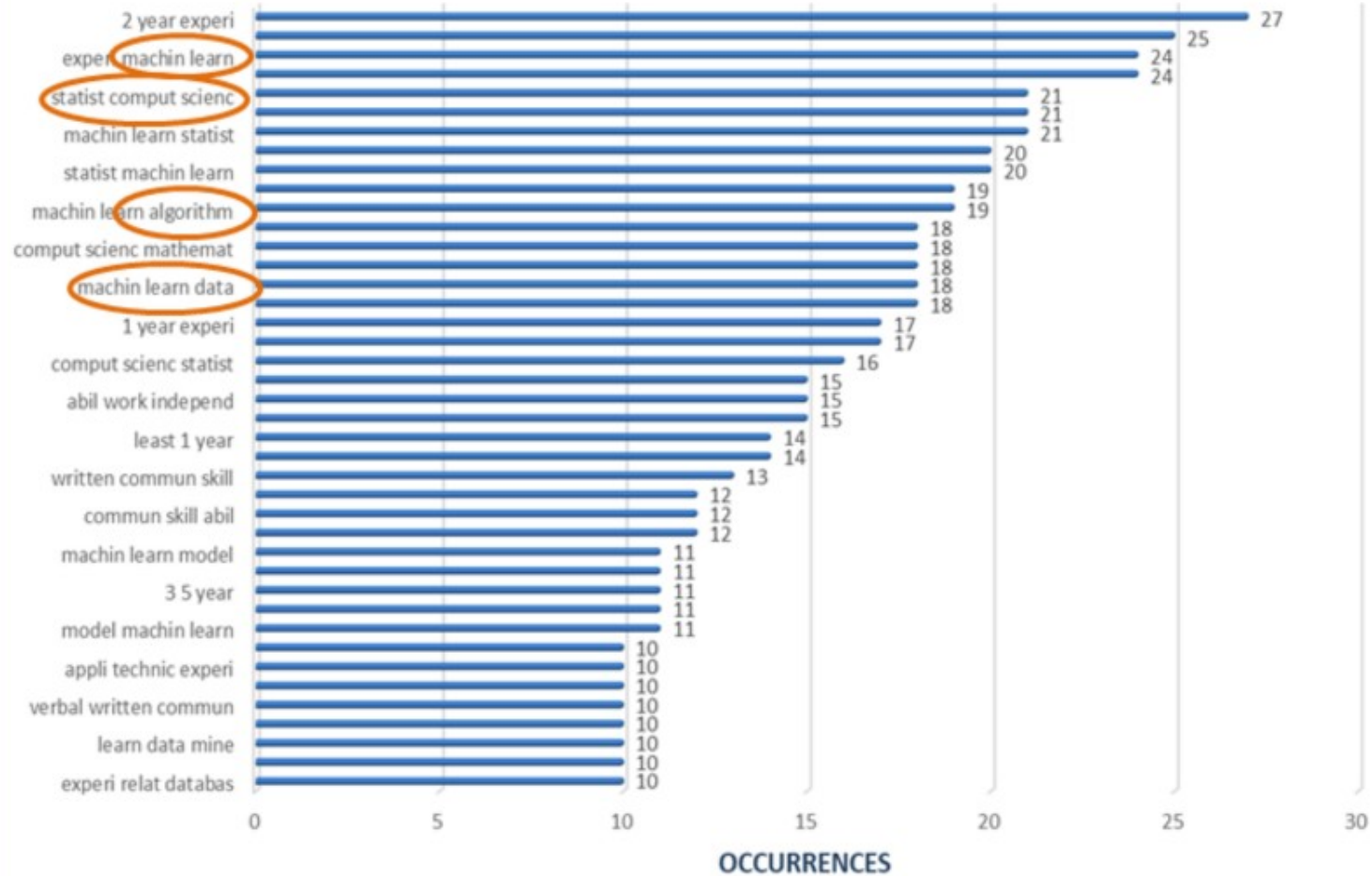
Google+

Free software utility which allows you to find the most frequent phrases and frequencies of words. Non-English language texts are supported. It also counts number of words, characters, sentences and syllables. Also calculates lexical density.

