

Appendix A: API

ADS 509: Final Project Code

Team 4

Zachariah Freitas and Brianne Bell

```
In [1]: # Loading necessary libraries
import pandas as pd
import numpy as np
import os
import re
import random
import time
import datetime

import nltk
from collections import Counter, defaultdict
from nltk.corpus import stopwords
from string import punctuation
from tqdm import tqdm
import sqlite3

# from https://github.com/soumik12345/multi-label-text-classification/blob/master/arxiv
import arxiv
```

```
In [2]: # doing similar set up with setting up keywords to focus on
## Alternative keywords can be used to attempt better model performance or for different
query_keywords = [
    "\\"representation learning\\",
    "\\"image generation\\",
    "\\"object detection\\",
    "\\"transformers\\",
    "\\"image segmentation\\",
    "\\"natural language\\",
    "\\"graph\\",
    "\\"colorization\\",
    "\\"depth estimation\\",
    "\\"point cloud\\",
    "\\"structured data\\",
    "\\"reinforcement learning\\",
    "\\"attention\\",
    "\\"tabular\\",
    "\\"unsupervised learning\\",
    "\\"semi-supervised learning\\",
    "\\"explainable\\",
    "\\"time series\\",
    "\\"molecule\\",
    "\\"physics\\",
    "\\"graphics\\"
]
```

```
In [3]: # https://github.com/soumik12345/multi-label-text-classification/blob/master/arxiv_scr

# We are pulling just the terms (topic), titles, and abstracts of articles
## with a limit of 20k for each to save time and space
# other queries can be pulled if needed/wanted but we are focusing on titles and abstr

client = arxiv.Client(num_retries=20, page_size=500)

def query_with_keywords(query):
    search = arxiv.Search(
        query=query,
        max_results=20000,
        sort_by=arxiv.SortCriterion.LastUpdatedDate
    )
    terms = []
    titles = []
    abstracts = []
    for res in tqdm(client.results(search), desc=query):
        if res.primary_category in ["cs.CV", "stat.ML", "cs.LG"]:
            terms.append(res.categories)
            titles.append(res.title)
            abstracts.append(res.summary)
    return terms, titles, abstracts
```

```
In [4]: # setting up save file
# if not os.path.isdir("arxiv_data") :
#     os.mkdir("arxiv_data")
```

```
In [5]: # setting up for pull
all_titles = []
all_summaries = []
all_terms = []

# timing:
start_time = datetime.datetime.now()

# pulling
for query in query_keywords:
    terms, titles, abstracts = query_with_keywords(query)
    all_titles.extend(titles)
    all_summaries.extend(abstracts)
    all_terms.extend(terms)

# seeing how long ^that took:
end_time = datetime.datetime.now()
print(end_time - start_time)
```

```
"representation learning": 6118it [01:41, 60.09it/s]
"image generation": 1978it [00:30, 64.01it/s]
"object detection": 6536it [02:01, 53.69it/s]
"transformers": 20000it [06:55, 48.10it/s]
"image segmentation": 2890it [00:43, 66.93it/s]
"natural language": 13021it [03:36, 60.18it/s]
"graph": 20000it [05:39, 58.89it/s]
"colorization": 20000it [05:37, 59.20it/s]
"depth estimation": 1218it [00:20, 60.71it/s]
"point cloud": 4308it [01:30, 47.41it/s]
"structured data": 1915it [00:35, 54.30it/s]
"reinforcement learning": 16211it [04:36, 58.58it/s]
"attention": 20000it [05:48, 57.44it/s]
"tabular": 1382it [00:21, 63.47it/s]
"unsupervised learning": 2763it [00:41, 66.14it/s]
"semi-supervised learning": 0it [00:03, ?it/s]
"explainable": 20000it [06:17, 52.97it/s]
"time series": 15302it [04:17, 59.39it/s]
"molecule": 20000it [05:29, 60.67it/s]
"physics": 20000it [08:39, 38.50it/s]
"graphics": 15861it [04:34, 57.71it/s]
```

