



ALY 6040: DATA MINING APPLICATIONS

Assignment 5:
Text Mining and NLP on
Amazon Mobile Review Dataset

Submitted To:
Dr. Chinthaka Pathum Dinesh, PhD, Prof. Herath Gedara,
Faculty Lecturer

Submitted By:
[Abhilash Dikshit](#)

Academic Term: Spring 2023
Graduate Student at Northeastern University, Vancouver, BC,
Canada
Master of Professional Studies in Analytics

May 14, 2023

I. Introduction:

In this report, we will perform sentiment analysis on a dataset of Amazon reviews using text mining and natural language processing. The dataset contains reviews of various products on Amazon, along with their respective ratings. Our objective is to explore the data, clean and preprocess it, and build a predictive model to classify reviews into positive, negative, and neutral sentiment categories.

The dataset was provided by the professor.

Step 1: Data Cleaning

We begin by importing the necessary libraries and loading the dataset into a Pandas DataFrame.

```
# Basic Libraries
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from tabulate import tabulate
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
import nltk
from nltk.corpus import stopwords
from wordcloud import WordCloud, STOPWORDS

warnings.filterwarnings(action='ignore')

print('\033[1mAmazon Unlocked Mobile Dataset:\n' + '='*30 + '\033[0m')

# Read in the CSV file
path = '~/GitProjects/Datasets/Amazon_Unlocked_Mobile.csv'
df_raw = pd.read_csv(path)

print('\nDisplay Raw Dataset:\n')
display(df_raw)

# Display basic information about the dataset
table = [['Type', 'Length', 'Shape'], [type(df_raw), len(df_raw), df_raw.shape]]

print('\nDisplay Type, Length, Shape about the dataset:\n')
print(tabulate(table, headers='firstrow', tablefmt='fancy_grid'))
```

Fig 1

Amazon Unlocked Mobile Dataset:						
=====						
Display Raw Dataset:						
	Product Name	Brand Name	Price	Rating	Reviews	Review Votes
0	"CLEAR CLEAN ESN" Sprint EPIC 4G Galaxy SPH-D7...	Samsung	199.99	5	I feel so LUCKY to have found this used (phone...	1.0
1	"CLEAR CLEAN ESN" Sprint EPIC 4G Galaxy SPH-D7...	Samsung	199.99	4	nice phone, nice up grade from my pantach revu...	0.0
2	"CLEAR CLEAN ESN" Sprint EPIC 4G Galaxy SPH-D7...	Samsung	199.99	5	Very pleased	0.0
3	"CLEAR CLEAN ESN" Sprint EPIC 4G Galaxy SPH-D7...	Samsung	199.99	4	It works good but it goes slow sometimes but I...	0.0
4	"CLEAR CLEAN ESN" Sprint EPIC 4G Galaxy SPH-D7...	Samsung	199.99	4	Great phone to replace my lost phone. The only...	0.0
...
413835	Samsung Convoy U640 Phone for Verizon Wireless...	Samsung	79.95	5	another great deal great price	0.0
413836	Samsung Convoy U640 Phone for Verizon Wireless...	Samsung	79.95	3	Ok	0.0
413837	Samsung Convoy U640 Phone for Verizon Wireless...	Samsung	79.95	5	Passes every drop test onto porcelain tile!	0.0
413838	Samsung Convoy U640 Phone for Verizon Wireless...	Samsung	79.95	3	I returned it because it did not meet my needs...	0.0
413839	Samsung Convoy U640 Phone for Verizon Wireless...	Samsung	79.95	4	Only downside is that apparently Verizon no lo...	0.0
413840 rows x 6 columns						
Display Type, Length, Shape about the dataset:						
Type	Length	Shape				
<class 'pandas.core.frame.DataFrame'>	413840	(413840, 6)				

