

Like 4chan found the Shannon Entropy: An information theoretic analysis of online discourse

Colton Cunov
(Dated: March 9, 2021)

I. INTRODUCTION

Reddit's r/WallStreetBets is an online forum dedicated to discussing stock and options trading and noted for its role in the meteoric rise of GameStop Corp. stock price in early 2021, which resulted in short position losses to US firms of \$70 billion [1]. The forum is associated with a unique culture of colorful language, self-deprecating humor, and pride in extreme capital gains and losses. At the start of the year it consisted of about two million users, exploding to over nine million due to media coverage of the so-called short squeeze.

To examine the effects of this drastic change in user base from an information theoretic perspective, I calculated the correlation information based on post comments for various correlation lengths before, during, and after the rise in user base. My hypothesis was that a drastic increase in the number of participants fueled by a single event would result in less information at various correlation lengths. Additionally, I generated sample comments with remarkable integrity.

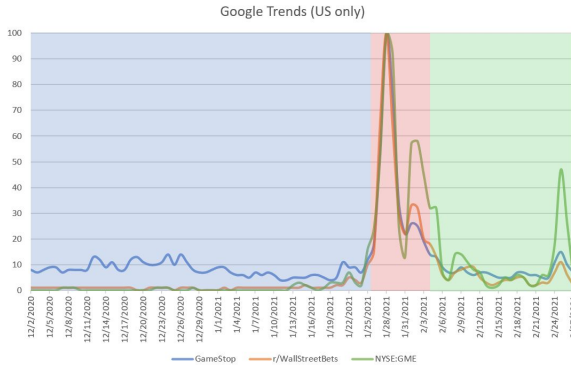


FIG. 1. Google Trends results for the selected keywords, separated into intervals defined by January 25, 2021 and February 5, 2021.

II. DATA COLLECTION AND PRE-PROCESSING

I began by examining Google Trends for the keywords: GameStop, r/WallStreetBets, and NYSE:GME. I chose to bucket the data into intervals separated by January 25, 2021 and February 4, 2021 (Fig. 1). In what follows, *Before* refers to data prior to January 25, *During* refers to data between January 25 and February 4, and

After refers to data after February 4. To gather a representative sample, I chose to collect highly upvoted (a measure of favorability) data. I used Reddit's API to gather the top posts from the past year and selected the most upvoted comments. Because the focus of comments are more likely to adhere to the topic of the post if they are higher in the comment tree, I only collected top-level comments. This resulted in 1.35 ± 0.05 million words in each range from over 1,300 posts.

Reddit uses Markdown to format text, so strong consideration was given to pre-processing the data. I removed all URLs and Markdown syntax, as well as corrected for missing punctuation and first letter capitalization. While emojis and special characters can add important context on the forum, they are often used to excess and were removed for simplicity. Resource and time limitations restricted the removal of some potentially impactful items (e.g. ASCII art).

III. METHODS

Words were used as symbols to capture syntactical information, which resulted in $46,500 \pm 3,000$ distinct symbols in each dataset. Additionally, punctuation (e.g. ., ?, (, ") were considered symbols.

To ensure an appropriate sample size was collected, resampling was performed on all datasets by separating them into 25 smaller samples. The average and median block entropies reflected the entropy for the unrestricted datasets.

Block Entropy of correlation length m was calculated as $S_m = -\sum_{x_1, \dots, x_m} p(x_1, \dots, x_m) \log_2 \frac{1}{p(x_1, \dots, x_m)}$. This was used in the calculation of correlation information $k_m = -S_m + 2S_{m-1} - S_{m-2}$. I set $S_0 = 0$ and $k_1 = \log_2 \nu - S_1$, where ν is the number of symbols.

To generate novel text, I combined the datasets and used a Markov Chain. The probability of choosing the next symbol x_{m+1} given the previous m symbols $x_1 \dots x_m$ is $p(x_{m+1} | x_1, \dots, x_m)$, with x_{m+1} being chosen randomly with probability α . I experimented with different lengths m , starting seeds, and randomness parameters α .

Due to the use of words as symbols α can have a magnified effect if a seldomly occurring word is selected, resulting in low comprehensibility for the rest of the generated text. A very low α is appropriate for this dataset, as

is $m < 4$ to prevent reproducing individual comments due to the large decrease in correlation information for $m \geq 4$.

-
- [1] S. Rao and T. Adinarayan, Losses top \$70 billion on short positions in u.s. firms - ortex data, Yahoo! Finance (2021).
 - [2] T. J. Gray, C. M. Danforth, and P. S. Dodds, Hahahahaha,

duuuuude, yeeessss!: A two-parameter characterization of stretchable words and the dynamics of mistypings and misspellings, PloS one **15**, e0232938 (2020).