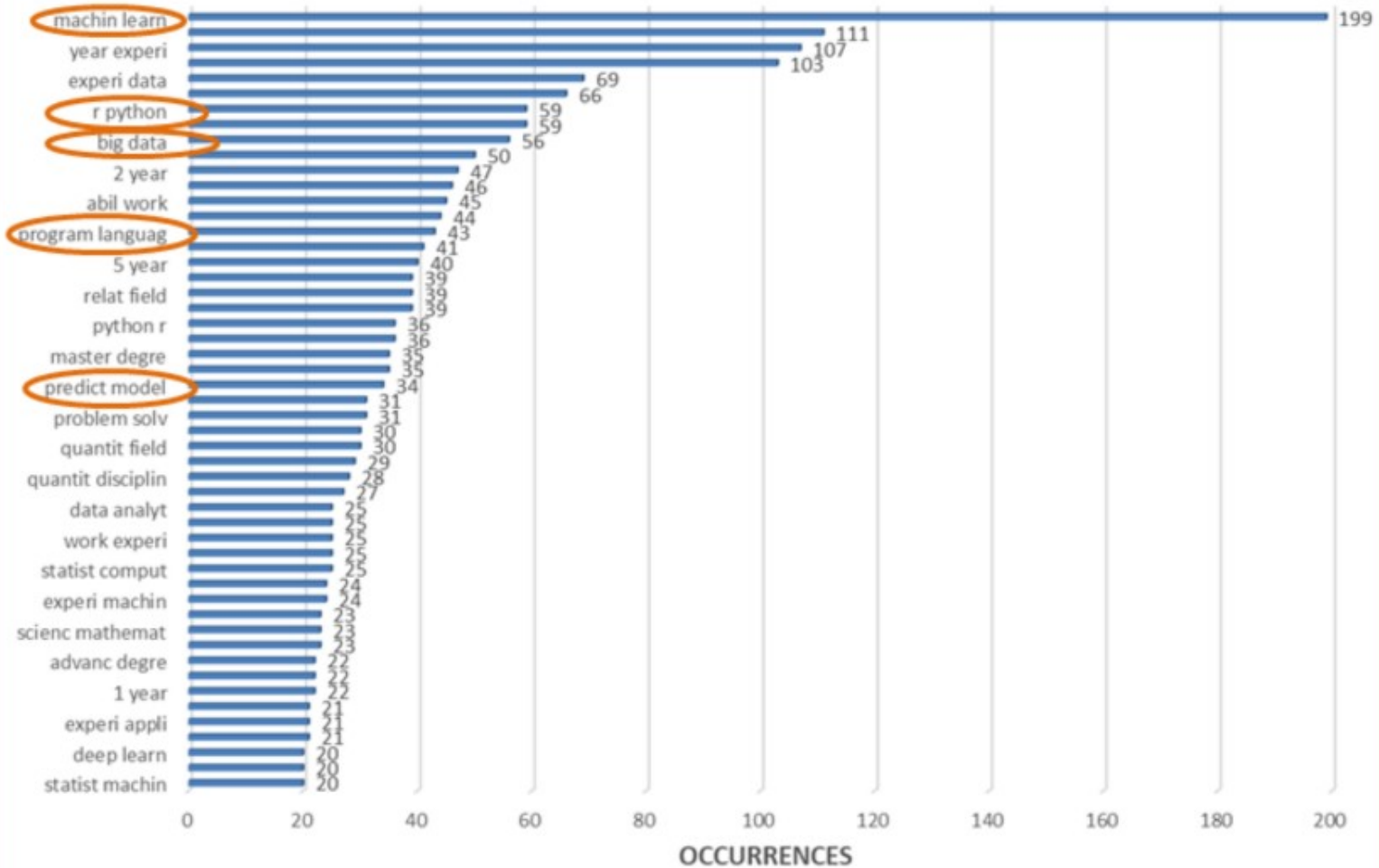
Number of characters (including spaces) :	126926
Number of characters (without spaces) :	107381
Number of words :	19264
Lexical Density :	9.9408
Number of sentences :	1
Number of syllables :	36985

