

Trojan Horses in Amazon’s Castle: Understanding the Incentivized Online Reviews

Soheil Jamshidi, Reza Rejaie, Jun Li
Technical Report
Department of Computer and Information Science,
University of Oregon
{Soheilj, Reza, Lijun}@cs.uoregon.edu

Abstract—During the past few years, sellers have increasingly offered discounted or free products to selected reviewers of e-commerce platforms in exchange for their reviews. Such incentivized (and often very positive) reviews can improve the rating of a product which in turn sways other users’ opinions about the product. Despite their importance, the prevalence, characteristics, and the influence of incentivized reviews in a major e-commerce platform have not been systematically and quantitatively studied.

This paper examines the problem of detecting and characterizing incentivized reviews in two primary categories of Amazon products. We describe a new method to identify Explicitly Incentivized Reviews (EIRs) and then collect a few datasets to capture an extensive collection of EIRs along with their associated products and reviewers. We show that the key features of EIRs and normal reviews exhibit different characteristics. Furthermore, we illustrate how the prevalence of EIRs has evolved and been affected by Amazon’s ban. Our examination of the temporal pattern of submitted reviews for sample products reveals promotional campaigns by the corresponding sellers and their effectiveness in attracting other users. Finally, we demonstrate that a classifier that is trained by EIRs (without explicit keywords) and normal reviews can accurately detect other EIRs as well as implicitly incentivized reviews. Overall, this analysis sheds an insightful light on the impact of EIRs on Amazon products and users.

I. INTRODUCTION

As the popularity of online shopping has rapidly grown during the past decade, the shoppers have increasingly relied on the online reviews and rating provided by other users to make more informed purchases. In response to shoppers’ behavior, product sellers have deployed various strategies to attract more positive reviews for their products as this could directly affect popularity of these products among users and thus their ability to sell more products online. Several prior studies have examined different aspects of online reviews including fake or spam [9], [13], [10], [17], [12], [2] and also biased and paid reviews [21], [22], [23], [18], [6] in different online shopping platforms.

The importance of online reviews has also prompted major e-commerce sites (e.g., Amazon) to implement certain policies to ensure that the provided user reviews and ratings are legitimate and unbiased to maintain the trust of online shoppers. In response to these policies, seller’s strategies for boosting their product rating have further evolved. In particular, in the past few years, some sellers have increasingly offered discounted or free products to selected online shoppers in

exchange for their (presumably positive) reviews. We refer to these reviews as *incentivized reviews*. Major e-commerce sites such as Amazon require reviewers to disclose any financial or close personal connection to the brand or the seller of the reviewed products [3]. However, it is unlikely that average shoppers who solely rely on product ratings notice the biased nature of such reviews. Intuitively, the reviewers who provide incentivized reviews may behave differently than other reviewers for the following reasons: (i) they might feel obligated to post positive reviews as the products are provided for free or with a considerable discount, (ii) their expectations might be lower than other users as they do not pay the full price, and (iii) they do not often consider the long-term usage of the product (e.g., product return or customer service) in their reviews. The presence of such incentivized reviews in Amazon has been reported in 2016 [20], however, to our knowledge, the prevalence of incentivized reviews, their characteristics, and their impact on the ecosystem of a major e-commerce site have not been systematically and quantitatively studied. Although Amazon has officially banned submission of incentivized reviews in October of 2016 [1], it is important to study such reviews to be able to determine whether Amazon’s new policy solved the issue or just forced reviewers to go under cover.

To tackle this important problem, this paper focuses on capturing and characterizing several aspects of incentivized reviews in the Amazon.com environment. We leverage the hierarchical organization of Amazon products into categories/subcategories and collect all the information for top-20 best-seller products in all subcategories of two major categories. The first contribution of this paper is a method to identify explicitly incentivized reviews (EIRs) on Amazon. We identify a number of textual patterns that indicate explicitly incentivized reviews. We carefully capture and fine-tune these textual patterns using a regular expression. We then use these patterns to identify a large number of EIRs along with their associated products and reviewers.

