

The title is framed by a large dashed rectangle. At the top-right corner, a dashed line extends horizontally to the right and then curves downwards as an arrow. At the bottom-right corner, a dashed line extends horizontally to the right and then curves upwards as an arrow. At the bottom-left corner, a dashed line extends horizontally to the left and then curves upwards as an arrow.

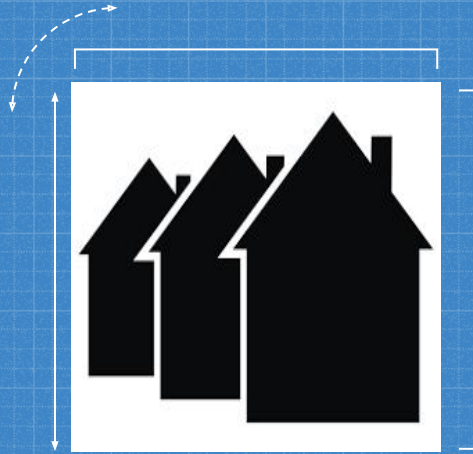
Project 3 - Subreddits

My Name is Cam Smugereski !

I am the Marketing
Analyst for Pet Supplies
Plus in Bedford, NH.

My goal is help you
better market dog or cat
specific to the right
customers.





Problem Statement

This is a binary classification problem asking if we can predict the Subreddit, cats or dogs, based on selftext or the title, or both.

The ideal outcome is to provide the marketing team with top 'Buzz Words' they can use to directly appeal to either dog or cat lovers. We also aim to be able to identify cat or dog owners based on the way they post on Reddit.

Let's Start With Our Most Notable Discovery:



Missing Values

Dogs

```
subreddit      0  
selftext       8  
title          0  
created_utc    0  
dtype: int64
```

Cats

```
subreddit      0  
selftext      820  
title          0  
created_utc    0  
dtype: int64
```




82% of our 'cats' Subreddit has missing values in
the selftext column

Could this be an indicator of
who is a cat or dog person?