1:10:04.621341

Interpreting the scraping results:

- "representation learning": went through 6118 iterations at 60.09it/s.
- "image generation": went through 1978 iterations at 64.01it/s]
- "object detection": went through 6536 iterations at 53.69it/s.
- "transformers": went through 20000 iterations at 48.10it/s.
- "image segmentation": went through 2890 iterations at 66.93it/s.
- "natural language": went through 13021 iterations at 60.18it/s.
- "graph": went through 20000 iterations at 58.89it/s.
- "colorization": went through 20000 iterations at 59.20it/s.
- "depth estimation": went through 1218 iterations at 60.71it/s.
- "point cloud": went through 4308 iterations at 47.41it/s.
- "structured data": went through 1915 iterations at 54.30it/s.
- "reinforcement learning": went through 16211 iterations at 58.58it/s.
- "attention": went through 20000 iterations at 57.44it/s.
- "tabular": went through 1382 iterations at 63.47it/s.
- "unsupervised learning": went through 2763 iterations at 66.14it/s.
- "semi-supervised learning": went through 0 iterations at ?it/s.
- "explainable": went through 20000 iterations at 52.97it/s.
- "time series": went through 15302 iterations at 59.39it/s.
- "molecule": went through 20000 iterations at 60.67it/s.
- "physics": went through 20000 iterations at 38.50it/s.
- "graphics": went through 15861 iterations at 57.71it/s.

Of particular note, seven of the 21 keywords maxed out the number of iterations. They are: "transformers", "graph", "colorization", "attention", "explainable", and "physics". At the other

end of the spectrum is "semi-supervised learning" which went through zero iterations either because it is not in the archive or it is but in a different format without the hypen.

```
In [6]: raw_data = pd.DataFrame({
    'titles': all_titles,
    'abstracts': all_summaries,
    'terms': all_terms
})

raw_data.head()
```

	titles	abstracts	terms
0	Reinforcement Learning from Multiple Sensors v...	In many scenarios, observations from more than...	[cs.LG]
1	Interventional Causal Representation Learning	Causal representation learning seeks to extrac...	[stat.ML, cs.LG]
2	Self-Supervised Node Representation Learning v...	Self-supervised node representation learning a...	[cs.LG]
3	Out-of-Distribution Representation Learning fo...	Time series classification is an important pro...	[cs.LG, cs.AI]
4	Trading Information between Latents in Hierarc...	Variational Autoencoders (VAEs) were originall...	[stat.ML, cs.CV, cs.IT, cs.LG, math.IT]

```
In [7]: raw_data.shape #(64573, 3)
```

```
Out[7]: (64573, 3)
```

```
In [8]: # saving to csv file because pulling data takes a long while

# well I didn't get it to the file folder but it did write.

# This is how you save a pandas dataframe and all the details associated with it. You
raw_data.to_pickle('G:\\My Drive\\ADS-509_Final_Team_Project\\arxiv_data_2023_02_13.pkl')

# Nextime you want to use it. All you have to do is read it as a pickle file. The file
# raw_data = pd.read_pickle('G:\\My Drive\\ADS-509_Final_Team_Project\\arxiv_data_2023_02_13.pkl')
```

Appendix B: Descriptive Statistics

Assignment 6.2: Preparing Data for Final Team Project

ADS 509: Final Project Code

Team 4

Zachariah Freitas and Brianne Bell

Instructions:

1. Set up a GitHub repository for your final project. Make it public, or add your instructor as a collaborator. Include in this repository the code to assemble your final project data set. The best practice is to place this data creation code in its own notebook.
2. In a separate notebook, calculate descriptive statistics on your final-project data. You are welcome to reuse code from earlier modules.