The second contribution of this paper is the characterization of key features of EIRs and associated reviewers and products. Our analysis demonstrates the effect of Amazon ban on the prevalence of EIRs as well as the difference between the features of EIRs and normal reviews. We also examine the temporal pattern of EIR, and non-EIR reviews that a product

receives and a reviewer produces to address two questions: (i) how the arrival pattern of EIRs for a specific product affects the level of interest (*i.e.*, rate of non-EIRs and their assigned rating) among other users, and (ii) how individual reviewers over time become engaged in providing EIRs. Finally, given an apparent gap between features of normal reviews and EIRs, we examine whether machine learning techniques can detect these differences to identify both explicitly or implicitly incentivized reviews. We show that such a technique can indeed detect other incentivized reviews.

The rest of this paper is organized as follows: We describe our data collection technique and our datasets in Section II. Section III presents our method for detecting EIRs. We characterize several aspects of EIRs and their associated products and reviews in Section IV. Section V discusses the temporal patterns of EIRs and non-EIRs that are submitted for individual products or produced by individual reviewers. Section VI presents our effort for automated detection of other explicitly or implicitly incentivized reviews using machine learning techniques. We present a summary of most relevant prior work and how they differ from this study in Section VII. Finally, Section VIII concludes the paper and summarizes our future plans.

II. DATA COLLECTION AND DATASETS

This section summarizes some of the key challenges with data collection and then describes our methodology for collecting representative datasets that we capture and use for our analysis. Amazon web site organizes different products into categories that are further divided into smaller sub-categories. Each product is associated with a specific seller. A user who writes one (or multiple) review(s) for any product is considered a reviewer of that product. For each entity (*i.e.*, user, review or product), we crawled all the available attributes on Amazon as follows:

- Reviews’ attributes: review id, reviewer id, product id, Amazon Verified Purchase (AVP) tag, date, rating, helpful votes, title, text, and link to images.
- Products attributes: product id, seller id, price, category, rating, and title.
- Reviewers’ attributes: reviewer id, rank, total helpful votes, and publicly available profile information.

In particular, *AVP tag* of a review indicates whether the corresponding reviewer has purchased this product through Amazon and without deep discount or not [4].

There are a few challenges for proper collection and parsing of this information from Amazon. First, there is a very large number of product categories where the format, available fields for products, and tendency of users to offer reviews widely vary across different categories. Furthermore, we need to comply with the ethical guidelines as well as the enforced rate limits by Amazon servers for crawlers which makes it impossible to collect the reviews for all products within a reasonable window of time. To cope with these challenges, we collect three datasets where each one provides representative samples of products, reviews and reviewers.

TABLE I
BASIC FEATURES OF OUR DATASETS

| | Products (DS1) | EIRs (DS2) | Normal Reviews | Reviewers (DS3) |
|-----------|-------------------|---------------|-------------------|--------------------|
| Reviews | 3,797,575 | 100,086 | 100,086 | 217,000 |
| Reviewers | 2,654,048 | 39,886 | 98,809 | 2,627 |
| Products | 8,383 | 1,850 | 1,641 | 184,124 |

Sample Products (DS1): We focus on two popular categories of products, namely *Electronics* and *Health & Personal Care* since they have a large number of sub-categories and products that receive many reviews. To make the data collection manageable and given the skewed distribution of reviews across products, we only capture all the information for the top-20¹ best seller products in each sub-category in the above two categories from *Amazon.com*. While these products represent a small fraction of all products in these two categories, the top-20 products receive most of the attention (#reviews) from users and enable us to study incentivized reviews. We refer to this product-centric dataset as *DS1*.

Sample EIRs (DS2): Using our technique for detecting Explicitly Incentivized Reviews (EIR) that is described in Section III, we examine all the reviews associated with products in DS1 and identify any EIRs among them. We refer to this set of EIRs as DS2 dataset.

Normal Reviews: After excluding EIRs, we examine the remaining reviews for products in DS1 and consider each review as normal if it is not among EIRs and (i) associated with an Amazon Verified Purchase, (ii) submitted on the same set of products that received EIRs, and (iii) submitted by users who have not submitted any EIRs. We rely on this rather conservative definition of normal reviews to ensure that they are clearly not incentivized. We identified 1,214,893 normal reviews and then selected a random subset of them (the same number as EIRs). We refer to these selected reviews as our normal review dataset that serves as the baseline for comparison with EIRs in some of our analysis.

Incentivized Reviewers (DS3): To get a complete view of sample incentivized reviewers, we randomly select 10% of reviewers associated with the reviews in *DS2* dataset. For each selected reviewer with a public profile, we collect their profile information and all of their available reviews. Overall, we collect this information for 2,627 reviewers and only consider their reviews for our analysis.