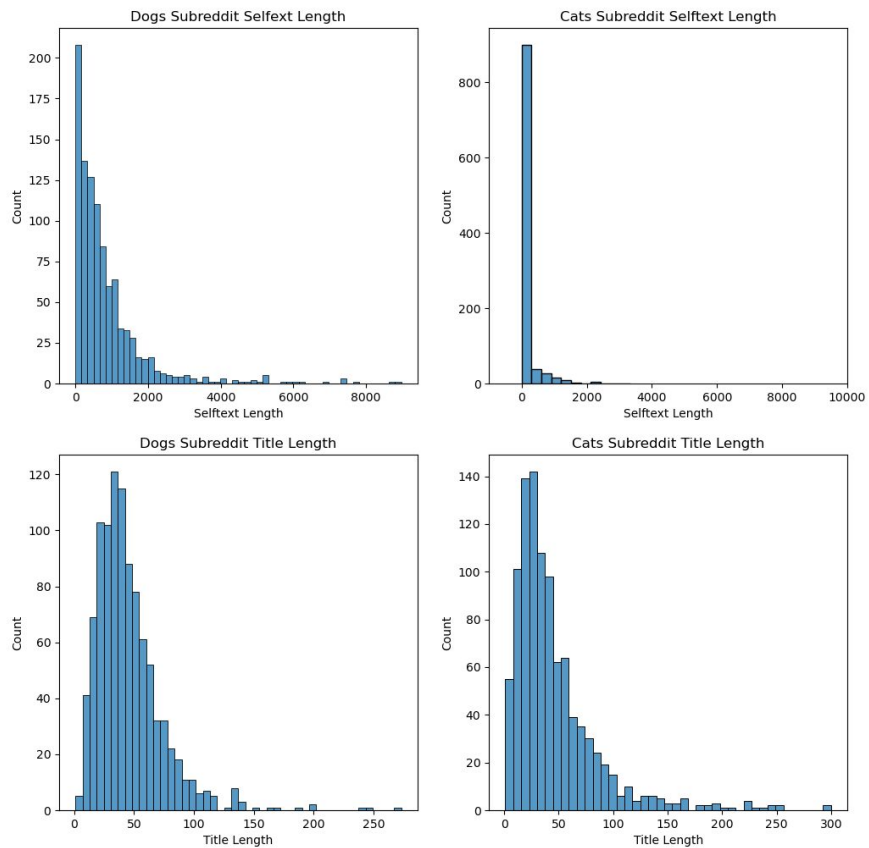




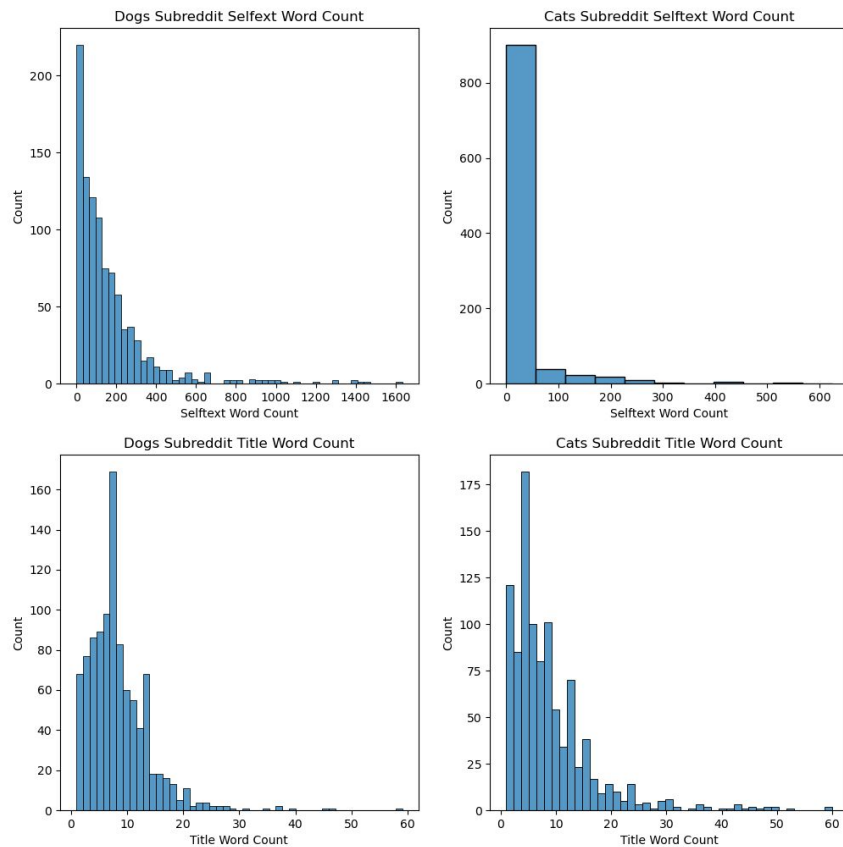
What Are the Posts Telling Us?

Let's look at some
simple stats!

Selftext and Title Lengths Based on Subreddit



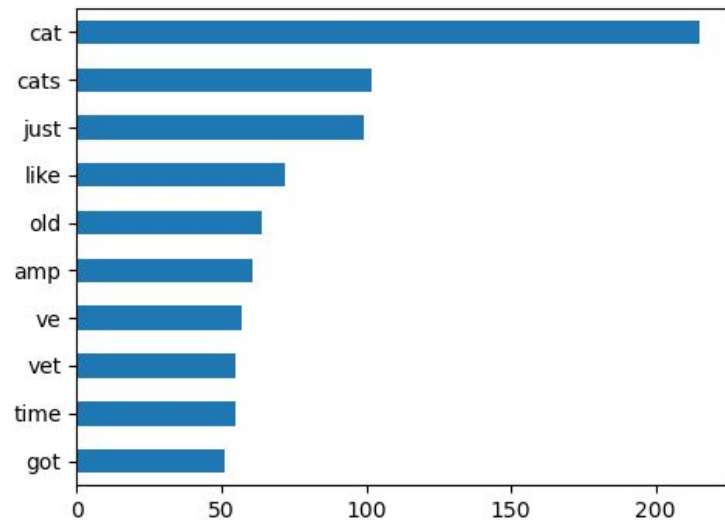
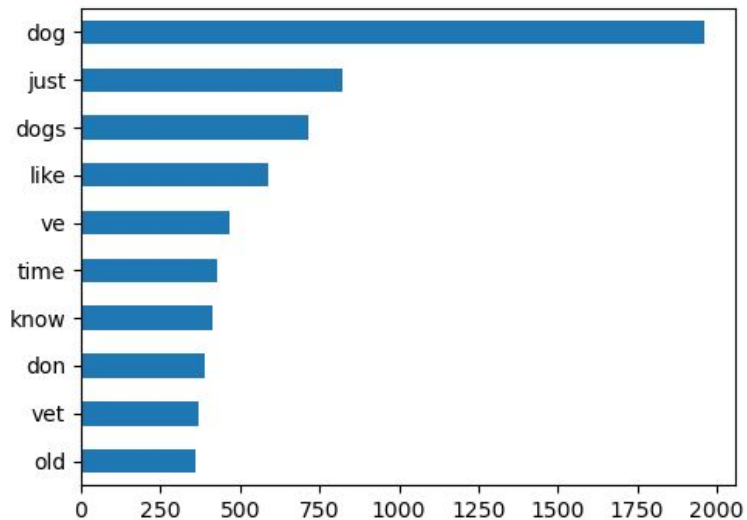
Selftext and Title Word Counts Based on Subreddit