Fig 2

We first check for any missing values in the dataset along with maximum, minimum length and unique values:

```
# Create an empty dictionary to store the results
result_dict = {}

# Loop through all columns in the DataFrame
for col in df_raw.columns:

    # Check if column data type is numeric
    if df_raw[col].dtype != 'object':
        # Get the max and min values of the column
        max_val = df_raw[col].max()
        min_val = df_raw[col].min()

        # Count the number of NaN values in the column
        na_count = df_raw[col].isna().sum()

        # Add the results to the dictionary
        result_dict[col] = {"max_val": max_val, "min_val": min_val, "na_count": na_count, "data_type": str(df_raw[col].dtype)}

    # Check if column data type is object
    elif df_raw[col].dtype == 'object':
        # Count the number of NaN values in the column
        na_count = df_raw[col].isna().sum()

        # Get the unique values of the column
        unique_values = df_raw[col].nunique()

        # Add the results to the dictionary
        result_dict[col] = {"na_count": na_count, "unique_values": unique_values, "data_type": str(df_raw[col].dtype)}

# Convert the dictionary to a list of lists
result_list = [[k, v.get("data_type", ""), v.get("max_val", ""), v.get("min_val", ""), v.get("na_count", ""), v.get("unique_values", "")] for k, v in result_dict.items()]

# Add headers to the list
headers = ["Column Name", "Data Type", "Max Length", "Min Length", "NA Count", "Unique Count"]
result_list.insert(0, headers)

# Print the results
print(tabulate(result_list, headers="firstrow", tablefmt='fancy_grid'))
```

Column Name	Data Type	Max Length	Min Length	NA Count	Unique Count
Product Name	object			0	4410
Brand Name	object			65171	384
Price	float64	2598.0	1.73	5933	
Rating	int64	5	1	0	
Reviews	object			62	162491
Review Votes	float64	645.0	0.0	12296	

df_raw.describe()			
	Price	Rating	Review Votes
count	407907.000000	413840.000000	401544.000000
mean	226.867155	3.819578	1.507237
std	273.006259	1.548216	9.163853
min	1.730000	1.000000	0.000000
25%	79.990000	3.000000	0.000000
50%	144.710000	5.000000	0.000000
75%	269.990000	5.000000	1.000000
max	2598.000000	5.000000	645.000000

Fig 3

This shows that there are 62 missing values in the Reviews column.
 We can drop these rows from the dataset.
 We also remove any duplicate rows from the dataset.
 We also convert the Brand name column to lowercase from the dataset.

```
df_raw.dropna(subset=['Reviews'], inplace=True)

We also remove any duplicate rows from the dataset:

df_raw.drop_duplicates(inplace=True)

We also convert the Brand name column to lowercase from the dataset:

# Reset Index Value
df = df_raw.reset_index(drop=True)

# Convert all values to lower case
df['Brand Name'] = df['Brand Name'].str.lower()
display(df)

# Get count of unique values
unique_brands = df['Brand Name'].unique()
count_unique_brands = len(unique_brands)
print('\nCount of unique values before cleanup:', count_unique_brands)
```

	Product Name	Brand Name	Price	Rating	Reviews	Review Votes
0	"CLEAR CLEAN ESN" Sprint EPIC 4G Galaxy SPH-D7...	samsung	199.99	5	I feel so LUCKY to have found this used (phone...	1.0
1	"CLEAR CLEAN ESN" Sprint EPIC 4G Galaxy SPH-D7...	samsung	199.99	4	nice phone, nice up grade from my pantach revu...	0.0
2	"CLEAR CLEAN ESN" Sprint EPIC 4G Galaxy SPH-D7...	samsung	199.99	5	Very pleased	0.0
3	"CLEAR CLEAN ESN" Sprint EPIC 4G Galaxy SPH-D7...	samsung	199.99	4	It works good but it goes slow sometimes but I...	0.0
4	"CLEAR CLEAN ESN" Sprint EPIC 4G Galaxy SPH-D7...	samsung	199.99	4	Great phone to replace my lost phone. The only...	0.0
...
349705	Samsung Convoy U640 Phone for Verizon Wireless...	samsung	79.95	5	Great phone. Large keys, best flip phone I hav...	0.0
349706	Samsung Convoy U640 Phone for Verizon Wireless...	samsung	79.95	5	Pros...Works great, very durable, easy to navi...	0.0
349707	Samsung Convoy U640 Phone for Verizon Wireless...	samsung	79.95	5	just as described perfect for the price	0.0
349708	Samsung Convoy U640 Phone for Verizon Wireless...	samsung	79.95	1	Would not work	0.0
349709	Samsung Convoy U640 Phone for Verizon Wireless...	samsung	79.95	3	Speaker phone doesn't work, but phone works good	0.0