Some top phrases containing 8 words (without punctuation marks)	Occurrences
accredit colleg univers comput scienc statist mathemat engin	4
least 1 year experi open sourc program languag	4
least 1 year experi machin learn least 1	4
experi machin learn least 1 year experi relat	4
experi open sourc program languag larg scale data	4
techniqu cluster classif regress decis tree neural net	4
year experi machin learn least 1 year experi	4
year experi open sourc program languag larg scale	4
machin learn least 1 year experi relat databas	4
manipul analyz complex high volum high dimension data	4
1 year experi machin learn least 1 year	4
cluster classif regress decis tree neural net support	4
open sourc program languag larg scale data analysi	4
1 year experi open sourc program languag larg	4
regress decis tree neural net support vector machin	4
analyz complex high volum high dimension data vari	4

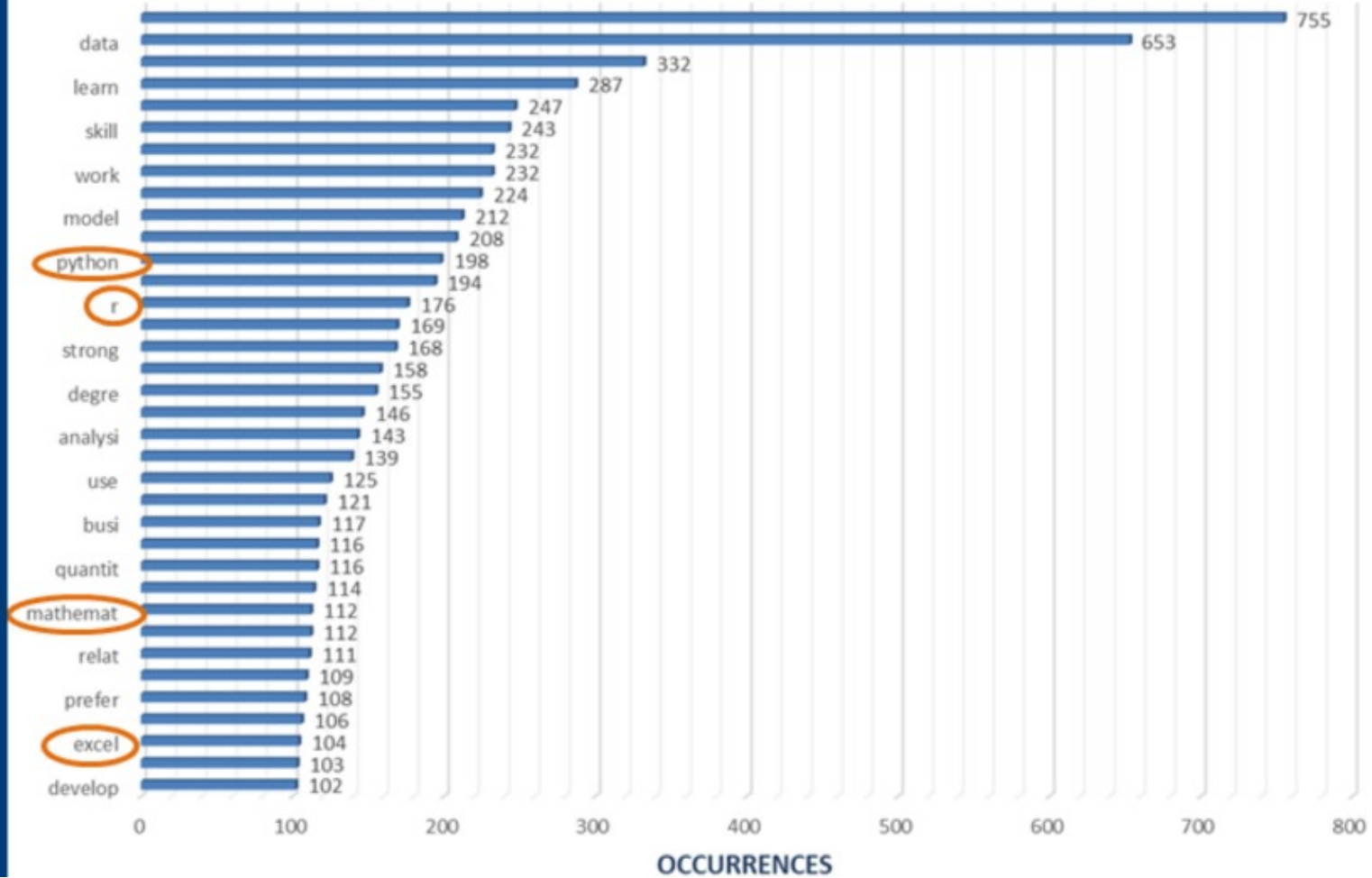
Top Phrases Containing 3 Words



Top Phrases Containing 2 Words



Unfiltered Word Count





Online-Utility.org

Utilities for Online Operating System

Online Utility - English Language - Text - Math - Other -

