



National Science Foundation (NSF) Awards Granted History Analysis

Prepared and presented by:

YARA SEIF

MOHAMMAD AL-SHAMI

FAIRUS TANZIM

ASIF KAMAL TURZO

RADUANUL CHOWDHURY

Abstract

The National Science Foundation (NSF) are the funding source for approximately 25 percent of all federally supported basic research conducted by America's colleges and universities. Our team is conducting an analysis on the NSF Award data shown over the years and visualize the data to conclude awards analysis granted, to do that we collect/export our XML formatted data directly from the provided link (<https://www.nsf.gov/awardsearch/download.jsp>) We parsed the data using Python (pandas data frame) and extracted the file as CSV to start our analysis. Then we used multiple software to demonstrate our findings such as (Tableau, PowerBI, and Jupiter Notebook.)

Post our visualization we were able to analyze the spread of the awards are in the East coast, West coast, and South of the United States, also showed the highest dollar amount of NSF grant was given to Directorate for Mathematical & Physical Science, a total of \$30 Billion. We also analyzed the abstract and title of the NSF proposal using the Natural Language Processing (NLP) technique. We identified that Machine Learning/Data analysis is the topic that has received the most grant.

Background

NSF funds research and education in most fields of science and engineering. They do this through grants and cooperative agreements to more than 2,000 colleges, universities, K-12 school systems, businesses, informal science organizations and other research organizations throughout the U.S. The Foundation considers proposals submitted by organizations on behalf of individuals or groups for support in most fields of research. Interdisciplinary proposals also are eligible for consideration. Awardees are chosen from those who send us proposals asking for a specific amount of support for a specific project.

The goal of our class study is to analyze awards grants given by the NSF, National Science Foundation, to the different programs throughout the years. In our project we are examining more than 10 years of data, and we can conclude that the NSF awards different fellowship programs based on their importance relevant to the year we look at. Also, our dataset provides us with the national spread of the awards by state, also by money amount. Therefore, we can conclude that the more educationally developed the state or city is the more money amount we see, due to more accredited institutions for example, based on our heat map visual, we seem to find that NSF had given more award money in Massachusetts than it has to Idaho.

Methodology

In the study we have conducted on the NSF Awards data, we first collected the data in its XML schema file formats and parsed data in all years in the study and consolidated it into one CSV file format to be analyzed. We have used Python Jupyter and treating the dataset as a data frame using pandas, we used the code below:

```
1 import requests
2 import os
3 import csv
4 from zipfile import ZipFile
5 import xml.etree.ElementTree as ET
6
7
8 def file_download():
9     path = "datasets"
10    isExist = os.path.exists(path)
11    if not isExist:
12        os.makedirs(path)
13        print("The new directory is created!")
14
15    for i in range(2000,2022):
16        URL = "https://www.nsf.gov/awardsearch/download?DownloadFileName={}&All=true".format(i)
17        # Download the data behind the URL
18        response = requests.get(URL)
19        # Open the response into a new file
20        file_name = "datasets/{0}.zip".format(i)
21
22        open(file_name, "wb").write(response.content)
23        print("File download complete for year {0}".format(i))
24
25
26 def extract_zip():
27     with os.scandir('datasets/') as entries:
28         for entry in entries:
29             if(".zip" in entry.name):
30                 with ZipFile('datasets/'+entry.name, 'r') as zipObj:
31                     strn = entry.name.split(".")
32                     zipObj.extractall('datasets/'+strn[0])
33
34
35 def parseXML(xmlfile,headers = []):
36
37     tree = ET.parse(xmlfile)
38     root = tree.getroot()
39     file_name="dataset.csv"
40     exists = os.path.exists(file_name)
41     row_header=[]
42     row_value=[]
43     if exists:
44         row_header_2=[None] * len(headers)
45
46     for i in range(len(root[0])):
47         if root[0][i].text is not None:
48             if not root[0][i].text.strip():
49                 for j in range(len(root[0][i])):
50                     if root[0][i][j].text is not None:
51                         if not root[0][i][j].text.strip():
52                             for k in range(len(root[0][i][j])):
53                                 #print(str(root[0][i][j][k].tag)+ " : "+str(root[0][i][j][k].text))
54                                 if not exists:
55                                     row_value.append(str(root[0][i][j][k].text))
56                                     row_header.append(str(root[0][i][j].tag)+"_"+str(root[0][i][j][k].tag))
57                                 else:
58                                     s_match=str(root[0][i][j].tag)+"_"+str(root[0][i][j][k].tag)
59                                     ind=0
60                                     try:
61                                         ind=headers.index(s_match)
62                                     except ValueError:
63                                         ind=-1
64                                     if ind!=-1:
65                                         row_header_2[ind]=str(root[0][i][j][k].text)
66                             else:
67                                 #print(str(root[0][i][j].tag)+ " : "+str(root[0][i][j].text))
68                                 if not exists:
69                                     row_value.append(str(root[0][i][j].text))
70                                     row_header.append(str(root[0][i].tag)+"_"+str(root[0][i][j].tag))
71                             else:
72                                 s_match=str(root[0][i].tag)+"_"+str(root[0][i][j].tag)
73                                 ind=0
```

```

74         try:
75             ind=headers.index(s_match)
76         except ValueError:
77             ind=-1
78             if ind!=-1:
79                 row_header_2[ind]=str(root[0][i][j].text)
80     else:
81         #print(str(root[0][i].tag)+ " : "+str(root[0][i].text))
82         if not exists:
83             row_value.append(str(root[0][i].text))
84             row_header.append(str(root[0][i].tag))
85         else:
86             s_match=str(root[0][i].tag)
87             ind=0
88             try:
89                 ind=headers.index(s_match)
90             except ValueError:
91                 ind=-1
92                 if ind!=-1:
93                     row_header_2[ind]=str(root[0][i].text)
94
95     if not exists:
96         with open(file_name, 'w', newline='') as f:
97             write = csv.writer(f)
98             write.writerow(row_header)
99
100     with open(file_name, 'a', newline='') as f:
101         write = csv.writer(f)
102         if not exists:
103             write.writerow(row_value)
104         else:
105             write.writerow(row_header_2)
106         #print("writing complete for file: "+xmlfile)
107
108     if len(headers)==0:
109         return row_header
110     else:
111         return None
112
113
114 def create_dataset():
115     rootdir = 'datasets/'
116     file_name="dataset.csv"
117     flag=0
118     headers=[]
119     exists = os.path.exists(file_name)
120     if exists:
121         os.remove(file_name)
122     list=sorted(os.listdir(rootdir),reverse=True)
123     for file in list:
124         d = os.path.join(rootdir, file)
125         if os.path.isdir(d):
126             #print(d)
127             file_list = os.listdir(d)
128             for file in file_list:
129                 try:
130                     if flag==0:
131                         headers=parseXML(d+"/"+file)
132                         flag=1
133                     else:
134                         parseXML(d+"/"+file,headers)
135                 except:
136                     continue
137
138 #file_download()
139 #extract_zip()
140 create_dataset()
141
142

```

We parsed the data into useful fields (shown below) to be able to begin our text analysis, and visualization using Tableau and PowerBI:

A	B	C	D	E	F	G	H	I	J	K
AwardTitle	AGENCY	AwardEffectiveDate	AwardExpirationDate	AwardTotalInAmount	AwardAmount	AwardInstrumentValue	OrganizationCode	Directorate Abbreviation	Directorate LongName	Division Abbreviation
1	Collaborative NSF	10/1/2022	9/30/2025	214373	69413	Continuing Grant	6040200	GEO	Directorate For Geoscience OCE	
2	The Role of Int NSF	10/1/2022	12/31/2022	318962	34403	Standard Grant	7020000	ENG	Directorate For Engineering CBT	
3	CNS Core Small NSF	10/1/2022	6/30/2023	484041	484041	Standard Grant	5050000	CSH	Direct For Computer & Info CSH	
4	GEM Modeling NSF	10/1/2022	6/30/2024	449411	388928	Standard Grant	6020200	GEO	Directorate For Geoscience AGS	
5	RAPID Critical NSF	11/1/2022	4/30/2023	49970	49970	Standard Grant	6030000	GEO	Directorate For Geoscience EAR	
6	Collaborative NSF	11/1/2022	7/31/2023	399803	130038	Continuing Grant	6040300	GEO	Directorate For Geoscience OCE	
7	RUI SpecEEC NSF	10/1/2022	12/31/2023	250000	179126	Standard Grant	7010000	ENG	Directorate For Engineering ECOS	
8	Collaborative NSF	10/1/2022	12/31/2025	250000	250000	Standard Grant	7010000	ENG	Directorate For Engineering ECOS	
9	CAREER Deep NSF	10/1/2022	4/30/2024	550000	911351	Continuing Grant	5020000	CSH	Direct For Computer & Info IIS	
10	Conference: H NSF	11/1/2022	10/31/2023	9000	9000	Standard Grant	3010000	MPS	Direct For Mathematical & IPHY	
11	Iterative Algor NSF	10/1/2022	6/30/2023	300000	259625	Continuing Grant	3040000	MPS	Direct For Mathematical & IDMS	
12	FET Small: Opt NSF	10/1/2022	9/30/2024	389995	389995	Standard Grant	5010000	CSH	Direct For Computer & Info CCF	
13	Collaborative NSF	10/1/2022	11/30/2024	368331	333135	Standard Grant	8010000	BIO	Direct For Biological Scienc DEB	
14	Research Initit NSF	10/1/2022	8/31/2024	199998	199998	Standard Grant	7050000	ENG	Directorate For Engineering EEC	
15	Collaborative NSF	11/1/2022	10/31/2023	24559	24559	Standard Grant	6030000	GEO	Directorate For Geoscience EAR	
16	Collaborative NSF	11/1/2022	10/31/2023	21750	21750	Standard Grant	6030000	GEO	Directorate For Geoscience EAR	
17	RR: The Validat NSF	10/1/2022	8/31/2023	704337	149554	Standard Grant	4050000	SBE	Direct For Social, Behav & SES	
18	Collaborative NSF	11/1/2022	9/30/2025	251472	251472	Standard Grant	5020000	CSH	Direct For Computer & Info IIS	
19	Collaborative NSF	10/1/2022	9/30/2026	275000	275000	Standard Grant	7010000	ENG	Directorate For Engineering ECOS	
20	RUI Simulatio NSF	11/1/2022	8/31/2023	282833	120814	Standard Grant	3090000	MPS	Direct For Mathematical & ICH	
21	CNS Core Small NSF	10/1/2022	10/31/2023	500000	440034	Standard Grant	5050000	CSH	Direct For Computer & Info CSH	
22	Collaborative NSF	10/1/2022	9/30/2024	160241	137683	Standard Grant	5010000	CSH	Direct For Computer & Info CCF	
23	CAREER Round NSF	10/1/2022	7/31/2025	328649	483306	Standard Grant	7030000	ENG	Directorate For Engineering CMMI	
24	CHI: SaTe-Phy NSF	6/30/2022	6/30/2023	174428	152994	Standard Grant	5020000	CSH	Direct For Computer & Info CSH	
25	CAREER: Distrib NSF	10/15/2022	5/31/2026	489750	189676	Continuing Grant	5010000	CSH	Direct For Computer & Info CCF	
26	Collaborative NSF	10/1/2022	10/31/2023	250000	231303	Standard Grant	7030000	ENG	Directorate For Engineering CMMI	
27	CAREER Cognit NSF	11/1/2022	10/31/2025	731165	145597	Continuing Grant	11090000	EDU	Directorate For STEM Educat EDC	
28	Collaborative NSF	11/1/2022	7/31/2025	48579	48579	Standard Grant	6040200	GEO	Directorate For Geoscience OCE	
29	CAREER: Optim NSF	11/1/2022	12/31/2023	701000	537452	Continuing Grant	3090000	MPS	Direct For Mathematical & ICH	
30	Collaborative NSF	10/1/2022	10/31/2023	79996	58089	Standard Grant	7010000	ENG	Directorate For Engineering ECOS	
31	CAREER: Tubul NSF	11/1/2022	12/31/2026	679989	598844	Continuing Grant	3090000	MPS	Direct For Mathematical & ICH	
32	Collaborative NSF	10/15/2022	8/31/2024	288583	133303	Standard Grant	6090300	GEO	Directorate For Geoscience OPP	
33	CS&E: Co NSF	10/1/2022	7/31/2023	270802	203480	Standard Grant	5090000	CSH	Direct For Computer & Info CAC	
34	Collaborative NSF	11/1/2022	10/31/2023	12007	12007	Standard Grant	6010000	GEO	Directorate For Geoscience EAR	
35	Collaborative NSF	11/1/2022	10/31/2023	7990	7990	Standard Grant	6010000	GEO	Directorate For Geoscience EAR	
36	NSF-BSP: Af- Sr NSF	10/1/2022	9/30/2024	433382	430132	Standard Grant	5010000	CSH	Direct For Computer & Info CCF	
37	NSF Student Tr NSF	11/1/2022	10/31/2023	5000	5000	Standard Grant	5090000	CSH	Direct For Computer & Info CAC	
38	CAREER: Educat NSF	11/15/2022	3/31/2023	495593	97162	Continuing Grant	3090000	MPS	Direct For Mathematical & ICH	
39	Collaborative NSF	10/15/2022	6/30/2023	944037	652363	Continuing Grant	6010000	GEO	Directorate For Geoscience RISE	

Natural Language Processing (NLP) Task Visualization:

NLP can be applied to the title of the projects that have received NSF grants to understand interesting characteristics. As well as we have access to the NSF-funded project's abstract narration. So, we can apply topic modeling techniques to find what are different topics that have received NSF grants. For achieving this, we collected titles of the projects that have received NSF grants from the years 2000 to 2023. To understand which keywords are frequent in NSF grant titles, we draw a WordCloud of the titles. Before drawing the WordCloud, we removed the stop words and unwanted character sequences from the titles. The WordCloud is presented below:



From WordCloud, we can visualize the most frequent words appearing in the NSF grant's title. We further identified some unwanted stop words such as 'sbir' from WordCloud and removed those stop words.

Then for identifying the underlying topics from the collection of titles, we used an unsupervised topic modeling technique. The algorithm we have used is called Latent Dirichlet Allocation (LDA) [1]. LDA is an unsupervised machine learning technique used for identifying topics from a text collection. LDA can also be used for classifying documents and grouping the documents based on topic similarity. We used the python package ‘gensim’ for implementing the LDA algorithm.

By analyzing the titles of NSF-funded projects, we can identify what are the topics that receive the most grants as well as visualizing which keywords appear frequently in NSF grant titles. Moreover, we have analyzed the abstracts of the funded NSF grants to understand the dominant topics and understand which document contains which topics.

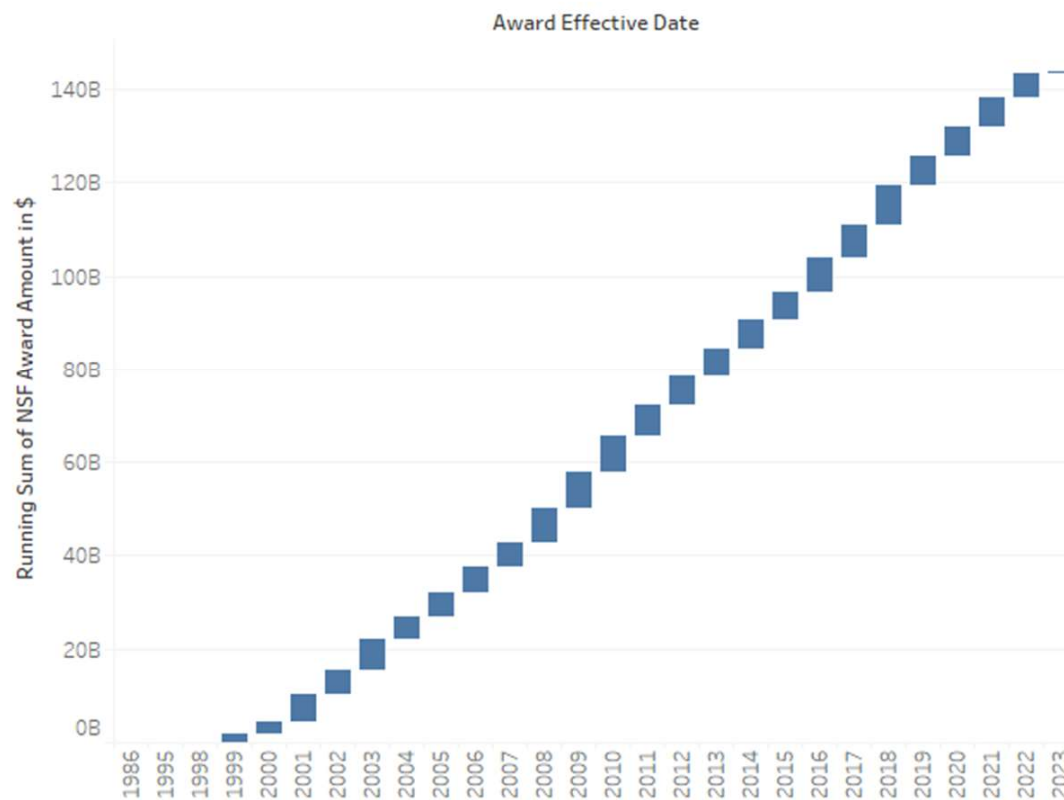
Result

Tableau Visualization

We have made several visualizations in the tableau desktop to analysis the data and made some recommendations.

