



# WEB APIs & CLASSIFICATION

...

PROJECT 3

# PROBLEM STATEMENT

To **classify** Reddit posts from r/nosleep and r/thetruthishere using Natural Language Processing (**NLP**) and **Classification** modelling. The model evaluated with **Accuracy** scores is then expected to predict to which subreddit a given post belongs to.

The model should then help Reddit data science team to advise their advertisers on spending forecast for targeted marketing campaigns of products and services like online gaming targeting members of r/nosleep, calming supplements & therapy/ counselling services for members of r/thetruthishere depending on the predicted subreddit.

# BACKGROUND

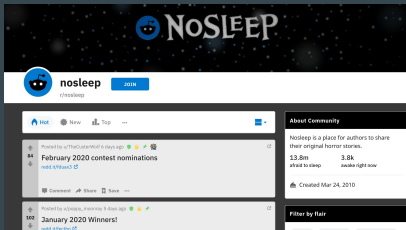
The objective is to find out the attributes of each group and their online behaviour.

With the information gathered, business can;

- Decide how much advertising dollars to spend and on which group to obtain better ROI
- Types of products and services to be advertised
- Frequency of advertisement
- Type of advertisement eg, static, pop-up, etc
- KPI expected from Reddit

# DATA COLLECTION

## HTML Data (Reddit's API)



## Collect Data (Python)

# Python's Requests Library



## Save Data (JSON)

## Parse Data (Python)

URL:

"https://www.reddit.com/r/nosleep.json"

\* 25 posts per request

```
headers = {'User-agent': 'Pony Inc 1.0'}
response = requests.get(url, headers
= headers)
Add timer: time.sleep()
```

## response.json() JavaScript Object Notation

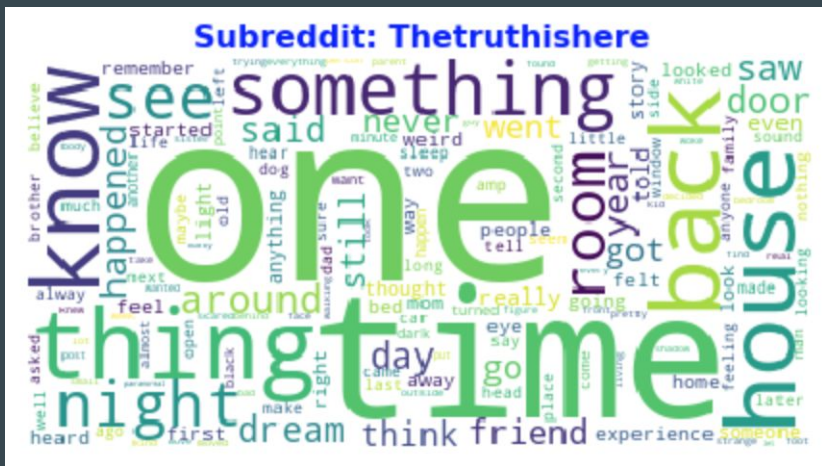
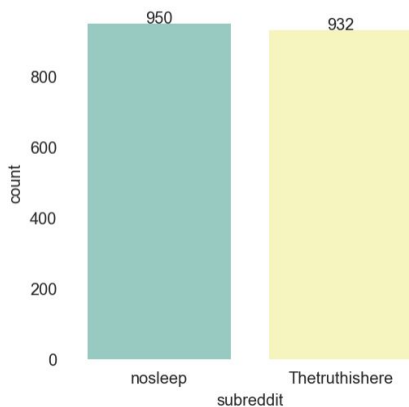
[illegible]

subreddit	selftext	author	fullname	saved	mod_reason_line	gilded	clicked	title	link_hair	content	subreddit	...
Theriotishness	is planning to kill himself	12_4qgdtg	False		N/A	0	False	Demonic sounds	█		r7H	...
Theriotishness	Does anyone else have this thing ...	12_4w3hdms	False		N/A	0	False	Anyone else?	█		r7H	...
Theriotishness	H there I thought you got the subreddit ...	12_3fH2t2	False		N/A	0	False	Demonic type of strains and other strange stuff...	█		r7H	...