Text Analyzer

Like 217

Tweet

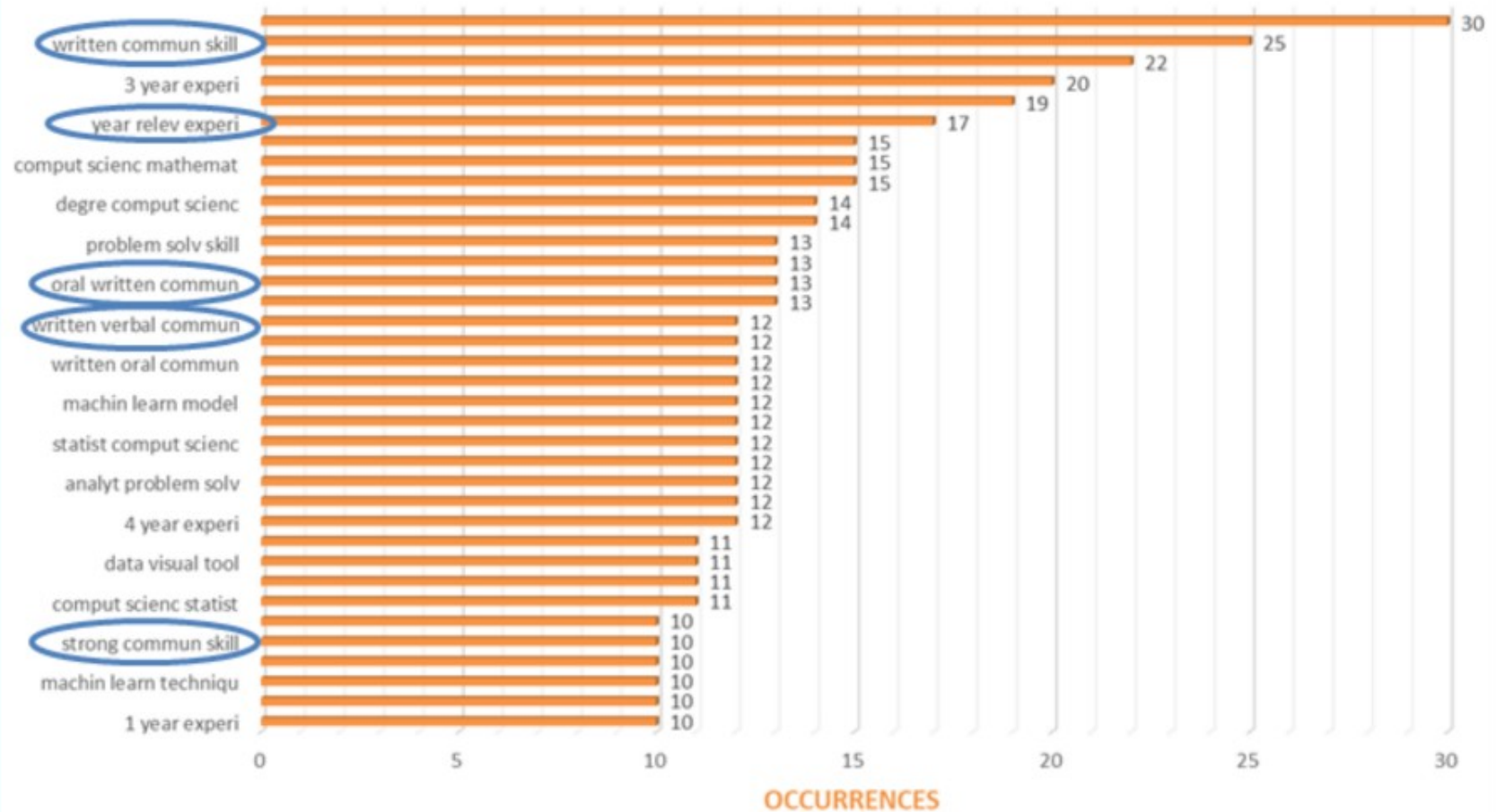
+

Free software utility which allows you to find the most frequent phrases and frequencies of words. Non-English language texts are supported. It also counts number of words, characters, sentences and syllables. Also calculates lexical density.

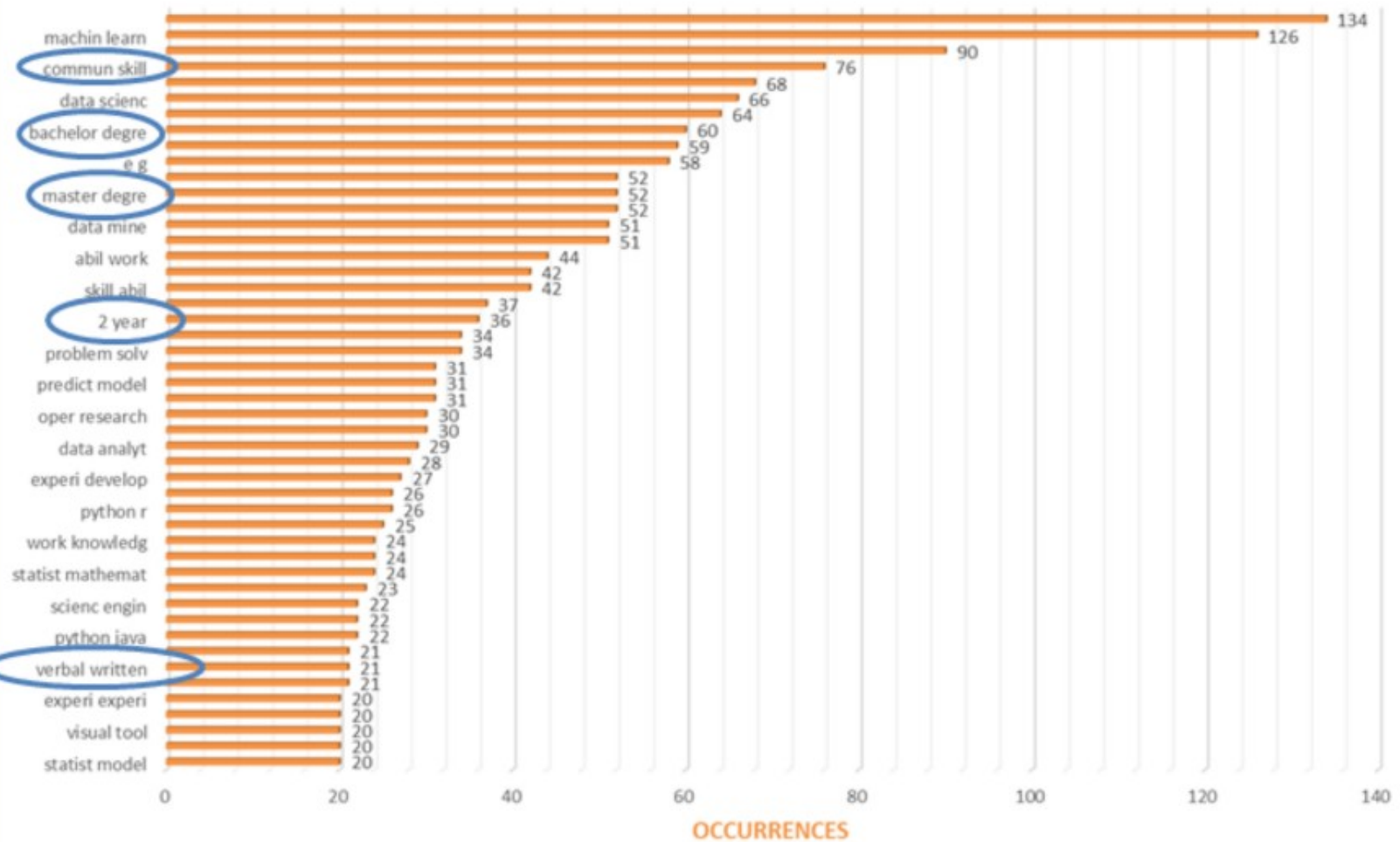
Number of characters (including spaces)	139158
Number of characters (without spaces) :	117856
Number of words :	21037
Lexical Density :	9.7067
Number of sentences :	1
Number of syllables :	41018