At first, this Waterfall Chart is showing the running total amount of NSF grants awarded from the year 1986 to 2023. 3. Total of \$140B has been awarded over 35 years and highest amount of NSF grant was awarded in 2008, 2009, and 2010 and again in 2017, 2018, and 2019, which was more than \$7B.

Waterfall Chart for Total NSF Grant Awarded Each Year

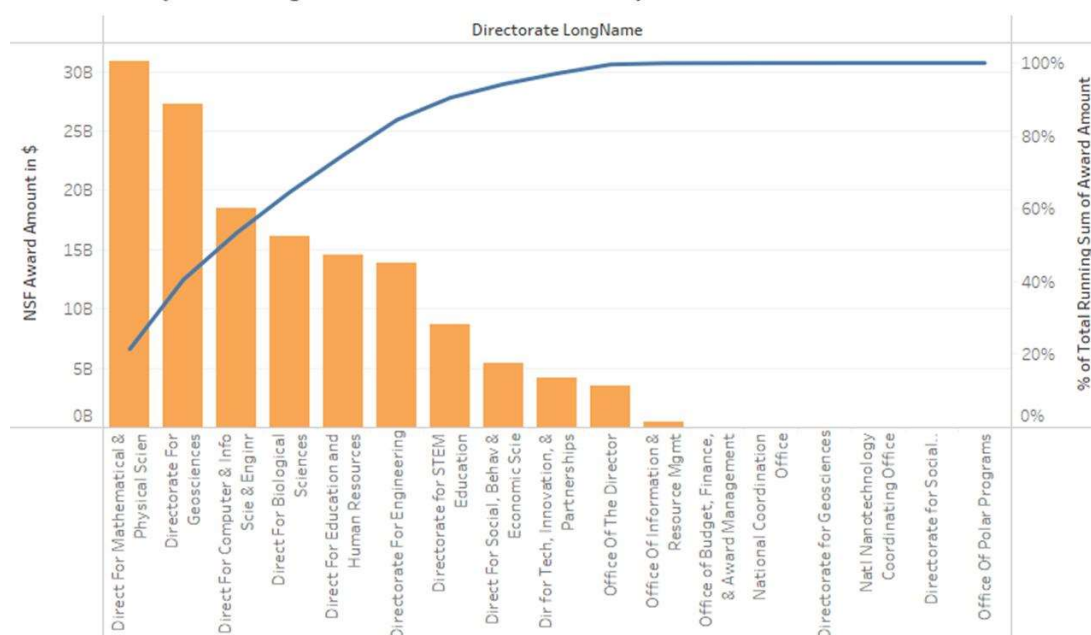


Then, we tried to find out which Directorate Field of studies received a significant amount of NSF grant. Pareto analysis shows that

- Directorate for Mathematical & Physical Science
- Directorate for Geoscience
- Directorate for Computer & Information Science & Engineering
- Directorate for Biological science
- Directorate for Education and human resource
- Directorate for Engineering

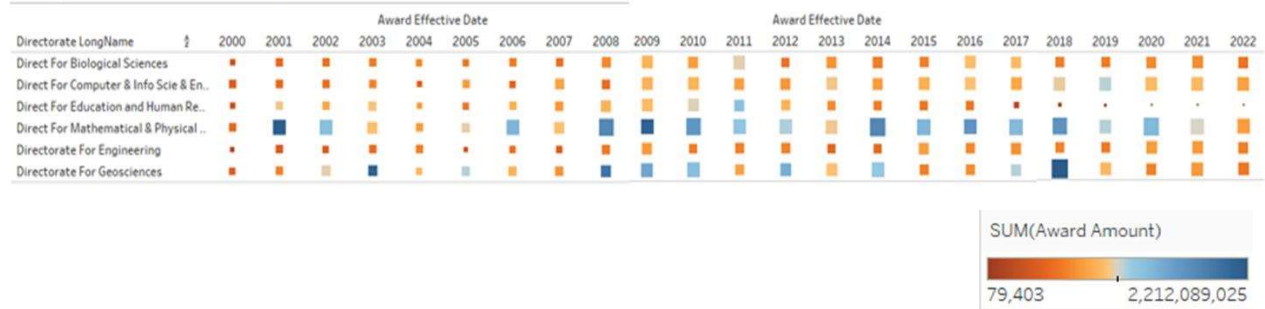
These fields received 80% of the total NSF grant over these 35 years.

Pareto Analysis for Highest NSF Grant Awarded by Field

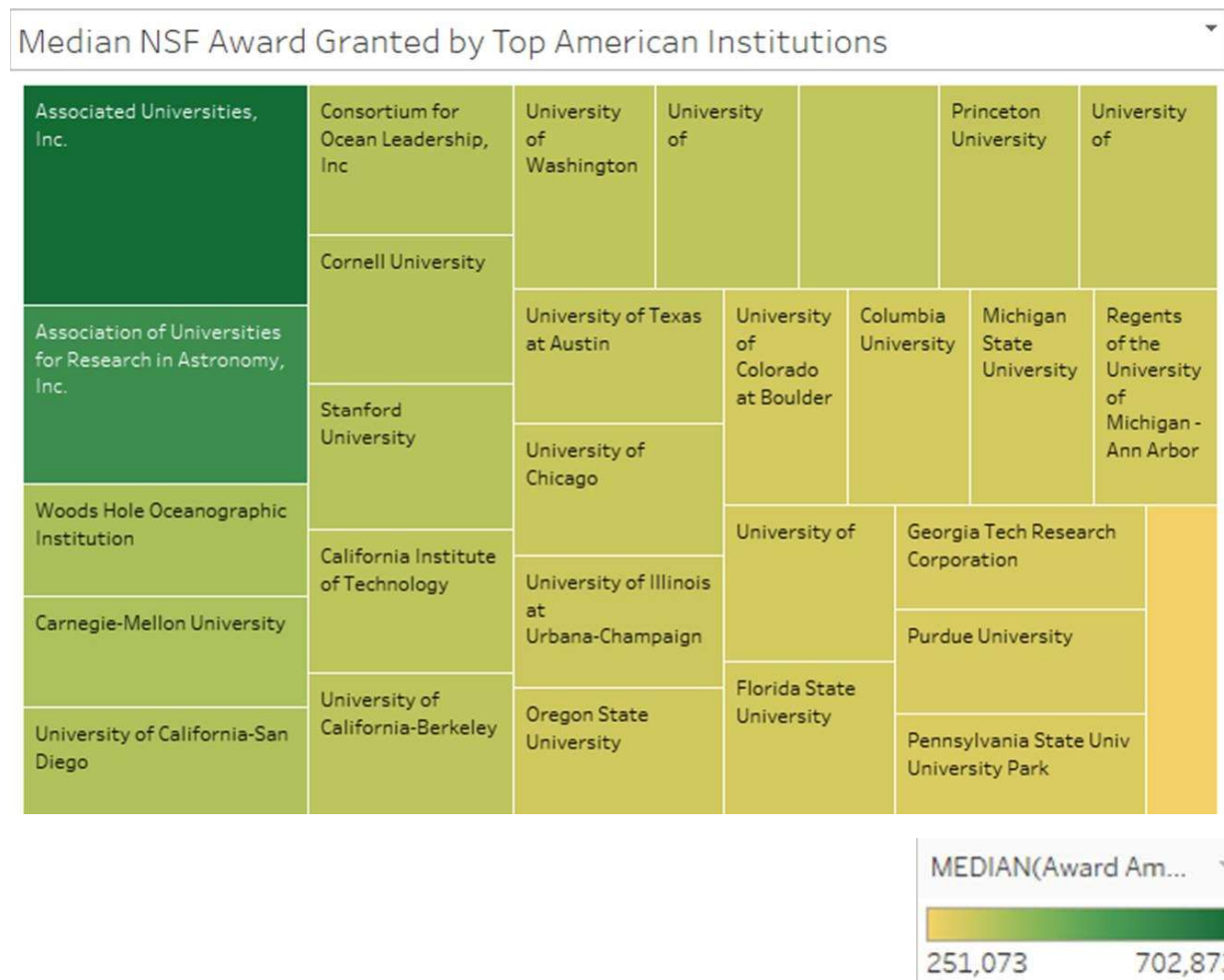


After that, we tried to see which Directorate field of study received the highest amount of total NSF Grant each year over last 22 years. The highest amount of NSF grants has been awarded to Directorate of Mathematical & Physical Science Directorate for last 22 years constantly. Whereas, Directorate of Geoscience has received over \$2B during 2018 and Grant for Directorate for Computer & Information Science & Engineering has been increased since 2018. On the other hand, Directorate for Education and Human Resource has been losing NSF grant since 2016.

