

THIS IS A JOKE

An Exploration of Natural Language Processing

By Molly Baird

PROBLEM STATEMENT

Can we use natural language processing to distinguish between two joke-themed subreddits?

- /r/Jokes
- /r/DadJokes
- /r/MommaJokes
- /r/CleanJokes
- /r/DirtyJokes



JOKES VS DADJOKES

I'VE RUN OUT OF TOILET PAPER AND STARTED USING
OLD NEWSPAPERS INSTEAD...

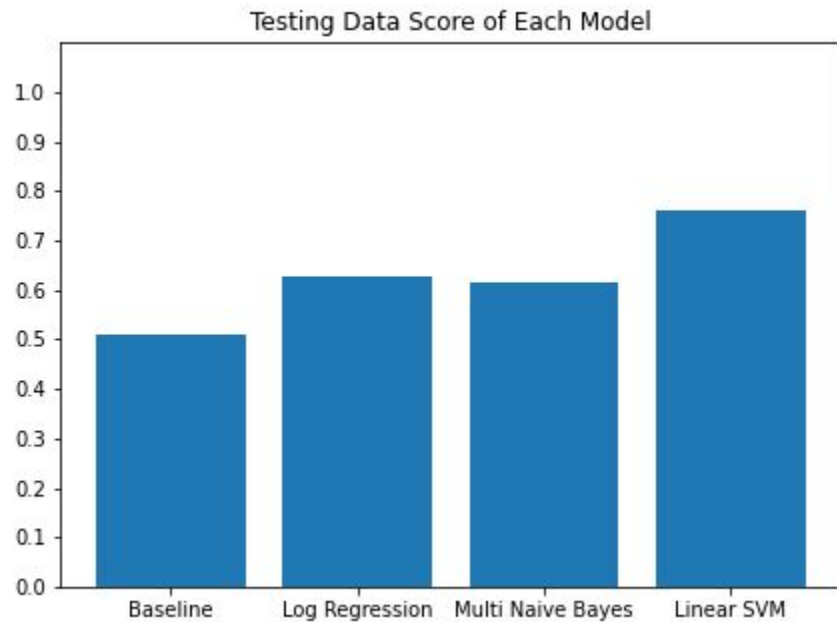
I'VE RUN OUT OF TOILET PAPER AND STARTED USING
OLD NEWSPAPERS INSTEAD...

...THE TIMES ARE ROUGH

Posted on /r/DadJokes by u/HellsJuggernaut

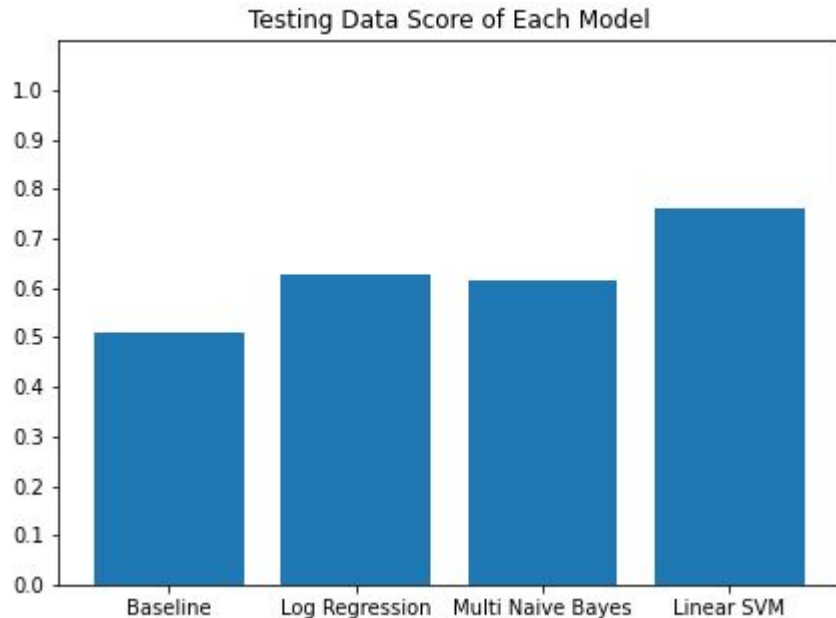
MODEL PERFORMANCE

- Size of dataset: 9,126
- Submission title and selftext only
- No stopword removal
- N-gram range (1,1)
- Max features 500
- Lemmatizing/Stemming had little to no effect
- CountVectorize vs Tfidf had little to no effect
- Linear Support Vector Machine 76% accurate



MODEL PERFORMANCE

- Size of dataset: 9,126
- Submission title and selftext only
- No stopword removal
- N-gram range (1,1)
- Max features 500
- Lemmatizing/Stemming had little to no effect
- CountVectorize vs Tfidf had little to no effect
- Linear Support Vector Machine 76% accurate



WHY ISN'T THIS MODEL BETTER?

MY FRIEND AND I JUST STARTED A BUSINESS
WHERE WE WEIGH TINY OBJECTS...