The DS1, DS2, and Normal reviews datasets were collected in December 2016, and the Reviewers dataset (DS3) was collected in January 2018.

III. DETECTING EXPLICIT INCENTIVIZED REVIEWS

Automated identification (or labeling) of incentivized reviews requires a reliable indicator in such reviews. To this end, we first focus on reviews in which the reviewer *explicitly* indicates her intention for writing the review in exchange

¹<https://www.amazon.com/gp/bestsellers/>

for a free or discounted product. Such an indication must be provided in the reviews since Amazon requires that reviewers disclose any incentive they might have received from the sellers [3]. Furthermore, these reviewers also include such incentives in their reviews² to attract more sellers to offer them similar incentives in exchange for their reviews to promote their products. For these recruit these reviewers for promotional campaigns and following the stylistic features of their precedents [15] could be among the reasons they explicitly disclose their motivation. Our manual inspection of a large number of reviews revealed that many reviewers indeed explicitly state their incentive for writing their reviews. These reviews contain some variants of the following statements: “I received this product at a discount in exchange for my honest/unbiased review/feedback.” To capture all variants of such statements, we select any review that matches the following regular expression in a single sentence of the review:

```
'(sent|receive|provide)[^\.!?]*
(discount|free|in-trade|in-exchange)[^\.!?]*
(unbiased|honest)[^\.!?]*
(review|opinion|feedback|experience)'
```

This will find the reviews that have the above combination of words in a single sentence. Among all the 3.79M reviews in the DS1 dataset, 100,086 reviews submitted by 39,886 users on 1,850 products match some variants of the above regular expression in one sentence. We consider these 100,086 reviews as EIRs and group them in our DS2 dataset. We also considered a more relaxed setting where reviews could have the above regular expression across multiple sentences. This strategy tags 325,043 reviews from 210,198 users on 7,059 products as EIR. However, our careful inspection of many of these newly-identified EIRs by this more flexible strategy revealed that some of them are non-incentivized reviews that happen to match the regular expression. Here is an example

Received my PS4 on Friday (...) just like Xbox Live they discount games and have free offerings (...) to be honest I am not a fan of (...) UPDATE: although my review was (...)

To avoid any such a false-positive in our EIRs, we adopt a conservative strategy and only consider a review as EIR if the desired pattern detected within a single sentence.

EIR-Aware Reviews: Our extensive manual inspection of the identified EIRs also revealed that in a tiny fraction (only 30 reviews) the reviewer simply refers to other EIRs to complain and indicate its awareness and inform other users of such incentivized reviews. However, these reviews are not incentivized themselves. To exclude these reviews, we manually checked random samples of reviews and found that aware reviews contain one of the following terms (*who received—with the line “i received—which say they received—their so-called “honest”*). Here is an example of such a review:

Lately I started questioning the honesty of Amazon reviews because so many reviews say “I received this product at

²These reviewers also adopt a certain style and level of details (including pictures) in their review as well.

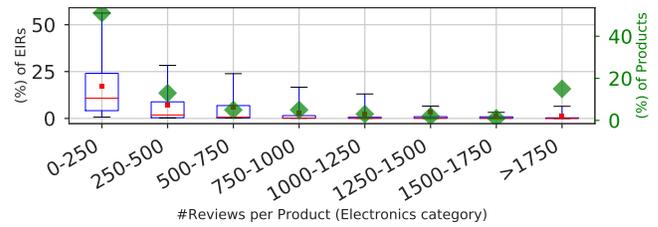


Fig. 1. Distribution of fraction of EIRs per product in *Electronics* category

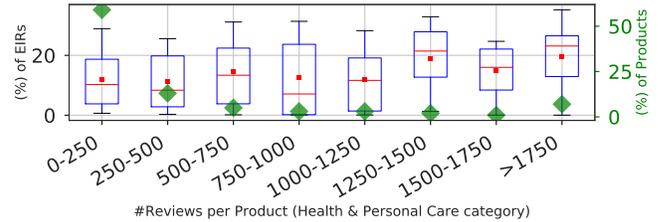


Fig. 2. Distribution of fraction of EIRs per product in *Health* category

the discount price in exchange for an honest opinion”
How can we trust this kind of reviews? I don t know whom to believe anymore.
I feel like I got scammed. If you think “I received this product at a discount (free) for an honest review” is really going, to be honest, think again.
STOP SENDING FREE PRODUCTS FOR 5-STAR REVIEWS. Honestly, almost every 5-star review has a disclaimer at the end saying that they received the product free in exchange for an “honest” review. Well how come there is no person who received it for free that rated it below 5????