Deliverable:

- When you have completed this notebook, run all cells, print the notebook as a PDF, and submit the PDF in Blackboard.
- Commit and push your code, so your repo is up to date.
- Enter your GitHub link as “online text” in the Blackboard assignment.

```
In [1]: # Import Libraries
import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
from collections import Counter
from nltk.corpus import stopwords
from string import punctuation
from wordcloud import WordCloud
import textacy.preprocessing as tprep
from lexical_diversity import lex_div as ld

# Some punctuation variations
punctuation = set(punctuation) # speeds up comparison

# Stopwords
sw = stopwords.words("english")
```

```
In [2]: # Get Imported Data
raw_data = pd.read_pickle('G:\\My Drive\\ADS-509_Final_Team_Project\\arxiv_data_2023_6
raw_data
```

Out[2]:

	titles	abstracts	terms
0	Reinforcement Learning from Multiple Sensors v...	In many scenarios, observations from more than...	[cs.LG]
1	Interventional Causal Representation Learning	Causal representation learning seeks to extract...	[stat.ML, cs.LG]
2	Self-Supervised Node Representation Learning v...	Self-supervised node representation learning a...	[cs.LG]
3	Out-of-Distribution Representation Learning fo...	Time series classification is an important pro...	[cs.LG, cs.AI]
4	Trading Information between Latents in Hierarc...	Variational Autoencoders (VAEs) were originall...	[stat.ML, cs.CV, cs.IT, cs.LG, math.IT]
...
64568	Plot 94 in ambiance X-Window	<PLOT > is a collection of routines to draw su...	[cs.CV, cs.GR]
64569	Automatic Face Recognition System Based on Loc...	We present an automatic face verification syst...	[cs.CV]
64570	Convexity Analysis of Snake Models Based on Ha...	This paper presents a convexity analysis for t...	[cs.CV, cs.GR, I.4; I.4.6;I.4.8]
64571	Semi-automatic vectorization of linear network...	A system for semi-automatic vectorization of l...	[cs.CV, cs.MM, I.4.6]
64572	Digital Color Imaging	This paper surveys current technology and rese...	[cs.CV, cs.GR, A.1;I.4,I.3.3,I.2.10;I.3.7;B.4.2]

64573 rows × 3 columns

In [3]:

```
# Helper Function - Descriptive Statistics
def descriptive_stats(tokens, top_n_tokens = 5, verbose=True) :
    """
        Given a list of tokens, print number of tokens, number of unique tokens,
        number of characters, lexical diversity (https://en.wikipedia.org/wiki/Lexical

diversity

) and num_tokens most common tokens. Return a list with the number of tokens, nu
```

```

top_words = Counter(tokens).most_common(top_n_tokens)

# InLine Printing
if verbose:
    print(f"There are {num_tokens} tokens in the data.")
    print(f"There are {num_unique_tokens} unique tokens in the data.")
    print(f"There are {num_characters} characters in the data.")
    print(f"The average token length in the data is {avg_token_len:.3f}.")
    print(f"The lexical diversity is {lexical_diversity:.3f} in the data.")

# print the five most common tokens
print(f"\n\nThe top {top_n_tokens} most common tokens")
print(top_words)

# Return Dictionary
results = { 'tokens' : num_tokens,
            'unique_tokens' : num_unique_tokens,
            'avg_token_length' : avg_token_len,
            'lexical_diversity': lexical_diversity,
            'num_characters': num_characters,
            'top_words': top_words}

return(results)