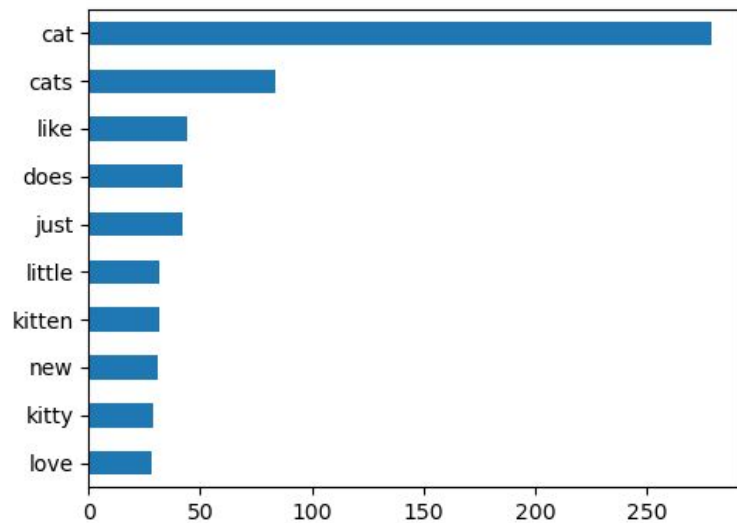
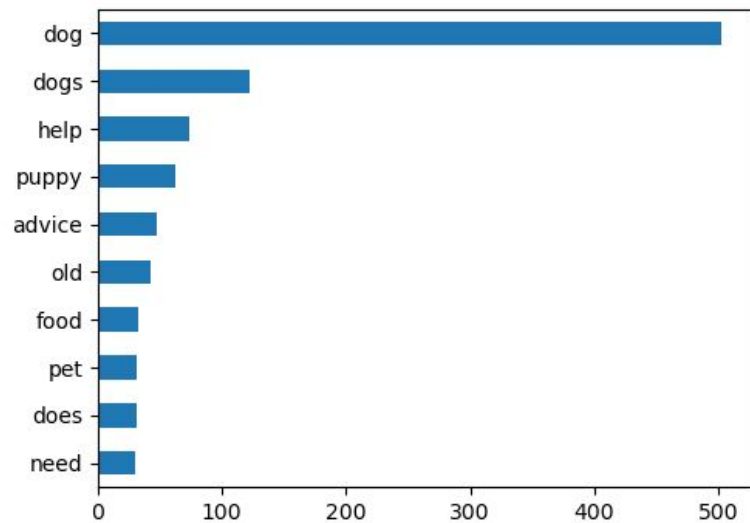