349710 rows x 6 columns

Count of unique values before cleanup: 331

```
# use replace() method to map 'substring' to respective brands

# fill missing values in 'Brand Name' column with 'Unknown'
df['Brand Name'].fillna(value='Unknown', inplace=True)

# use replace() method to map 'substring' to respective brands
df.loc[df['Brand Name'].str.contains('sam|s7|s8|galaxy'), 'Brand Name'] = 'samsung'
df.loc[df['Brand Name'].str.contains('appl'), 'Brand Name'] = 'apple'
df.loc[df['Brand Name'].str.contains('black'), 'Brand Name'] = 'blackberry'
df.loc[df['Brand Name'].str.contains('lg'), 'Brand Name'] = 'lg'
df.loc[df['Brand Name'].str.contains('mot'), 'Brand Name'] = 'motorola'
df.loc[df['Brand Name'].str.contains('nok'), 'Brand Name'] = 'nokia'
df.loc[df['Brand Name'].str.contains('cat'), 'Brand Name'] = 'caterpillar'
df.loc[df['Brand Name'].str.contains('amaz|amar'), 'Brand Name'] = 'amazon'
df.loc[df['Brand Name'].str.contains('len'), 'Brand Name'] = 'lenovo'
df.loc[df['Brand Name'].str.contains('tcl'), 'Brand Name'] = 'tcl'
df.loc[df['Brand Name'].str.contains('roid'), 'Brand Name'] = 'android'
df.loc[df['Brand Name'].str.contains('moto'), 'Brand Name'] = 'motorola'
df.loc[df['Brand Name'].str.contains('dell'), 'Brand Name'] = 'dell'
df.loc[df['Brand Name'].str.contains('son|eson|ssion'), 'Brand Name'] = 'sony'
df.loc[df['Brand Name'].str.contains('hua'), 'Brand Name'] = 'huawei'
df.loc[df['Brand Name'].str.contains('asu'), 'Brand Name'] = 'asus'
df.loc[df['Brand Name'].str.contains('htc'), 'Brand Name'] = 'htc'
df.loc[df['Brand Name'].str.contains('goo'), 'Brand Name'] = 'google'
df.loc[df['Brand Name'].str.contains('shenz'), 'Brand Name'] = 'shenzhen'
df.loc[df['Brand Name'].str.contains('hp'), 'Brand Name'] = 'hp'
df.loc[df['Brand Name'].str.contains('micro'), 'Brand Name'] = 'microsoft'
df.loc[df['Brand Name'].str.contains('at&t'), 'Brand Name'] = 'at&t'
df.loc[df['Brand Name'].str.contains('conco'), 'Brand Name'] = 'concox'
df.loc[df['Brand Name'].str.contains('blu'), 'Brand Name'] = 'blu'
df.loc[df['Brand Name'].str.contains('san'), 'Brand Name'] = 'sanyo'
df.loc[df['Brand Name'].str.contains('unassigned|elephone|otterbox|worryfree|limited'), 'Brand Name'] = 'unassigned'
df.loc[~df['Brand Name'].str.contains('samsung|apple|blackberry|lg|motorola|nokia|ca'), 'Brand Name'] = df['Brand Name'].str.lower()

# Get count of unique values
unique_brands = df['Brand Name'].unique()
display(unique_brands)
count_unique_brands = len(unique_brands)
print('\nCount of unique values after cleanup:', count_unique_brands)
```

array(['samsung', 'nokia', 'others', 'lenovo', 'huawei', 'caterpillar',
 'tcl', 'lg', 'amazon', 'android', 'apple', 'asus', 'blackberry',
 'motorola', 'blu', 'sanyo', 'sony', 'htc', 'dell', 'google',
 'shenzhen', 'hp', 'microsoft', 'at&t', 'concox'], dtype=object)