...IT'S A SMALL SCALE OPERATION

POSTED ON /R/JOKES BY U/PORICHOYGUPTO

WHAT'S THE DIFFERENCE BETWEEN JUVENILE HUMOR
AND DAD JOKES?

A DAD JOKE IS FULL GROAN

POSTED ON /R/JOKES BY U/BEENTHEREONCE2

JOKES VS MOMMAJOKES

YO MOMMA SO FAT...

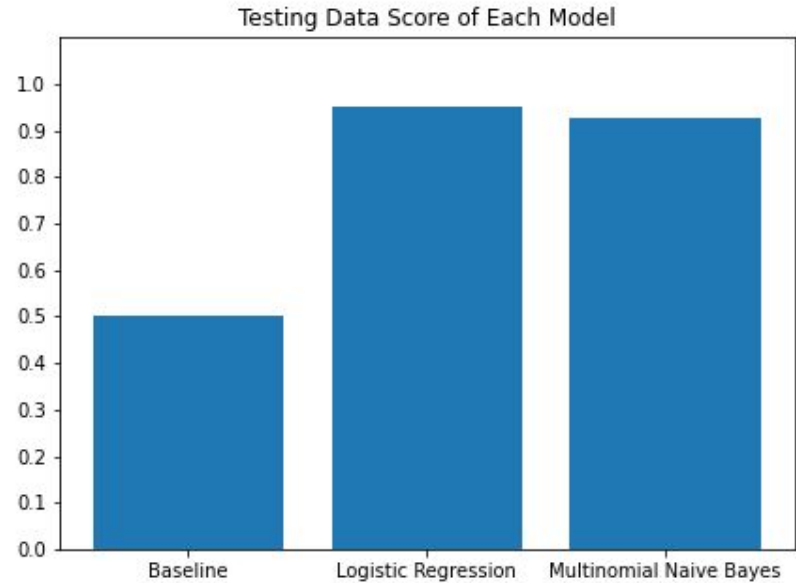
YO MOMMA SO FAT...

...SHE GOT COUNTED TWICE IN THE CENSUS

Posted on /r/MommaJokes by u/luxdesigns

MODEL PERFORMANCE

- Size of dataset: 2,744
- Submission title and selftext only
- No stopword removal
- N-gram range (1,2)
- Max features 500



CLEANJOKES VS DIRTYJOKES

WHY DO PEOPLE IN ATHENS HATE GETTING UP SO EARLY?



WHY DO PEOPLE IN ATHENS HATE GETTING UP SO EARLY?

BECAUSE DAWN IS TOUGH ON GREECE.

Posted on /r/CleanJokes by u/Crazycrafter97531

WHAT KIND OF BEES MAKE MILK?

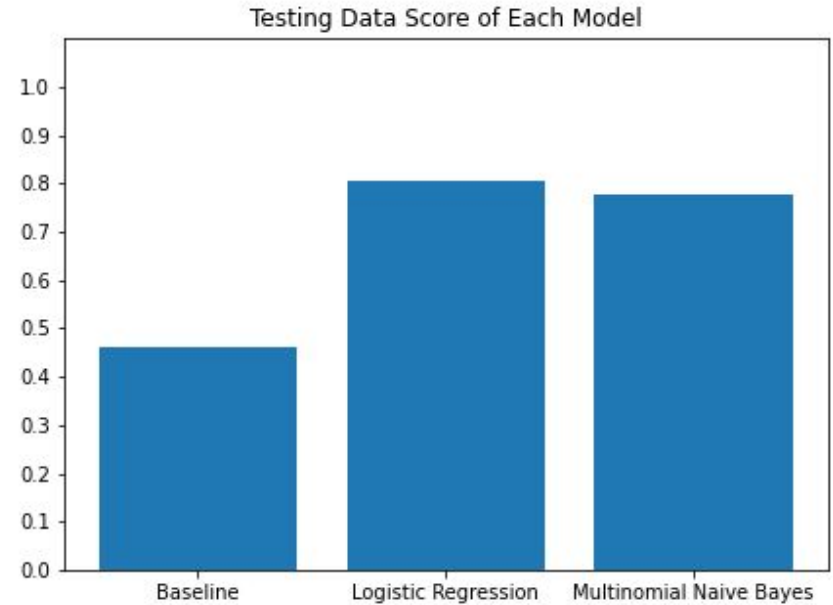
WHAT KIND OF BEES MAKE MILK?

BOOBIES.

Posted on /r/DirtyJokes by u/TheDisappointment101

MODEL PERFORMANCE

- Size of dataset: 9,276
- Submission title and selftext only
- English stopword removal
- N-gram range (1,1)
- Max features 500



CONCLUSIONS AND CONTINUATIONS

- Jokes with standard phrasing are easier to classify
- Dad joke is not well defined
- Stemming/lemmatizing didn't really help much
- Want to incorporate comments
- Bigger datasets?
- Want to know which jokes were misclassified
- Want to know which words got the most weight
- Multiclass model with all joke subreddits

THANK YOU!