Top Words

Selftext



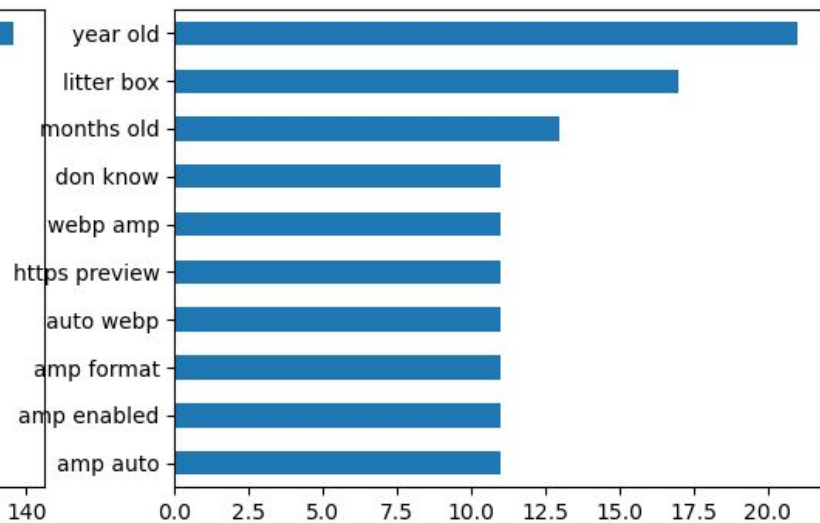
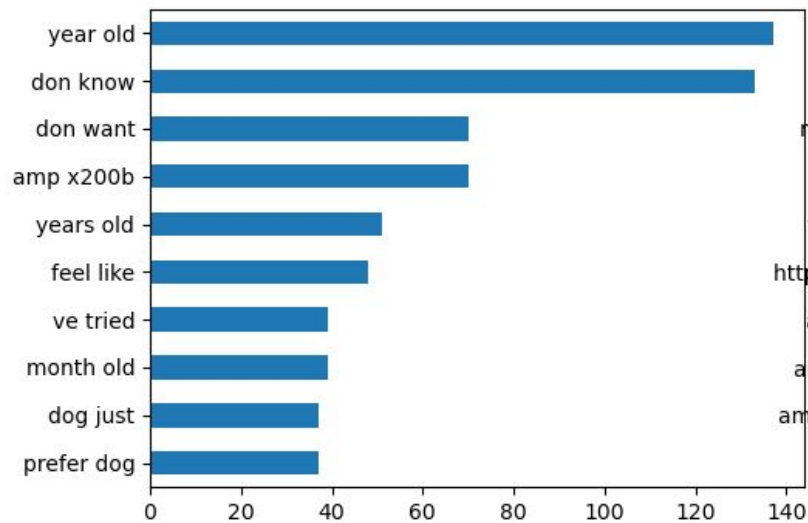
Title



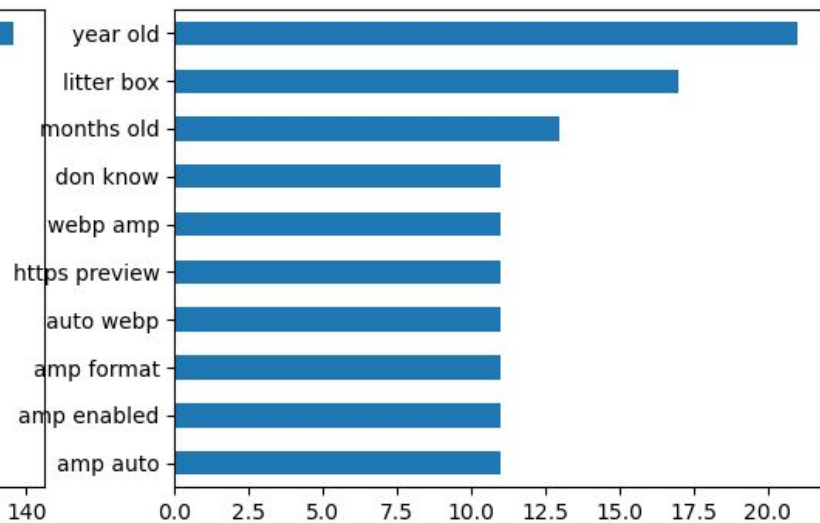
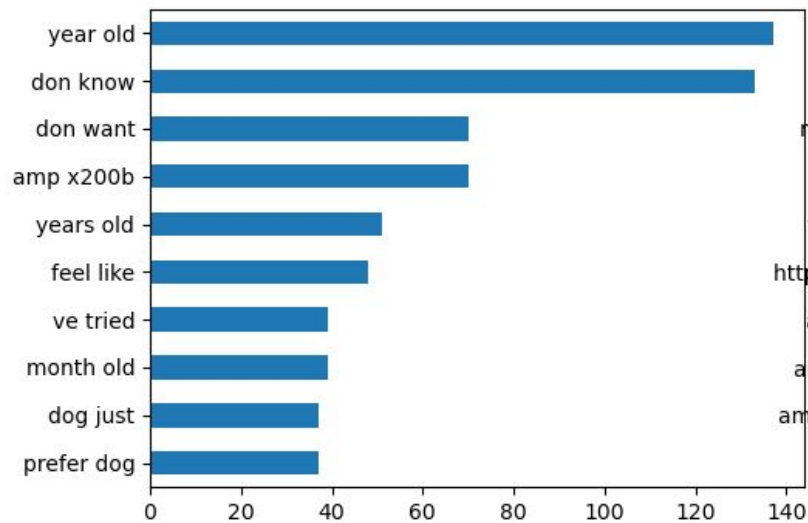