Yearly Total NSF Grant Awarded to Most Field



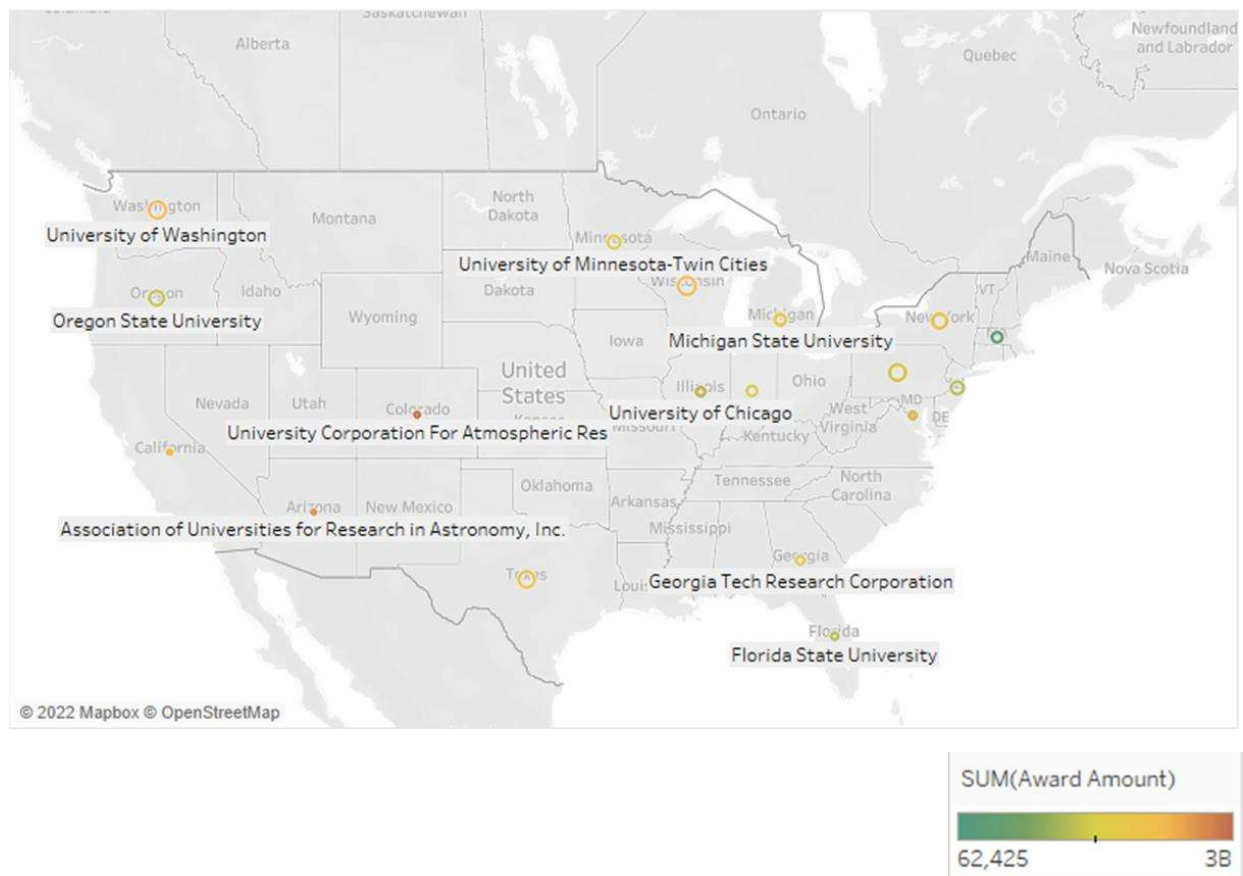
We also tried to analysis the institutions which received the highest median NSF grant each year over these 6 Directorate fields of studies. We listed 29 Institutions which received the highest Median NSF grant each year. These institutions are working on a focused 6 Directorate field of studies.



Finally, we tried to locate those 29 institutions in United States Map. We have found that all the institutions that are receiving highest median NSF grant are in East coast, West coast and South of United States. There are few institutions in Colorado and Arizona that received a good amount of NSF grant too. But the institutions from the Midwest of USA like Montana, Idaho, Wyoming, North and South Dakota and other states are not receiving competitive NSF grant each year.

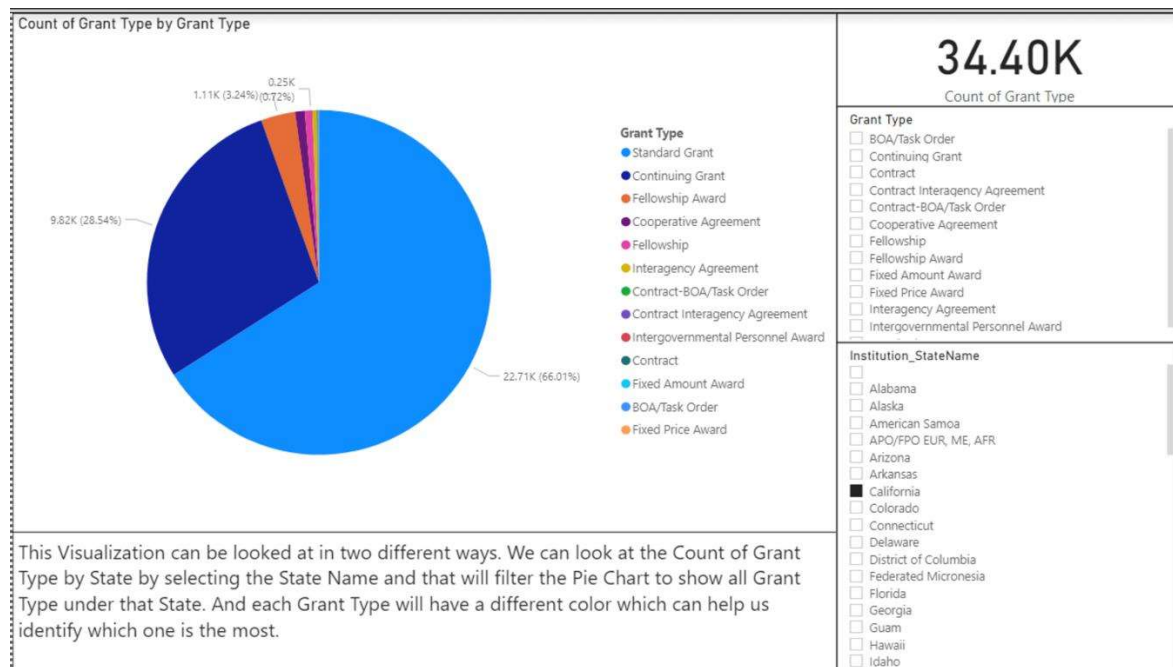
Proper steps need to be taken for the institutes from the Midwest of the USA to make them competitive with the institutions from the East and West coast of the USA. That is how overall educational improvement will happen.

Most NSF Grant Awarded Institution's State.