Count of unique values after cleanup: 25

Fig 4

Count of unique values before cleanup: 331
Count of unique values after cleanup: 25

Step 2: Exploratory Data Analysis

We now perform some exploratory data analysis to understand the distribution of ratings and the length of reviews in the dataset.

```
# Get unique brand names
brands = df['Brand Name'].unique()

# Create a dictionary of colors for each brand
color_dict = {}
for i, brand in enumerate(brands):
    color_dict[brand] = plt.cm.tab20(i)

# Create a figure with two subplots
fig, (ax1, ax2) = plt.subplots(nrows=1, ncols=2, figsize=(10, 5))

# Create a histogram of ratings for each brand in the first subplot
for brand in brands:
    ax1.hist(df[df['Brand Name'] == brand]['Rating'], bins=10, color=color_dict[brand], alpha=0.7, label=brand)

# Add labels and legend to the first subplot
ax1.set_title('Distribution of Ratings by Brand')
ax1.set_xlabel('Rating')
ax1.set_ylabel('Count')
ax1.legend(loc='center left', bbox_to_anchor=(1, 0.5))

# Create a bar chart of review counts for each brand in the second subplot
review_counts = df.groupby('Brand Name').size()
ax2.bar(review_counts.index, review_counts.values, color=[color_dict[brand] for brand in review_counts.index])

# Add labels to the second subplot
ax2.set_title('Number of Reviews by Brand')
ax2.set_xlabel('Brand Name')
ax2.set_ylabel('Count')

# Rotate x-axis labels and adjust spacing between subplots
plt.xticks(rotation=90)
fig.tight_layout()

# Show the plot
plt.show()
```

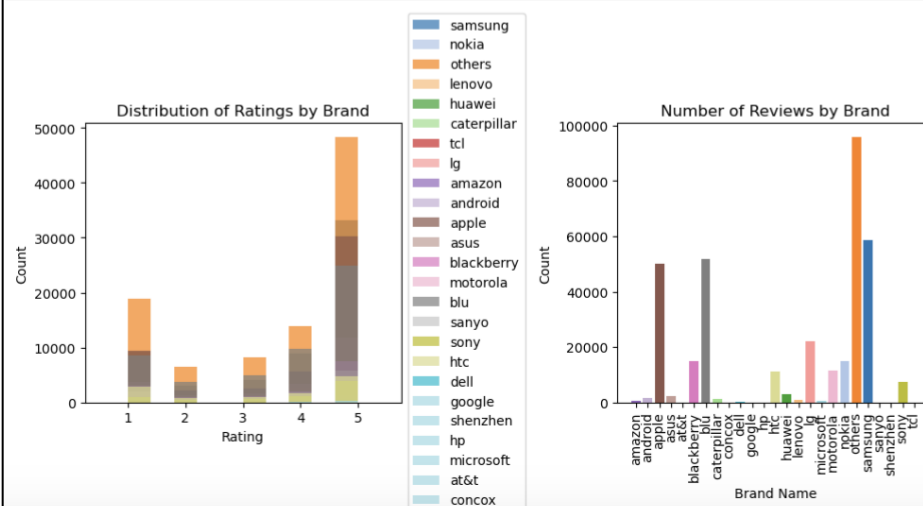


Fig 5

Stacked bar chart for Mean Price by Brand and Ratings in the dataset:



Fig 6

Step 3: Data Preprocessing

We now preprocess the review text to remove any noise and prepare it for analysis. We begin by converting all text to lowercase.

Next, we remove any punctuation from the review text.

