# Reddit Classification

John DSI 22

# Introducing subreddits

| | Similarities | Difference |
|---|---|---|
| r/ValueInvesting | Both investment forum | • Focuses fundamental analysis<br>• Looks up to Warren Buffet |
| r/wallstreetbets | Common topics: Investment questions, Investment pitch | • "Social momentum" investing 🚀<br>• Looks up to Elon Musk |

# How Michael Burry figured out the 2007 crash, simple (own repost from Burryology)

Books

I have been reading the book: The oil factor by Stephen Leeb written in 2004. He talks about the inverse relation between (rapid) increase in oil prices, lowering supply and high demand, but he takes a detour. The dotcom bubble dropped sp500 -40%, nasdaq -80%, 16trillion USD wealth went to 7 trillion. The fed lowered rates to 0.75%, boosted borrowing and home prices served as a healthy collateral, which can only go up right? US was highly in debt before the bust, but after... oh with low rates causing booms in home prices, more debt. In this 2004 books he says, if home prices would fall it would be taking down the banking system (1:6 leverage at that time so 18% default was needed to make the banks insolvent, we know later the leverage was 1:20 so 5% default was enough). What would cause home prices to fall? Policies to curb inflation, aaaand when did the fed start to raise rates? Yes, early 2007. No more cheap refinancing causing defaults (subprime etc), and booooom.

Amazing book btw on oil, I would recommend it :) thought I would share my joy of finding this out, maybe Burry read this book also in 2004?

# 13,900 Shares ASO YOLO

**YOLO**

| Actions | Symbol | Name | Quality@Price | Market Cap | Available QTY | Cost | % Chg | |
|---------|--------|------|---------------|------------|---------------|------|-------|---|
| Trade | ASO | Academy Sports and Outdoors, Inc. | 8,400@42.41 | 356,244.00 | 8,400 | 35.848 | +18.30% | 74 |

# Ape interested in IPO's specifically DNUT (Krispy Kreme)

**Discussion**

Hey gang, I have not gotten in on any IPO's yet but browsing the upcoming IPO calendar on NASDAQ looks like good ole KRISPY KREME is on the list for the first of July!

I am blinded by my love and nostalgic memories of Krispy Kreme growing up in NC so I was wondering if any of you have experience with IPO's on opening day?

Seems like well known, and tech companies usually do well but they can also take a dive from the paper hands making a quick buck.

Bullish or Bearish on this one? or any other upcoming IPO suggestions? This is not a DD or financial advice

# Use case?

# Model workflow

| Data Collection | Data Cleaning | Pre-processing | Modeling |

**From Reddit API**
1070 r/valueinvesting post
1789 r/wallstreetbets post

**Drop null 'selftext'**
**Drop word length <20**
705 r/valueinvesting
825 r/wallstreetbets

**Steps taken**
- Remove URLs
- Remove symbols
- Convert to lowercase
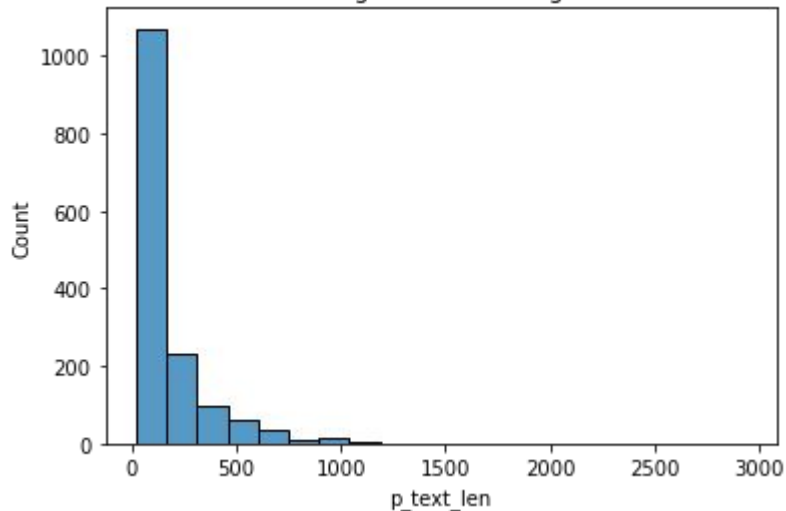- Remove stopwords
- Lemmatize

**Vectorizer:**
CountVectorizer
TfidfVectorizer

**Model:**
Logistic Regression
SVM
Naive Bayes

# Exploratory Data Analysis - Text Length


Histogram of text length

Mean word count:

r/valueinvesting: 103

r/wallstreetbets: 226

# Exploratory Data Analysis - Top 10 words

**Common top 10 words**

Company, Stock, Market, Year, Price, Like, Share

# Model Comparison

| Vectorizer | Model | Accuracy* |
|------------|-------------|-----------|
| Count | Logistic | 87.9% |
| TFIDF | Logistic | 90.3% |
| Count | SVM | 87.4% |
| TFIDF | SVM | 90.9% |
| Count | Naive Bayes | 90.3% |
| TFIDF | Naive Bayes | 89.5% |

*Baseline score: 54%

# Hyperparameters

C = 100

Kernel = rbf

Max_df = 0.8

Min_df = 2

Ngram_range = (1,2)

# Model Evaluation

**Overfitting**

Train score > Test score for all model

**Suboptimal hyperparameters**

Reduce hyperparameters in GridSearch due to processing time

**SVM blackbox**

Blackbox model means cannot get insights, to intuitively understand what features matters

# Side note: Logistic Regression Coefficient

A few notable words with high coefficients

**r/valueinvesting**

Ratio, operating, cash, performance, market cap

**r/wallstreetbets**

CLOV, short, meme, china, tl dr

# Thank you