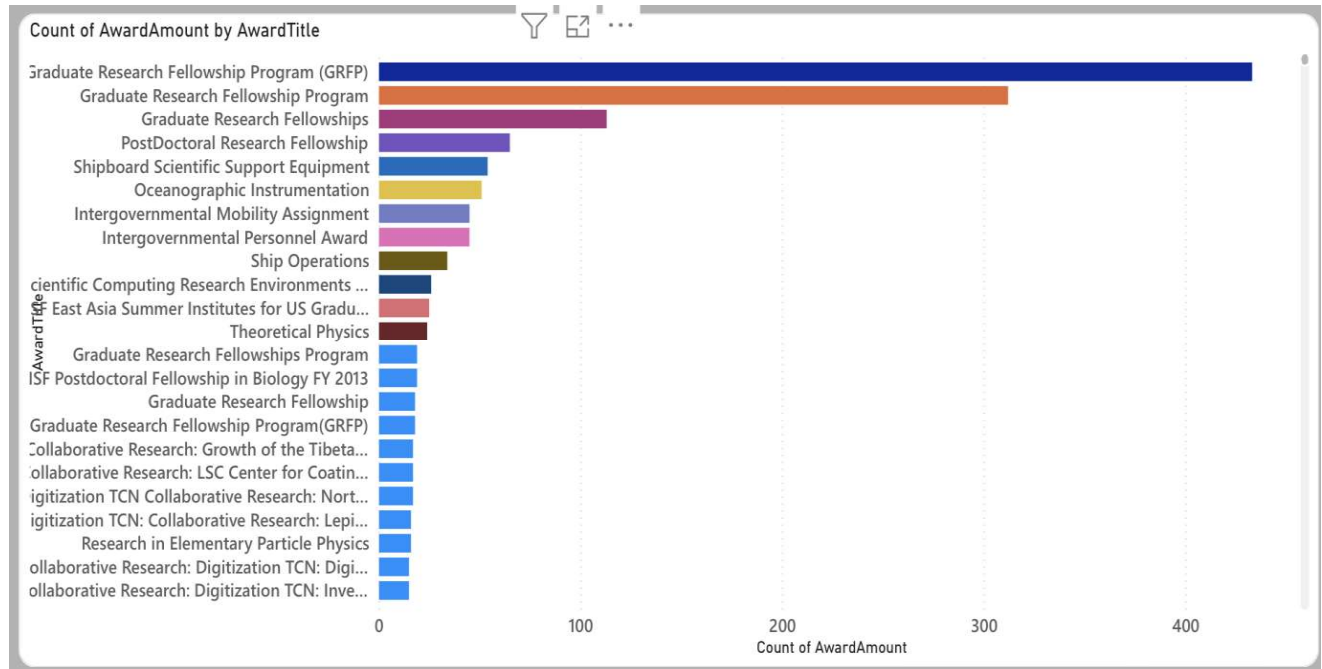


Power BI Visualizations

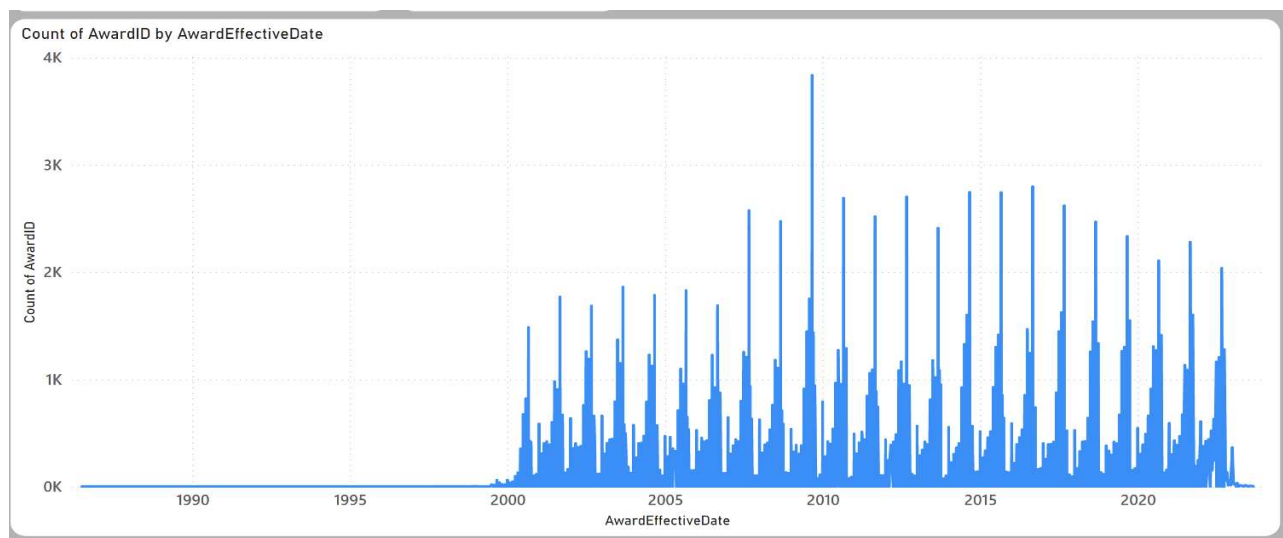
In the Visualization below, we are looking at the Count of Grant-by-Grant Type and State. When selecting the state, it will filter the Pie chart and display all Grant Types for that State. Each Grant Type has its own color which helps us quickly identify it. We can also look at it the other way around by filtering by Grant Type and that will show us the Count of that Grant in each State.



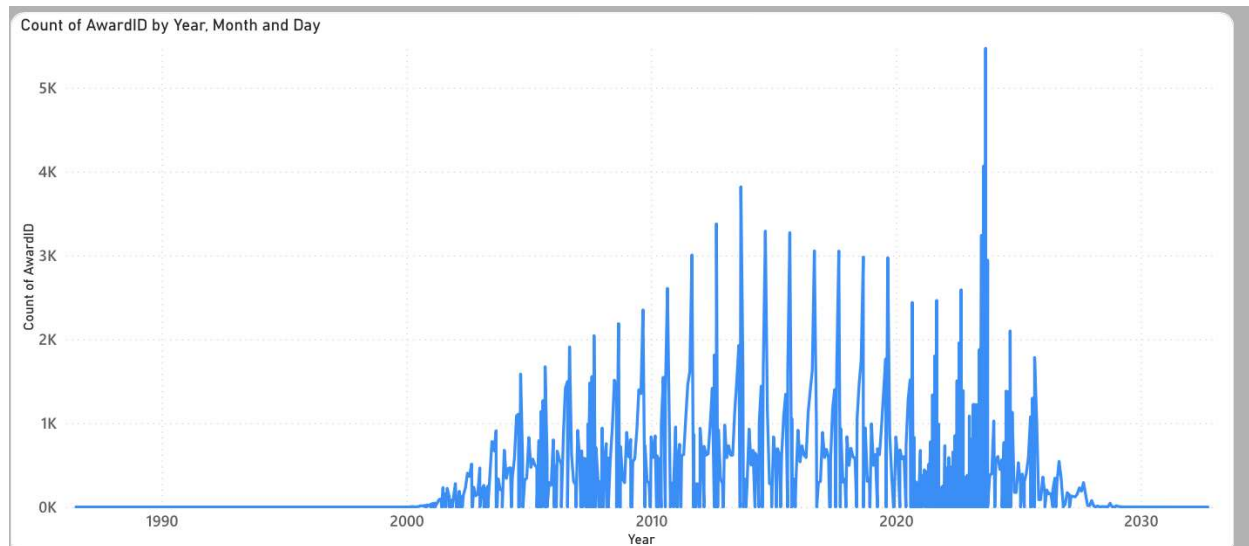
This dashboard displays the visual of Award amount by name, or title as it is displayed. The top Award granted in our National Science Foundation data is "Graduate Research Fellowship Program (GRFP)" with an award amount of 433. the Slicer provided to select the award title to help the data experts to analyze Award by single name.



This dashboard displays the bar chart of the Award effective date compared to the Award ID; we are showing when the Award was granted to those IDs in our data. Based on the visual below it's saying that the most awards granted were mostly clustered between 2009 and 2011 and it looks to continue to be in the high numbers of awards granted up to the year 2020.

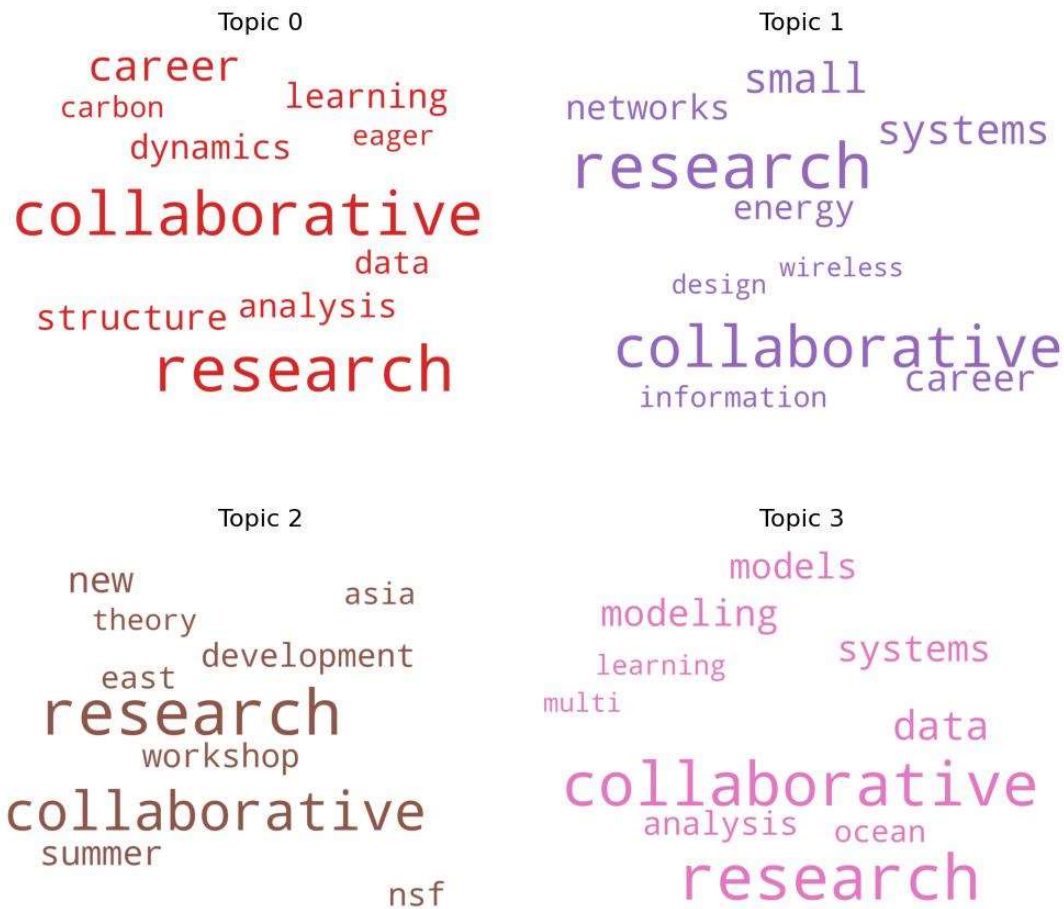


This dashboard displays the bar chart of the Award expiration date compared to the Award ID; we are showing when the expiration date was for the Award granted to those IDs in our data. Based on the visual below its saying that the most awards for expiration date for awards granted were mostly clustered to expire September 2023 and 2011 and to decline after to continue to fade in 2030.



