

# 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
```

of unique tokens, lexical diversity, and number of characters.

```
"""
    if (len(tokens) == 0):
        raise ValueError("Can't work with empty text object.")

    # These are placeholder values.
    num_tokens = 1
    num_unique_tokens = 0
    avg_token_len = 0.0
    lexical_diversity = 0.0

    # Calculate statistics
    num_tokens = len(tokens)
    num_unique_tokens = len(set(tokens))
    lexical_diversity = ld.ttr(tokens) # Simple TTR = Len(Counter(text))/Len(text)
    num_characters = sum([len(i) for i in tokens])
    avg_token_len = np.mean([len(w) for w in tokens])
```

```

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]:

	<b>titles</b>	<b>abstracts</b>	<b>terms</b>	<b>titles_abs</b>	<b>titles_tokens</b>	<b>abstra</b>
<b>0</b>	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
<b>1</b>	Interventional Causal Representation Learning	Causal representation learning seeks to extract...	[stat.ML, cs.LG]	Interventional Causal Representation LearningC...	[interventional, causal, representation, learn...]	repr learn
<b>2</b>	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
<b>3</b>	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
<b>4</b>	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.
...	...	...	...	...	...	...
<b>64568</b>	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
<b>64569</b>	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
<b>64570</b>	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
<b>64571</b>	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
<b>64572</b>	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()###

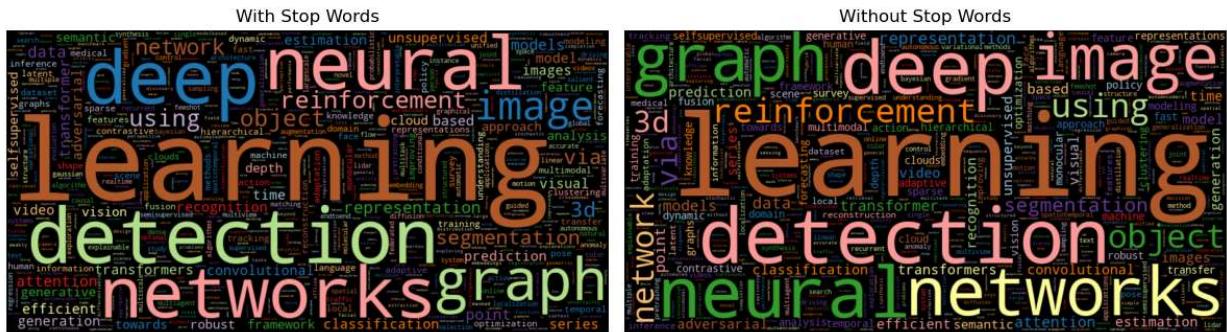
```

# 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.