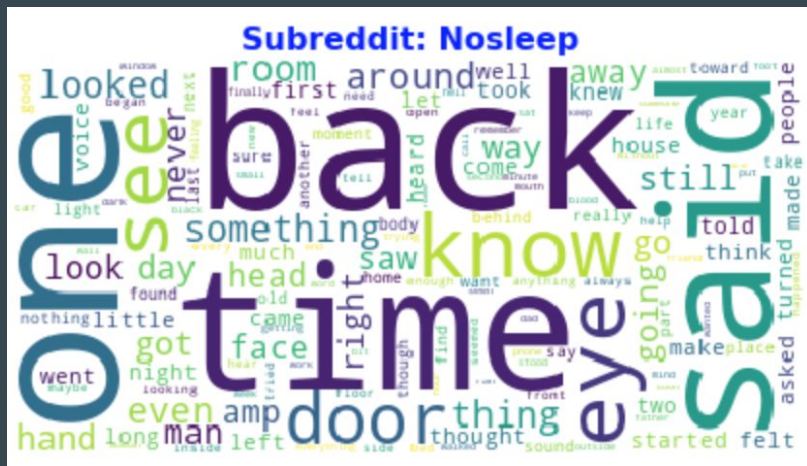
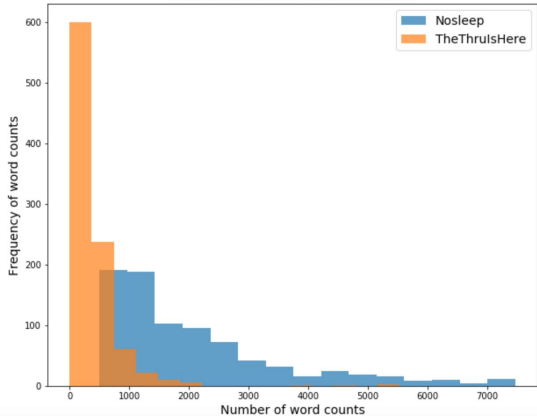


# EDA

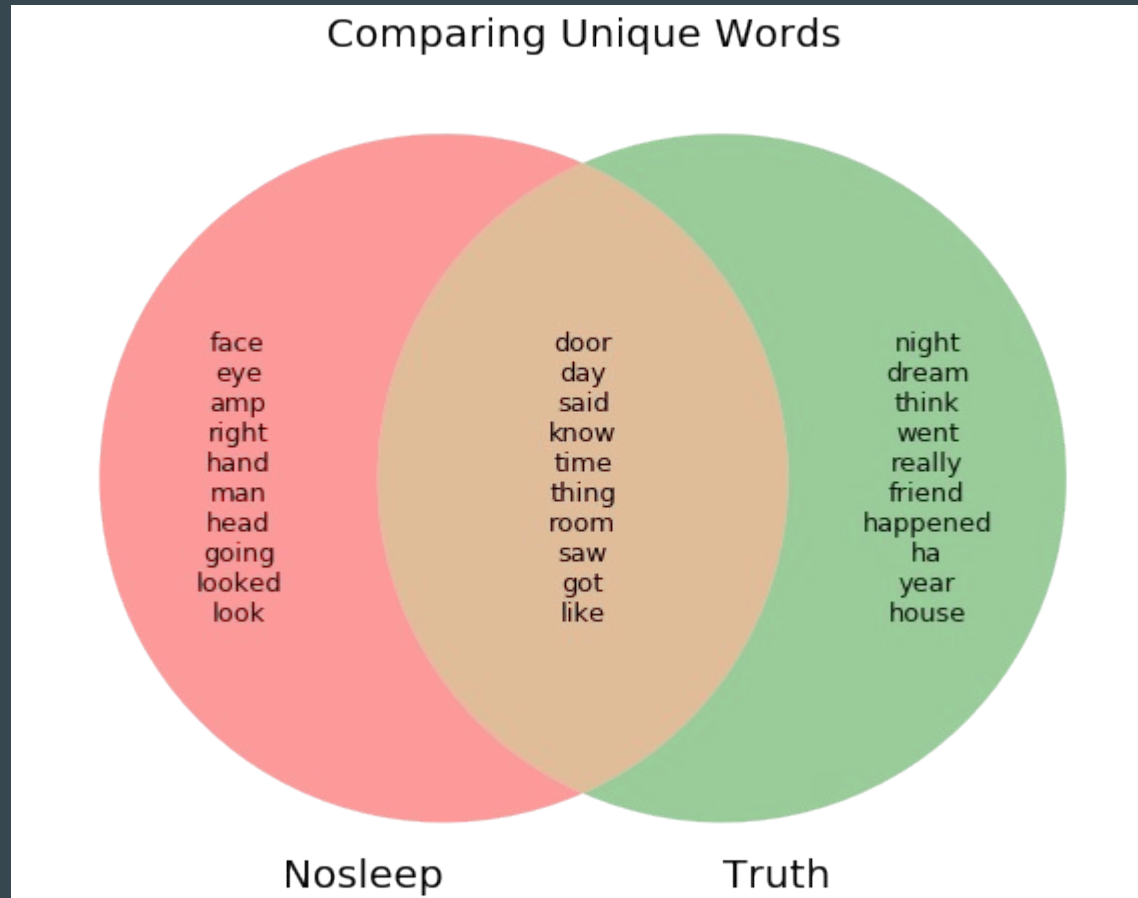
### Number of unique posts per subreddit