```

In [4]: # Helper Functions - Cleaning data

```

def normalize(text):
    text = tprep.normalize.hyphenated_words(text)
    text = tprep.normalize.quotation_marks(text)
    text = tprep.normalize.unicode(text)
    text = tprep.remove.accents(text)
    return text

def remove_punctuation(text, punct_set=punctuation) :
    """This function removes punctuation from a string."""
    return"".join([ch for ch in text if ch not in punct_set]))

def tokenize(text) :
    """ Splitting on whitespace rather than the book's tokenize function. That
        function will drop tokens like '#hashtag' or '2A', which we need for Twitter.

    # modify this function to return tokens
    return text.lower().strip().split()

def remove_stop(tokens) :
    """This function removes stopwords from a list of tokens."""
    return[t for t in tokens if t.lower() not in sw]

def prepare(text, pipeline) :
    """ This fuction manages and executes other functions like a pipeline. """
    tokens = str(text)

    for transform in pipeline :
        tokens = transform(tokens)

    return(tokens)

```

```
In [5]: # Helper Function - Flatten a List of Lists
def flatten_lists(list_of_lists):
    """This function flattens a list of lists into a single list."""
    return [i for s in list_of_lists for i in s]

In [6]: # Clean and Tokenize Data
df = raw_data.copy()

df["titles_abs"] = df["titles"] + df["abstracts"]

my_pipeline = [normalize, remove_punctuation, tokenize, remove_stop]
df["titles_tokens"] = df["titles"].apply(prepare, pipeline=my_pipeline)
df["abstract_tokens"] = df["abstracts"].apply(prepare, pipeline=my_pipeline)
df["titles_abs_tokens"] = df["titles_abs"].apply(prepare, pipeline=my_pipeline)

df
```

Out[6]:

	titles	abstracts	terms	titles_abs	titles_tokens	abstra
0	Reinforcement Learning from Multiple Sensors v...	In many scenarios, observations from more than...	[cs.LG]	Reinforcement Learning from Multiple Sensors v...	[reinforcement, learning, multiple, sensors, v...]	ob one, s
1	Interventional Causal Representation Learning	Causal representation learning seeks to extract...	[stat.ML, cs.LG]	Interventional Causal Representation LearningC...	[interventional, causal, representation, learn...]	repr learn
2	Self-Supervised Node Representation Learning v...	Self-supervised node representation learning a...	[cs.LG]	Self-Supervised Node Representation Learning v...	[selfsupervised, node, representation, learnin...]	[selfs repr
3	Out-of-Distribution Representation Learning fo...	Time series classification is an important pro...	[cs.LG, cs.AI]	Out-of-Distribution Representation Learning fo...	[outofdistribution, representation, learning, ...]	[ti cla
4	Trading Information between Latents in Hierarc...	Variational Autoencoders (VAEs) were originally...	[stat.ML, cs.CV, cs.IT, cs.LG, math.IT]	Trading Information between Latents in Hierarc...	[trading, information, latents, hierachical, ...]	[t auto vaes.
...
64568	Plot 94 in ambiance X-Window	<PLOT > is a collection of routines to draw su...	[cs.CV, cs.GR]	Plot 94 in ambiance X-Window<PLOT > is a colle...	[plot, 94, ambiance, xwindow]	[plot, rout s
64569	Automatic Face Recognition System Based on Loc...	We present an automatic face verification syst...	[cs.CV]	Automatic Face Recognition System Based on Loc...	[automatic, face, recognition, system, based, ...]	autor v
64570	Convexity Analysis of Snake Models Based on Ha...	This paper presents a convexity analysis for t...	[cs.CV, cs.GR, I.4; I.4.6;I.4.8]	Convexity Analysis of Snake Models Based on Ha...	[convexity, analysis, snake, models, based, ha...]	[paper
64571	Semi-automatic vectorization of linear network...	A system for semi-automatic vectorization of l...	[cs.CV, cs.MM, I.4.6]	Semi-automatic vectorization of linear network...	[semiautomatic, vectorization, linear, network...]	semi vec
64572	Digital Color Imaging	This paper surveys current technology and rese...	A.1;I.4,I.3.3,I.2.10;I.3.7;B.4.2]	Digital Color ImagingThis paper surveys curren...	[digital, color, imaging]	[paper t

64573 rows × 7 columns

In [7]:

```
# Statistics on Titles
title_results = descriptive_stats(flatten_lists(df['titles_tokens']))
```

There are 470005 tokens in the data.
There are 31401 unique tokens in the data.
There are 3829797 characters in the data.
The average token length in the data is 8.148.
The lexical diversity is 0.067 in the data.