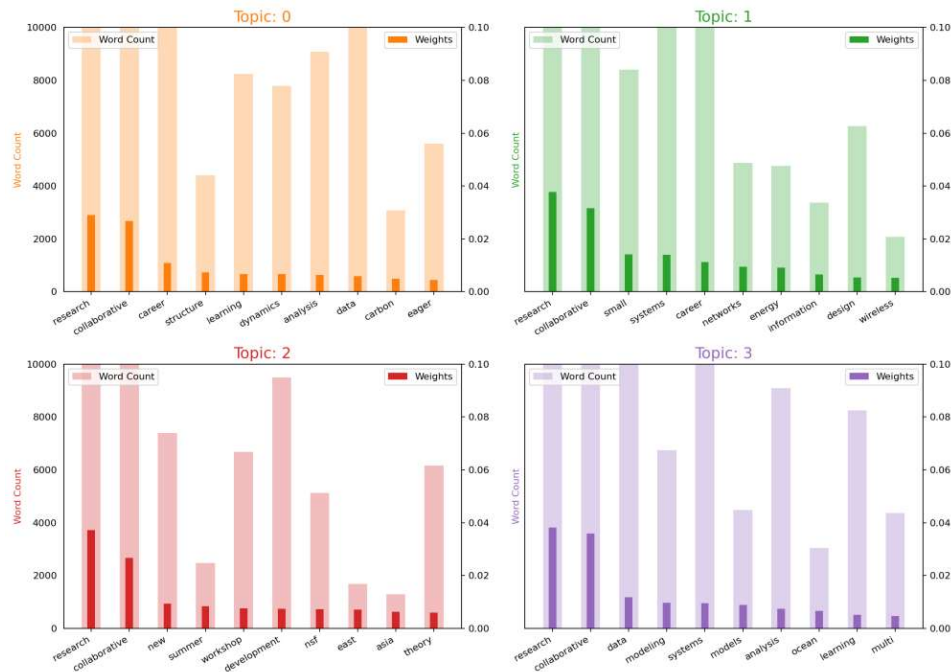
Visualize findings for NLP technique applied to NSF grant titles:

First, analyzing the titles of NSF grants, we have identified the topmost ten topics that received the most NSF grants. Below we present the topmost four topics that received the most NSF grant.



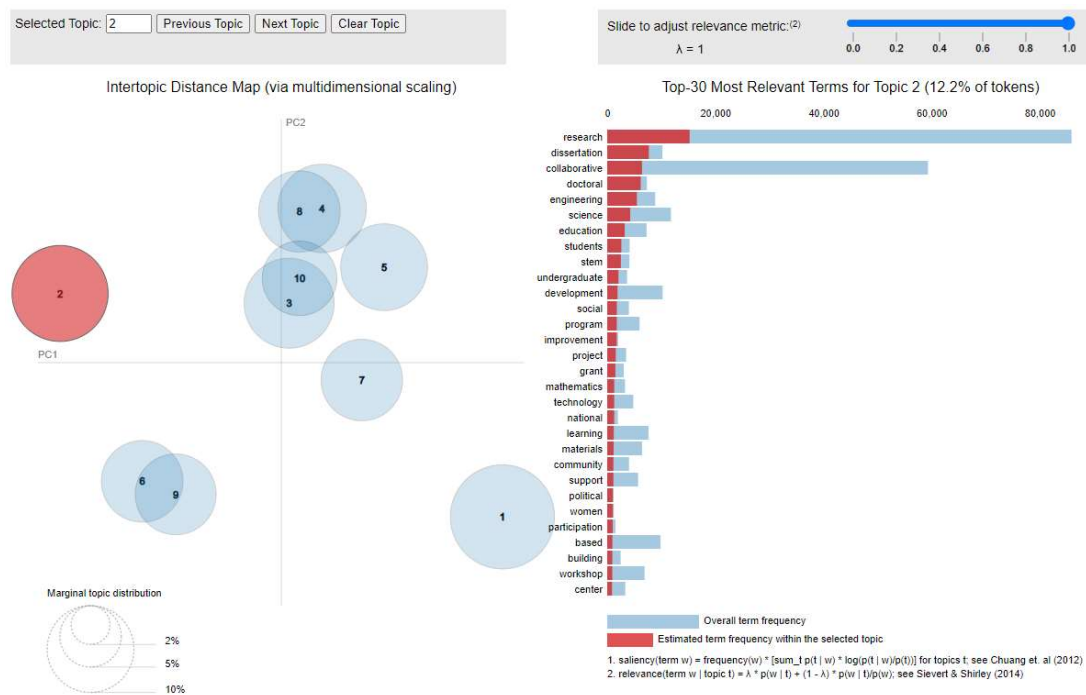
By visualizing the frequent words from each topic, we can easily identify what the topic is about. For example, in topic 0, some representative keywords are 'data', 'analysis', and 'learning'. So, we can guess that the topic that received the most grants is machine learning/data science topic.

Next, we analyzed the frequency of representative words from each topic to understand which keywords appear frequently in the NSF title. Below we present a bar chart to show the frequency of different words for each topic.



From the bar chart above, we can easily visualize the dominant keywords in NSF titles and understand their frequency.

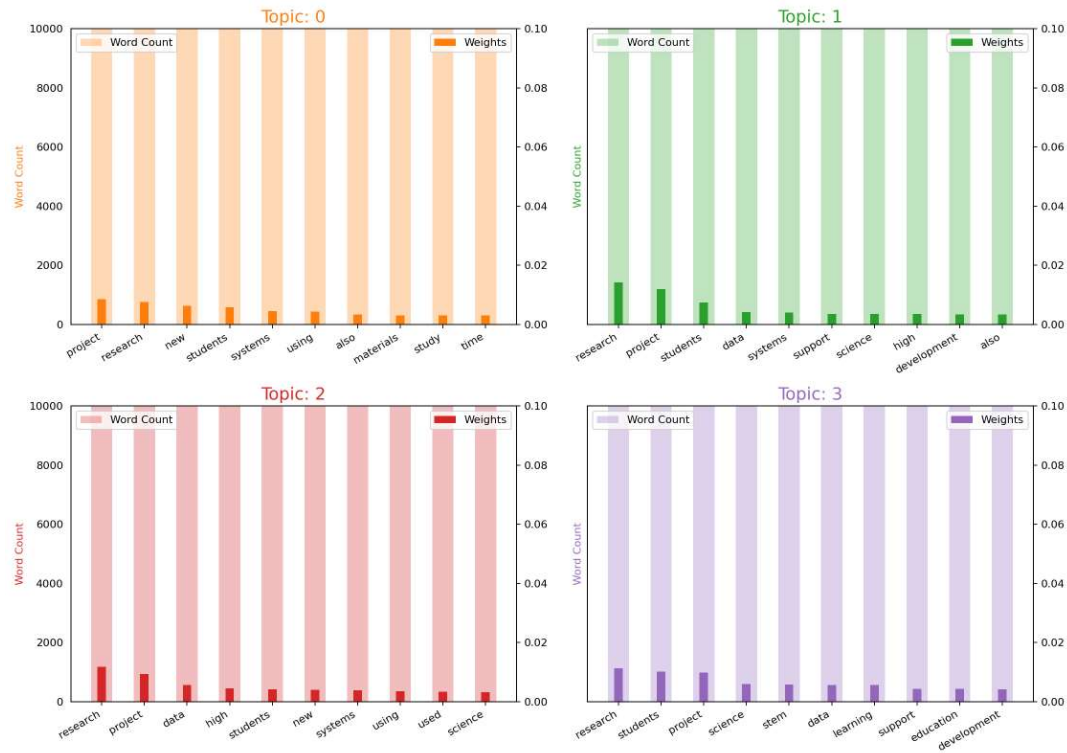
Our last visualization for the NSF title is an interactive visualization of LDA results. From this graph, we can understand what the keywords of the ten most important NSF-funded topics are. We can also visualize the most important keywords of each topic and also, the distance measure among the topics.