Word counts for Post's content



# PREPROCESSING AND MODELLING



# PREPROCESSING AND MODELLING

- Split our data into X and y

```
x = df_combined_clean[['combined']]  
y = df_combined_clean['subreddit']
```

- Split our data into training and testing sets
- Turn our text into features
  - CountVectorizer
  - TF-IDFVectorizer



# PREPROCESSING AND MODELLING

## CountVectorizer

	able	actually	ago	alys	amp	arm	asked	ay	bad	bed	...	week	weird	went	white	window	woman	word	work	world	year
0	0	1	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	1	0
2	0	0	0	2	0	1	2	4	0	0	...	0	0	1	1	0	5	0	2	0	1
3	0	1	4	0	53	7	6	12	4	13	...	1	0	7	2	2	8	2	1	0	15
4	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0

## TF-IDFVectorizer

	able	actually	ago	alys	amp	arm	asked	ay	bad	bed	...	week	weird	went	white	window	woman
0	0.0	0.349599	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.0	0.000000	0.000000	0.000000	0.000000
1	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.0	0.000000	0.000000	0.000000	0.000000
2	0.0	0.000000	0.000000	0.080697	0.000000	0.049083	0.082776	0.141408	0.000000	0.000000	...	0.000000	0.0	0.033910	0.052101	0.000000	0.275428
3	0.0	0.007470	0.026131	0.000000	0.471109	0.051864	0.037485	0.064037	0.030006	0.084334	...	0.007056	0.0	0.035831	0.015729	0.01404	0.066522
4	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.0	0.000000	0.000000	0.000000	0.000000

# PREPROCESSING AND MODELLING

Classifier Models	
Naive Bayes Model (MultinomialNB)	Logistic Regression
Columns are all integer counts	Response default falls into one of two categories.

# EVALUATION

## BASELINE SCORE

	0	1
target	0.504782	0.495218

## ACCURACY

Scores	CVEC + LR	TDIF + LR
Training score	0.999	0.998
Testing score	0.953	0.958
Variance	-0.046	-0.040

## CONFUSION MATRIX

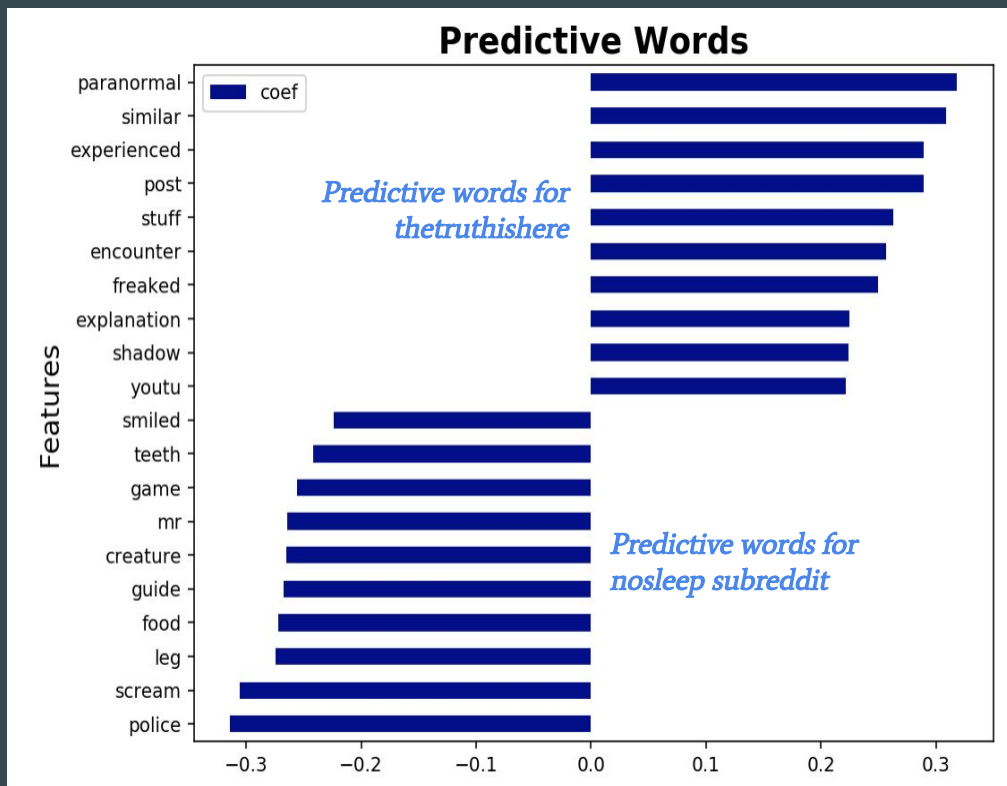
LR + CVEC	Predicted Nosleep	Predicted Thetruthishere
Actual Nosleep	228	10
Actual Thetruthishere	12	221