The top 5 most common tokens
[('learning', 19062), ('detection', 6548), ('deep', 6117), ('neural', 5997), ('networks', 5569)]

In [8]:

```
# Statistics on Abstracts
abstract_results = descriptive_stats(flatten_lists(df['abstract_tokens']))
```

There are 7296290 tokens in the data.
There are 123974 unique tokens in the data.
There are 55306651 characters in the data.
The average token length in the data is 7.580.
The lexical diversity is 0.017 in the data.

The top 5 most common tokens
[('learning', 79692), ('data', 57018), ('model', 52447), ('methods', 41571), ('models', 39819)]

In [9]:

```
# Statistics on Titles & Abstracts
abstract_results =
title_abstract_results = descriptive_stats(flatten_lists(df['titles_abs_tokens']))
```

There are 7726299 tokens in the data.
There are 155200 unique tokens in the data.
There are 59201115 characters in the data.
The average token length in the data is 7.662.
The lexical diversity is 0.020 in the data.

The top 5 most common tokens
[('learning', 92779), ('data', 58746), ('model', 54062), ('image', 44154), ('methods', 42002)]

Word Cloud Plotting

In [10]:

```
# Helper Function - Plot Word Cloud
```

```
def wordcloud(word_freq, title=None, max_words=200, stopwords=None):
    """
        Given a list of tokens, print number of tokens, number of unique tokens,
        number of characters, lexical diversity (https://en.wikipedia.org/wiki/Lexical

diversity

) and num_tokens most common tokens. Return a list with the number of tokens, number of unique tokens, lexical diversity, and number of characters.
    """
    wc = WordCloud(width=800, height=400,
                    background_color="black", colormap="Paired",
                    max_font_size=150, max_words=max_words)
```

```

# convert data frame into dict
if type(word_freq) == pd.Series:
    counter = Counter(word_freq.fillna(0).to_dict())
else:
    counter = word_freq

# filter stop words in frequency counter
if stopwords is not None:
    counter = {token:freq for (token, freq) in counter.items()
               if token not in stopwords}
wc.generate_from_frequencies(counter)

plt.title(title)

plt.imshow(wc, interpolation='bilinear')
plt.axis("off")

def count_words(df, column='tokens', preprocess=None, min_freq=2):
    """
    Given a list of tokens, print number of tokens, number of unique tokens,
    number of characters, lexical diversity (https://en.wikipedia.org/wiki/Lexical

diversity

) and num_tokens most common tokens. Return a list with the number of tokens, nu
    of unique tokens, lexical diversity, and number of characters.
    """
    # process tokens and update counter
    def update(doc):
        tokens = doc if preprocess is None else preprocess(doc)
        counter.update(tokens)

    # create counter and run through all data
    counter = Counter()
    df[column].map(update)

    # transform counter into data frame
    freq_df = pd.DataFrame.from_dict(counter, orient='index', columns=['freq'])
    freq_df = freq_df.query('freq >= @min_freq')
    freq_df.index.name = 'token'

    return freq_df.sort_values('freq', ascending=False)

def plot_wc(wordcloud_df):
    """
    Given a list of tokens, print number of tokens, number of unique tokens,
    number of characters, lexical diversity (diversity) and num_tokens most common tokens. Return a list with the number of tokens, nu
    of unique tokens, lexical diversity, and number of characters.
    """
    plt.figure(figsize=(12,4))
    plt.subplot(1,2,1)###
    wordcloud(wordcloud_df['freq'], max_words=1000)
    plt.title("With Stop Words")

    plt.subplot(1,2,2)###
    wordcloud(wordcloud_df['freq'], max_words=1000, stopwords=sw)
    plt.title("Without Stop Words")
    plt.tight_layout()###

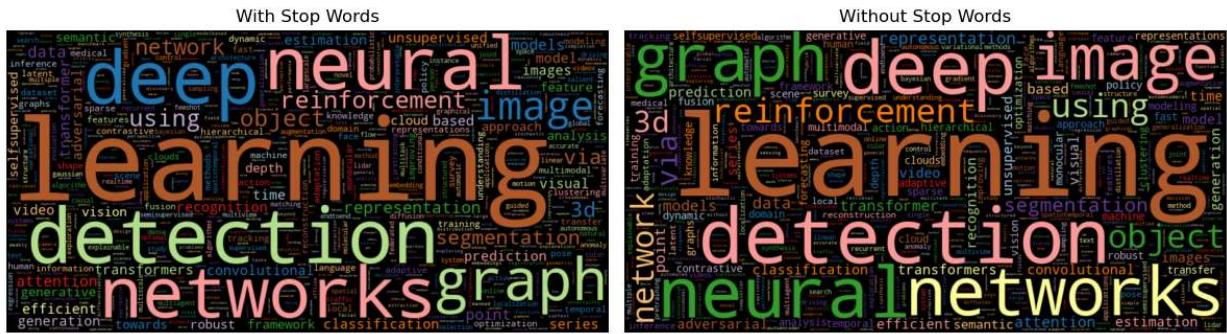
```