Top Bigrams

Selftext

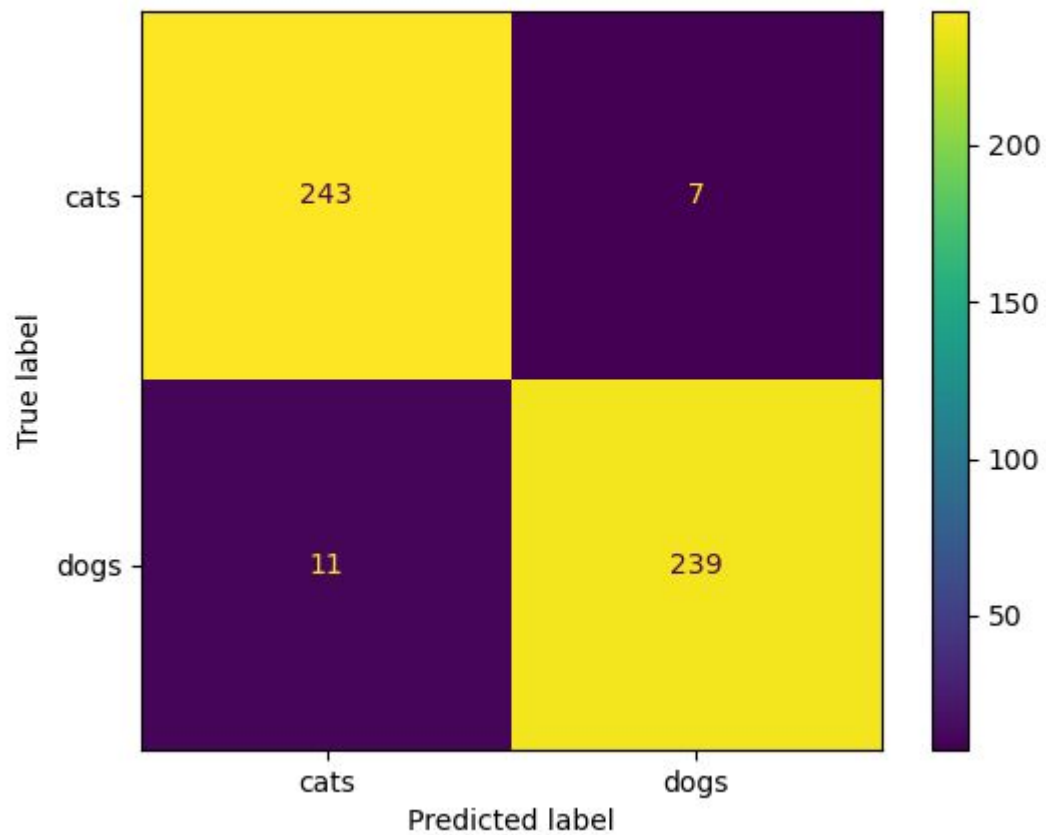


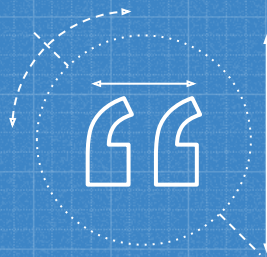
Title



Stacking

- Logistic Regression
- Naive Bayes
- Random Forest





Training Score: 99.3%

Test Score: 96.4%

Sensitivity: 95.6%



Sensitivity is our true positive rate and therefore tells us that 95% of our data that we predicted to be true, was actually true in our stacked model.

Conclusions and Recommendations



Takeaways

1. 82% of our 'cats' selftext is missing - What does this tell us?
2. I recommend using visuals for marketing cat products
3. Marketing for dog products can and should include some 'Buzz Words'
4. Our model and sensitivity scores are high, but this is likely due to how different each Subreddit is



Any Questions?