Here, we are showing the top 30 most relevant keywords for Topic 2. In the input box of 'Selected Topic', we can change the number to visualize keywords for different topics. We can also visualize the cluster distance among topics on the left side. For example, Topic 2 and Topic 1 have a large cluster distance, which indicates the two topics are completely different.

Visualize findings for NLP technique applied to NSF grant abstract:

Similarly, we applied LDA on the funded NSF grant abstract to identify the different keyword frequencies that appeared in the abstract of the NSF grant, cluster them into different topics, and to summarize each document's abstract. Below is the word frequency map for the top four topics.

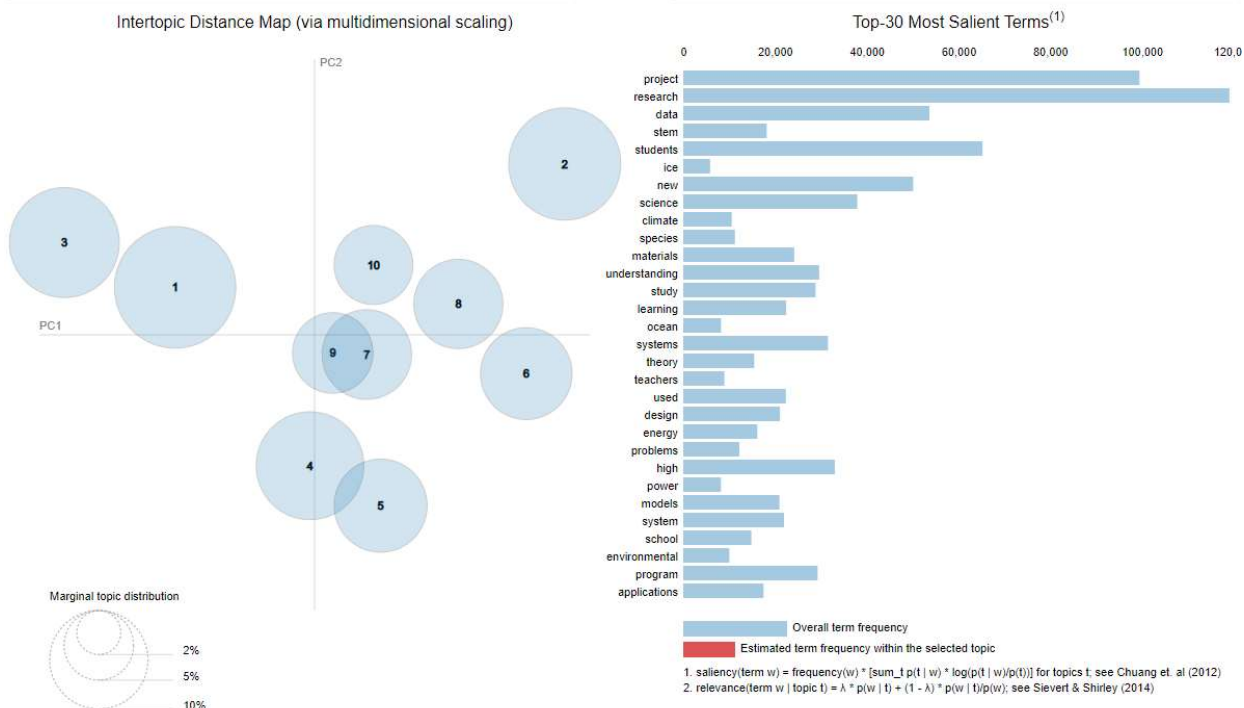


Next, we want to identify which topic is dominant within each abstract and text representation to summarize the document's abstract.

Document_No	Dominant_Topic	Topic_Perc_Contrib	Keywords	Text
0	0	0.0	0.9929 project, research, new, students, systems, usi...	[recent, surge, induced, seismicity, midwester...
1	1	6.0	0.9842 research, project, students, science, universi...	[objective, research, investigate, integration...
2	2	7.0	0.7434 research, project, new, understanding, data, m...	[simple, sub, ice, investigation, marine, plan...
3	3	9.0	0.9384 research, project, data, new, students, also, ...	[project, involves, research, empirical, metho...
4	4	8.0	0.4838 project, research, new, students, science, und...	[proposed, research, focuses, materials, growt...
5	5	6.0	0.6224 research, project, students, science, universi...	[gordon, award, supports, professor, mark, gor...
6	6	9.0	0.6650 research, project, data, new, students, also, ...	[broad, agreement, indigent, defense, counsel,...
7	7	5.0	0.3755 research, project, students, new, also, progra...	[assessing, assessments, historical, philosoph...
8	8	2.0	0.6124 research, project, data, high, students, new, ...	[yang, yong, conventional, cell, culture, meth...
9	9	7.0	0.7101 research, project, new, understanding, data, m...	[goal, research, understand, large, stand, rep...

In the table above, we can see which topic is dominant within each document, what are the keywords of the dominant topic, and last, the 'Text' column presents the representative text that summarizes the document.

Lastly, we present the LDA map summary for complete visualization of our LDA analysis result below.



Conclusion

Analyzing the NSF grant application from the year 2000 to 2022, we found which region in USA has received the most grant and uncovered that the highest dollar amount of NSF grant was given to Directorate for Mathematical & Physical Science which is a total of \$30 Billion. We also analyzed the abstract and title of the NSF proposal using the Natural Language Processing (NLP) technique.

We identified that Machine Learning/Data analysis is the topic that has received the most grant. We also identified descriptive keywords that appear most frequently in the NSF proposals and titles. We also presented a way to summarize each proposal's abstract and identified the underlying dominant topic of each proposal. The dataset collection techniques, dataset collection code, and NLP visualization code are attached to the document.

References

[1] Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." Journal of machine Learning research 3, no. Jan (2003): 993-1022.

[2] NSF Award Search. NSF AWARD SEARCH: Simple search. (n.d.). Retrieved December 10, 2022, from <https://www.nsf.gov/awardsearch>

Python code for data parsing: file:///Users/yara/Downloads/NSF_data_collection.py