Corpus Word Cloud

Title Word Cloud

With stopwords and without stopwords.

```
In [11]: wordcloud_df = count_words(df, column='titles_tokens')  
plot_wc(wordcloud_df)
```



Title word cloud differences between with stop words and without shows that removing stop words increases the number of recurring words. This can be seen above where there are more large font words in the 'Without Stop Words' wordcloud (~10) than there are in the wordcloud with stopwords (~6).

Abstract Word Cloud

With stopwords and without stopwords.

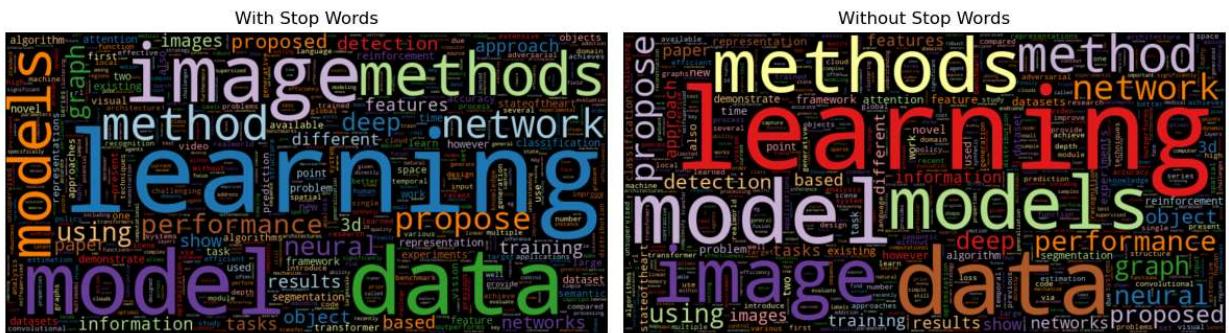
```
In [12]: wordcloud_df = count_words(df, column='abstract_tokens')  
plot_wc(wordcloud_df)
```



Title & Abstract Word Cloud

With stopwords and without stopwords.

```
In [13]: wordcloud_df = count_words(df, column='titles_abs_tokens')  
plot_wc(wordcloud_df)
```



Looking at the combination of abstract and title tokens, both the wordclouds (with and without stopwords) appear very similar. The most common tokens are 'learning', 'model', 'data', 'image', 'method', 'models', and 'network'. These are heavily leaning towards data science topics so it is likely that our models will be heavily populated by these and could impact modeling performance.