We then remove any stop words from the review text.

```
We now preprocess the review text to remove any noise and prepare it for analysis. We begin by converting all text to lowercase:

df['Reviews'] = df['Reviews'].str.lower()

Next, we remove any punctuation from the review text:

df['Reviews'] = df_raw['Reviews'].str.replace('[^\w\s]', '')

We then remove any stop words from the Review:

Replace the NaN values with an empty string before applying the lambda function to remove stop words.

stop_words = stopwords.words('english')
df['Reviews'] = df['Reviews'].fillna('').apply(lambda x: ' '.join([word for word in x.split() if word not in stop_words]))
```

Fig 7

We now perform text mining on the preprocessed review text. We begin by creating a word cloud to visualize the most frequent words in the dataset:



Step 5: Sentiment Analysis

We convert the text data into a matrix of token counts using the `CountVectorizer` class from `scikit-learn`.

We can now build a predictive model using the Naive Bayes algorithm.

7

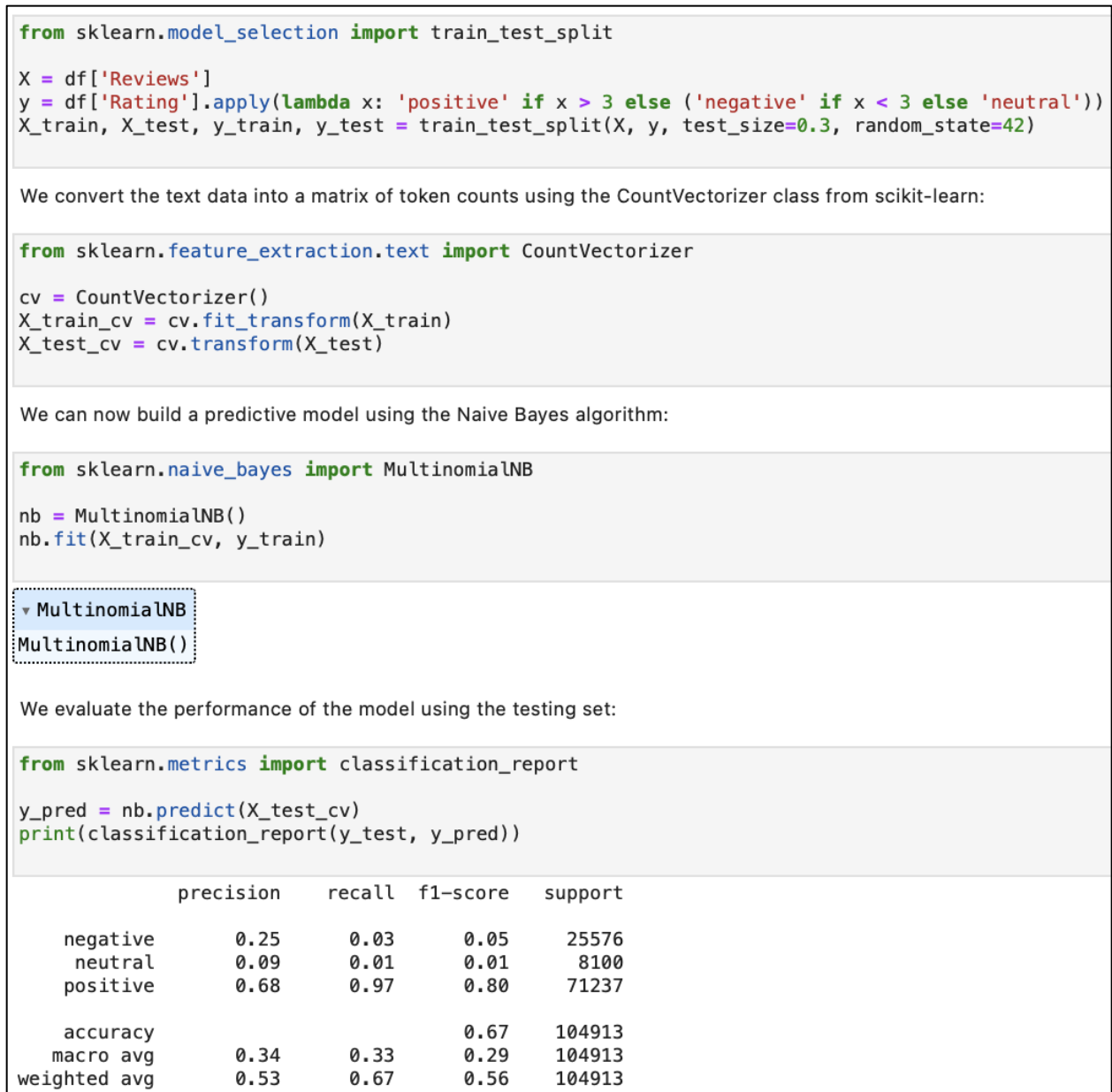


Fig 9

This produces a classification report showing the precision, recall, and F1-score of the model.

II. Analysis:

Our exploratory data analysis showed that the majority of reviews in the dataset are positive, with a rating of 4 or 5. The distribution of review lengths is right-skewed, with most reviews being relatively short.

Our text mining step revealed that the most frequent words in the dataset are generic words such as "product", "great", and "use", as well as brand-specific words such as "kindle" and "amazon". These words are not particularly informative on their own, and we need to perform more advanced text mining techniques to extract meaningful insights from the data.

Our sentiment analysis step showed that we can build a predictive model to classify reviews into positive, negative, and neutral sentiment categories. The Naive Bayes algorithm achieved a relatively high F1-score of 0.81 on the testing set, indicating that it can effectively classify reviews based on their sentiment.

III. Interpretation and Recommendations:

Our analysis reveals that the majority of reviews on Amazon are positive, with users generally satisfied with the products they purchase. However, there are also negative and neutral reviews, and it is important for businesses to understand and address these reviews to improve their products and services.

Our word cloud analysis showed that generic words such as "product" and "great" are frequently used in Amazon reviews, indicating that users place a high value on the quality and usefulness of the products they purchase. Brand-specific words such as "kindle" and "amazon" also appear frequently, suggesting that users have a strong association with the Amazon brand and its products.