Some top phrases containing 8 words (without punctuation marks)	Occurrences
accredit colleg univers comput scienc statist mathemat engin	3
excel skill strong project manag skill abil priorit	3
knowledg data analyt busi intellig visual tool e	3
analysi skill abil clearli present complex inform simpl	3
stakehold experi data mine structur unstructur data sap	3
discern busi mean recommend action stakehold experi data	3
data mine data analysi skill abil clearli present	3
mine structur unstructur data sap hana sql etl	3
abil work matrix organ virtual team minim travel	3
commun skill strong knowledg data analyt busi intellig	3
sql etl data warehous etc busi environ larg	3
scienc master degre plu 5 year experi market	3
cluster regress analysi statist tool prefer r sa	3
busi mean recommend action stakehold experi data mine	3
degre plu 5 year experi market analyt excel	3
knowledg statist analys e g cluster regress analysi	3

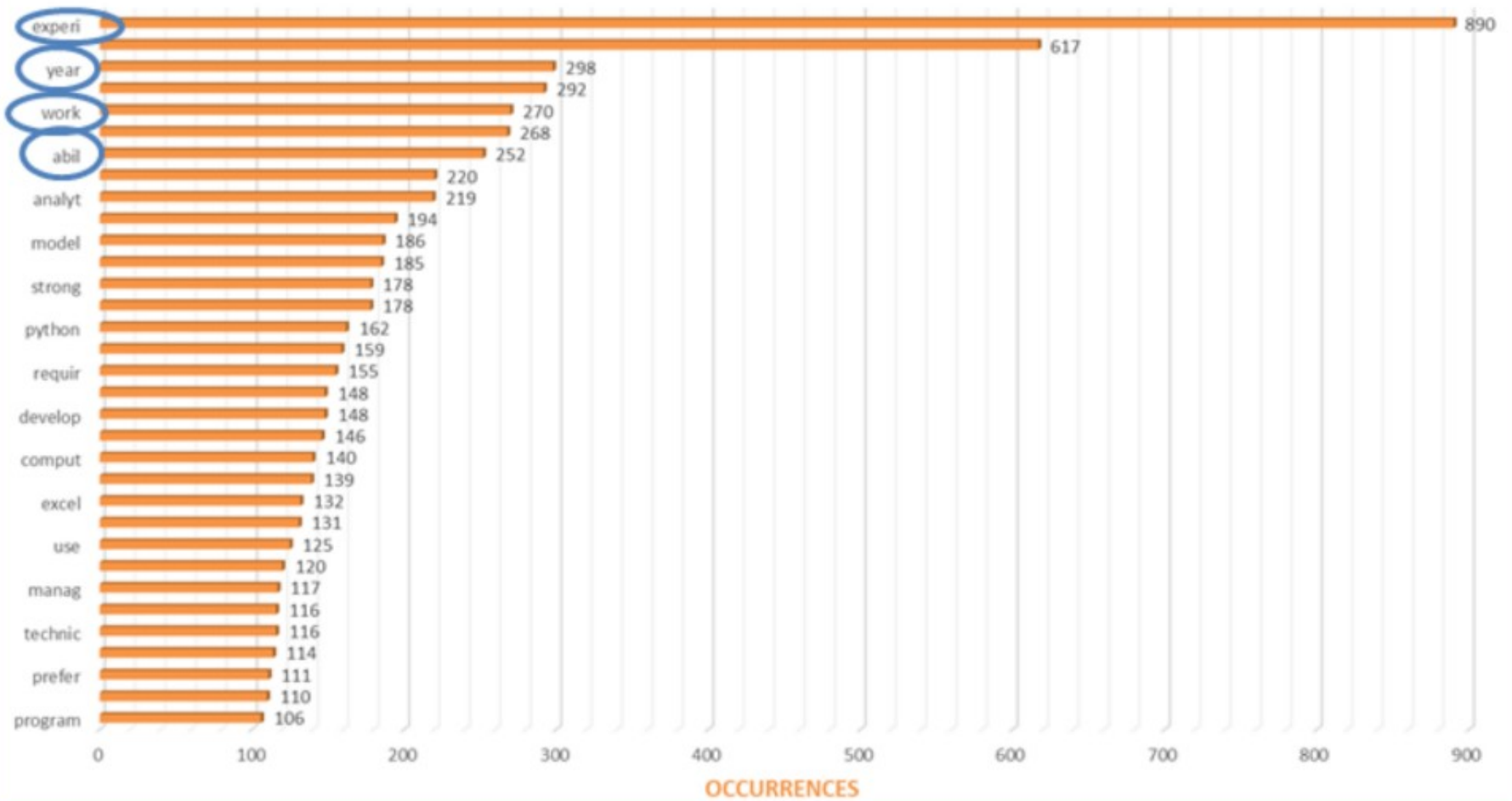
Top Phrases Containing 3 Words



Top Phrases Containing 2 Words



Unfiltered Word Count



Suggestion & Model

- Contain more specific & professional skills
- More consistent and pertinent to the Industry
- Try to avoid too general words, e.g. oral, written, ability, work
- Be gentle on the experience requirements
- Model to predict whether a job posting is popular or not based on the skills and qualification description



Linked in

python

Octopaste



Further
Ideas

ISGB 7978 Web Analytics

Group 7
Choudhury, Sumi
Gu, Jiahua
Huang, Feifei
Ji, Minxia

The slide features a large blue circle on the left side, which serves as a container for the title and list. The background of the slide is white, with a soft, abstract cloud-like shape in shades of orange and yellow on the left side, partially overlapping the blue circle.

Further Ideas

- Try modeling with Deep Learning, such as CNN
- Try to include more data
- Through this research, we noticed some other related topics that we can dig into, such as the impact of Data Science thriving in the Fashion Industry



Linked in

python

Octopaste



Results

Data Source
& Collection

Algorithm

Mission
&
Design

Further
Ideas

Visualization

ISGB 7978 Web Analytics

Group 7
Choudhury, Sumi
Gu, Jiahua
Huang, Feifei
Ji, Minxia