LR + TDIF	Predicted Nosleep	Predicted Thetruthishere
Actual Nosleep	234	4
Actual Thetruthishere	16	217

# EVALUATION

- 1) With logistic regression, we can identify predictive words that belong to each subreddit.
- 2) This will help in forecasting spending for the different products/ services that are targeting members of each subreddit.

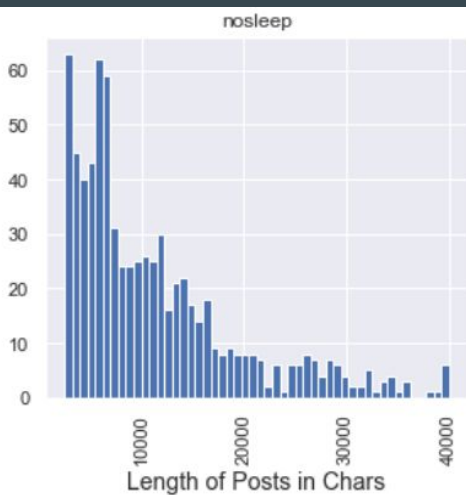
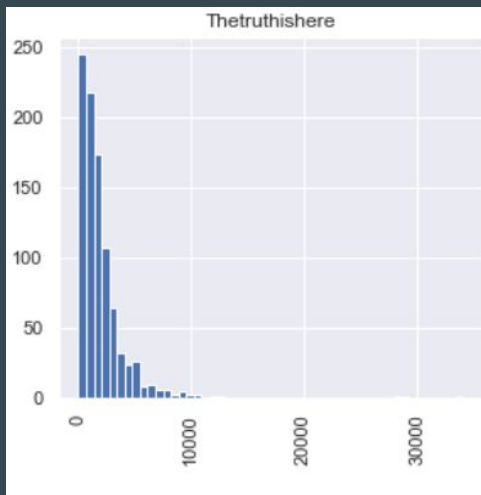


# LIMITATIONS OF THE MODEL

- 1) The hyper-parameters of the model has been tuned to include features/ words as many as 1500 to reduce the variance between testing and training scores.
- 2) But in reality, it may make sense to reduce the number of words.
- 3) The model can be further enhanced :
  - a) **To eliminate words of similar meanings**
  - b) **To include phrases instead of just words for prediction.**

# CONCLUSION

- Business should get data directly from Reddit using API or RESTful API which is significantly faster than scraping (limit on posts).
- Business may want to spend more advertising dollars on the Nosleep group as higher chance to convert users to customers; more bang for the buck



**About Community** ...

Nosleep is a place for authors to share their original horror stories.

**13.8m**  
afraid to sleep

**3.0k**  
awake right now

---

Created Mar 24, 2010

**CREATE POST**

---

**COMMUNITY OPTIONS** v

# RECOMMENDATIONS

- Offer products and services such as games, online gaming, movie subscriptions, and food delivery services
- Advertise counselling services, calming supplements to promote good night sleep to the Truth group
- Business can consider using our Topic modeling service to segment users into groups by say, sci-fi, horror, fantasy, etc for;
  - pop-up advertisement according to genre, and
  - Implement adaptive or real-time advertisements at Reddit
- Business can work with Reddit on advertisement placements, get Reddit to provide analytics such as 'click-rate'



5.0k

Posted by u/RichardSaxon 16 hours ago 

## 42 years ago we sent Voyager 1 into space to look for extraterrestrial life, today we found it.

42 years, 6 months and 3 days ago, on the 5th of September 1977, a space probe was launched from Earth, and sent on an endless journey through space. The probe, which was affectionately named Voyager 1, contains a multitude of information regarding humanity, including our language, our art, and in a more metaphorical way; Our souls.

Today, the little machine is 22 billion kilometers away from Earth, the furthest reach of our species, though not a manned vessel, it's still a part of us. Despite its distance, we're still keeping contact with it, and during its journey spanning almost half a century, it has given us an insight into the mysteries of the universe we thought we could only dream about.

My own father spent most of his life on the project, and I have since followed in his footsteps. I dreamed about taking over his work even as a twelve year old kid, and

 203 Comments  Share  Save ...

ADVERTISEMENT



 **GENERAL ASSEMBLY**

**Equip your team with  
the digital skills it needs  
to stay competitive**