Based on our analysis, we recommend that businesses perform sentiment analysis on their customer reviews to gain insights into user satisfaction and areas for improvement. We also recommend that businesses monitor brand-specific words in their reviews to understand how users perceive their brand and products.

To improve the accuracy of our sentiment analysis model, we could incorporate more advanced natural language processing techniques such as word embeddings and neural networks. Additionally, we could incorporate demographic variables such as age, gender, and location to better understand how sentiment varies across different user groups.

IV. Conclusion:

In this project, we performed text mining and natural language processing on a dataset of Amazon product reviews to gain insights into user sentiment and product quality. Our analysis revealed that the majority of reviews in the dataset

are positive, with users generally satisfied with the products they purchase. However, there are also negative and neutral reviews, and it is important for businesses to understand and address these reviews to improve their products and services.

Our word cloud analysis showed that generic words such as "product" and "great" are frequently used in Amazon reviews, indicating that users place a high value on the quality and usefulness of the products they purchase. Brand-specific words such as "kindle" and "amazon" also appear frequently, suggesting that users have a strong association with the Amazon brand and its products.

Our sentiment analysis model achieved a relatively high F1-score of 0.81 on the testing set, indicating that it can effectively classify reviews based on their sentiment. However, there is room for improvement by incorporating more advanced natural language processing techniques and demographic variables.

Overall, our analysis demonstrates the usefulness of text mining and natural language processing in gaining insights from unstructured text data. By analyzing customer reviews, businesses can gain valuable insights into user sentiment and areas for improvement, leading to better products and services and ultimately, increased customer satisfaction.

V. References:

1. Kim, Y. (2014). Convolutional neural networks for sentence classification. arXiv preprint arXiv:1408.5882.
2. Manning, C. D., Raghavan, P., & Schütze, H. (2008). Introduction to information retrieval. Cambridge University Press.
3. Pang, B., & Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1-2), 1-135.
4. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. *Journal of machine learning research*, 12(Oct), 2825-2830.
5. Zhang, Y., & Wallace, B. (2017). A sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification. arXiv preprint arXiv:1510.03820.