Appendices

NLP visualization Code:

```
1 import pandas as pd
2 import os
3 import csv
4
5 import gensim
6 from gensim.utils import simple_preprocess
7 import nltk
8 nltk.download('stopwords')
9 from nltk.corpus import stopwords
10 from wordcloud import WordCloud
11 import pyLDAvis
12 import pyLDAvis.gensim_models as gensimvis
13 import gensim.corpora as corpora
14 from pprint import pprint
15 from matplotlib import pyplot as plt
16 from wordcloud import WordCloud, STOPWORDS
17 import matplotlib.colors as mcolors
18 from collections import Counter
19
20
21 #load dataset
22 input_fd = open('dataset_nlp.csv', encoding="utf8", errors='ignore')
23 dataset=pd.read_csv(input_fd)
24
25 ##### data preprocessing
26
27 # Remove punctuation
28 dataset['AwardTitle_clean'] = dataset['AwardTitle'].str.replace('[^\w\s]','')
29
30 # Convert the titles to lowercase
31 dataset['AwardTitle_clean'] = dataset['AwardTitle_clean'].str.lower()
32
33 print(dataset['AwardTitle_clean'].head())
34
35 # Join the different processed titles together.
36 long_string=dataset['AwardTitle_clean'].str.cat(sep=' ')
37
38 # Create a WordCloud object
39 wordcloud = WordCloud(background_color="white", max_words=5000, contour_width=3, contour_color='steelblue')
40 # Generate a word cloud
41 wordcloud.generate(long_string)
42 # Visualize the word cloud
43 image=wordcloud.to_image()
44 image.save("word_cloud.png")
45
46 stop_words = stopwords.words('english')
47 stop_words.extend(['from', 'subject', 're', 'edu', 'use', 'th', 'us', 'ii', 'fy', 'sbir'])
48
49 def sent_to_words(sentences):
50     for sentence in sentences:
51         # deacc=True removes punctuations
52         yield(gensim.utils.simple_preprocess(str(sentence), deacc=True))
53
54 def remove_stopwords(texts):
55     return [[word for word in simple_preprocess(str(doc))
56             if word not in stop_words] for doc in texts]
```

```

57 data = dataset.AwardTitle_clean.values.tolist()
58 data_words = list(sent_to_words(data))
59
60
61 # remove stop words
62 data_words = remove_stopwords(data_words)
63 print(data_words[:2][0][:20])
64
65 # Create Dictionary
66 id2word = corpora.Dictionary(data_words)
67 # Create Corpus
68 texts = data_words
69 # Term Document Frequency
70 corpus = [id2word.doc2bow(text) for text in texts]
71 # View
72 print(corpus[:1][0][:30])
73
74 # number of topics
75 num_topics = 10
76 # Build LDA model
77 lda_model = gensim.models.LdaMulticore(corpus=corpus,
78                                       id2word=id2word,
79                                       num_topics=num_topics)
80 # Print the Keyword in the 10 topics
81 pprint(lda_model.print_topics())
82 doc_lda = lda_model[corpus]
83
84
85 cols = [color for name, color in mcolors.TABLEAU_COLORS.items()] # more colors: 'mcolors.XKCD_COLORS'
86
87 cloud = WordCloud(stopwords=stop_words,
88                  background_color='white',
89                  width=2500,
90                  height=1800,
91                  max_words=10,
92                  colormap='tab10',
93                  color_func=lambda *args, **kwargs: cols[i+1],
94                  prefer_horizontal=1.0)
95
96 topics = lda_model.show_topics(formatted=False)
97
98 fig, axes = plt.subplots(2, 2, figsize=(10,10), sharex=True, sharey=True)
99
100 for i, ax in enumerate(axes.flatten()):
101     fig.add_subplot(ax)
102     topic_words = dict(topics[i][1])
103     cloud.generate_from_frequencies(topic_words, max_font_size=300)
104     plt.gca().imshow(cloud)
105     plt.gca().set_title('Topic ' + str(i), fontdict=dict(size=16))
106     plt.gca().axis('off')
107
108
109 plt.subplots_adjust(wspace=0, hspace=0)
110 plt.axis('off')
111 plt.margins(x=0, y=0)
112 plt.tight_layout()
113
114
115 fig1 = plt.gcf()
116 plt.show()
117 fig1.savefig('topic_wordcloud.png')
118
119 topics = lda_model.show_topics(formatted=False)
120 data_flat = [w for w_list in data_words for w in w_list]
121 counter = Counter(data_flat)
122
123 out = []
124 for i, topic in topics:
125     for word, weight in topic:
126         out.append([word, i, weight, counter[word]])
127
128 df = pd.DataFrame(out, columns=['word', 'topic_id', 'importance', 'word_count'])
129
130 # Plot Word Count and Weights of Topic Keywords
131 fig, axes = plt.subplots(2, 2, figsize=(14,10), sharey=True, dpi=160)
132 cols = [color for name, color in mcolors.TABLEAU_COLORS.items()]
133
134 for i, ax in enumerate(axes.flatten()):
135     ax.bar(x='word', height='word_count', data=df.loc[df.topic_id==i, :], color=cols[i+1], width=0.5, alpha=0.3, label='Word Count')
136     ax_twin = ax.twinx()
137     ax_twin.bar(x='word', height='importance', data=df.loc[df.topic_id==i, :], color=cols[i+1], width=0.2, label='Weights')
138     ax.set_ylabel('Word Count', color=cols[i+1])
139     ax_twin.set_ylim(0, 0.10); ax.set_ylim(0, 10000)
140     ax.set_title('Topic: ' + str(i), color=cols[i+1], fontsize=16)
141     ax.tick_params(axis='y', left=False)
142     ax.set_xticklabels(df.loc[df.topic_id==i, 'word'], rotation=30, horizontalalignment='right')
143     ax.legend(loc='upper left'); ax_twin.legend(loc='upper right')
144
145 fig.tight_layout(w_pad=2)
146 fig.suptitle('Word Count and Weights of Topic Keywords', fontsize=22, y=1.05)
147 fig2 = plt.gcf()
148 plt.show()
149 fig2.savefig('topic_word_frequency.png')
150
151 pyLDAvis.enable_notebook()
152 vis = gensimvis.prepare(lda_model, corpus, dictionary=lda_model.id2word)
153 vis
154 pyLDAvis.save_html(vis, "title_visualiation.html")
155
156 dataset_abstract = dataset.sample(n = 50000)
157
158 dataset_abstract['AbstractNarration'] = dataset_abstract['AbstractNarration'].str.replace(['^\\w\\s'], ' ')
159
160 # Convert the titles to lowercase
161 dataset_abstract['AbstractNarration'] = dataset_abstract['AbstractNarration'].str.lower()
162
163 print(dataset_abstract['AbstractNarration'].head())
164
165 stop_words.extend(['it', 'gt', 'br', 'lt'])
166
167 def sent_to_words(sentences):
168     for sentence in sentences:
169         # deacc=True removes punctuations
170         yield(gensim.utils.simple_preprocess(str(sentence), deacc=True))

```



```

169
170 def remove_stopwords(texts):
171     return [[word for word in simple_preprocess(str(doc))
172             if word not in stop_words] for doc in texts]
173
174 data_abs = dataset_abstract.AbstractNarration.values.tolist()
175 data_words_abs = list(sent_to_words(data_abs))
176
177 # remove stop words
178 data_words_abs = remove_stopwords(data_words_abs)
179 print(data_words_abs[:2][0][:20])
180
181 # Create Dictionary
182 id2word_abs = corpora.Dictionary(data_words_abs)
183 # Create Corpus
184 texts_abs = data_words_abs
185 # Term Document Frequency
186 corpus_abs = [id2word_abs.doc2bow(text) for text in texts_abs]
187 # View
188 print(corpus_abs[:1][0][:30])
189
190 num_topics = 10
191 # Build LDA model
192 lda_model_abs = gensim.models.LdaMulticore(corpus=corpus_abs,
193                                           id2word=id2word_abs,
194                                           num_topics=num_topics)
195 # Print the Keyword in the 10 topics
196 pprint(lda_model_abs.print_topics())
197 doc_lda_abs = lda_model_abs[corpus_abs]
198
199 pyLDAvis.enable_notebook()
200 vis_abs = gensimvis.prepare(lda_model_abs, corpus_abs, dictionary=lda_model_abs.id2word)
201 vis_abs
202
203 def format_topics_sentences(ldamodel=None, corpus=corpus, texts=data):
204     # Init output
205     sent_topics_df = pd.DataFrame()
206
207     # Get main topic in each document
208     for i, row_list in enumerate(ldamodel[corpus]):
209         row = row_list[0] if ldamodel.per_word_topics else row_list
210         # print(row)
211         row = sorted(row, key=lambda x: (x[1]), reverse=True)
212         # Get the Dominant topic, Perc Contribution and Keywords for each document
213         for j, (topic_num, prop_topic) in enumerate(row):
214             if j == 0: # => dominant topic
215                 wp = ldamodel.show_topic(topic_num)
216                 topic_keywords = ", ".join([word for word, prop in wp])
217                 sent_topics_df = sent_topics_df.append(pd.Series([int(topic_num), round(prop_topic, 4), topic_keywords]), ignore_index=True)
218             else:
219                 break
220     sent_topics_df.columns = ['Dominant_Topic', 'Perc_Contribution', 'Topic_Keywords']
221
222     # Add original text to the end of the output
223     contents = pd.Series(texts)
224     sent_topics_df = pd.concat([sent_topics_df, contents], axis=1)
225
226
227
228
229
230 # Format
231 df_dominant_topic = df_topic_sents_keywords.reset_index()
232 df_dominant_topic.columns = ['Document_No', 'Dominant_Topic', 'Topic_Perc_Contrib', 'Keywords', 'Text']
233 df_dominant_topic.head(10)
234
235 topics = lda_model_abs.show_topics(formatted=False)
236 data_flat = [w for w_list in data_words_abs for w in w_list]
237 counter = Counter(data_flat)
238
239 out = []
240 for i, topic in topics:
241     for word, weight in topic:
242         out.append([word, i, weight, counter[word]])
243
244 df = pd.DataFrame(out, columns=['word', 'topic_id', 'importance', 'word_count'])
245
246 # Plot Word Count and Weights of Topic Keywords
247 fig, axes = plt.subplots(2, 2, figsize=(14, 10), sharey=True, dpi=160)
248 cols = [color for name, color in mcolors.TABLEAU_COLORS.items()]
249
250 for i, ax in enumerate(axes.flatten()):
251     ax.bar(x='word', height='word_count', data=df.loc[df.topic_id==i, :], color=cols[i+1], width=0.5, alpha=0.3, label='Word Count')
252     ax_twin = ax.twinx()
253     ax_twin.bar(x='word', height='importance', data=df.loc[df.topic_id==i, :], color=cols[i+1], width=0.2, label='Weights')
254     ax.set_ylabel('Word Count', color=cols[i+1])
255     ax_twin.set_ylim(0, 0.10); ax.set_ylim(0, 10000)
256     ax.set_title('Topic: ' + str(i), color=cols[i+1], fontsize=16)
257     ax.tick_params(axis='y', left=False)
258     ax.set_xticklabels(df.loc[df.topic_id==i, 'word'], rotation=30, horizontalalignment='right')
259     ax.legend(loc='upper left'); ax_twin.legend(loc='upper right')
260
261 fig.tight_layout(w_pad=2)
262 fig.suptitle('Word Count and Weights of Topic Keywords', fontsize=22, y=1.05)
263 fig3 = plt.gcf()
264 plt.show()
265 fig3.savefig('abstract_word_frequency.png')
266
267

```