We then exclude any identified EIR that matched these aware patterns. After extensive manual work on this step, we found only 30 aware reviews by 26 reviewers on 29 products that are excluded from DS2. Interestingly, all these aware reviews were collectively marked as helpful by 194 other users, indicating that many other reviewers felt the same way about the incentivized reviews. This illustrates how the presence of incentivized reviews could impact the trust of customers in the authenticity of Amazon reviews.

IV. BASIC CHARACTERIZATIONS

In this section, we examine a few basic characterizations of EIRs and their associated products and reviewers in order to shed some light on how these elements interact in Amazon.com.

Product Characteristics

One question is *what fraction of reviews for individual products are EIRs?* We use our all products in dataset (DS2) to examine several characteristics of products that receive at least one EIR. Figure 1 and Figure 2 present the summary distribution of the fraction of EIRs for different group of products based on their total number of reviews in each category. The second Y-axis on these figures shows the fraction of all products (per category) that are in each group. The

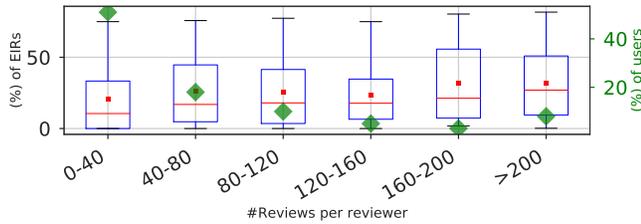


Fig. 7. Distribution of the fraction of provided EIRs per Reviewer

We perform sentiment analysis on both content and title of reviews using *textblob* library. The sentiment is measured by a value within the range of $[-1, 1]$ where 1 indicates positive, 0 neutral, and -1 the negative sentiment. Figure 9 presents the distribution of sentiment for the content of EIRs and normal reviews. We observe that 9.5% (9,498) of normal reviews have negative sentiment, 9.1% are neutral (*i.e.*, their sentiment measure is zero) and the rest of positive reviews are spread across the whole range with some concentration around 0.5, 0.8, and 1. In contrast, the sentiment of nearly all EIRs are positive, but more than 80% of them are between 0 to 0.5. In essence, the sentiment of normal reviews is widespread across the entire range while sentiments for EIR are mostly positive but more measured.

Similarly, less than half of the normal reviews and three-quarter of EIRs have titles with positive sentiments, *i.e.*, title of EIRs have more positive sentiments. Using *TextBlob* library, we also analyzed the review’s *Subjectivity*, which marks the presence of opinions and evaluations rather than using objective words to provide factual information. Figure 10 depicts the CDF of the subjectivity across EIRs and normal review datasets. This figure reveals that the subjectivity for 83% of EIRs are between 0.4 and 0.8 while the subjectivity is widely spread across the whole range for normal review, *i.e.*, $0 \leq 0.75$ subjectivity for 7.8% (84%) of normal reviews.

We use the *Gunning Fog index* [7] to measure the readability test for English writing in each group of reviews. This index estimates the number of years of formal education a person needs to understand the text on the first reading. For example, a Fog index of 12 requires the reading level of a U.S. high school senior. Figure 11 shows the CDF of the Fog index across EIRs and normal reviews. This result illustrates that the readability of EIRs requires at least 4 years of education and is 1.5 points higher than normal reviews in average (7.5 vs. 6 years of education). Also, the index exhibits much smaller variations across EIRs. In short, the writing of EIRs is more elaborate.

We can also assess the similarity of review content across all written reviews by individual reviewers that indicates to what extent she might repeat various phrases in different reviews. To this end, we calculate the *Jaccard index* on the 2-shingles (or bi-grams) between the content of all pairs of reviews written by each reviewer and use their average value as her average degree of self-similarity of reviews. The two group of users for this analysis are those who have provided at least two EIRs

vs those who provided at least two normal reviews. Figure 12 depicts the CDF of average self-similarity across all reviews (after removing the acknowledgment sentence in all EIRs) by individual reviewers for these two groups of reviewers. This figure demonstrates a measurably higher level of self-similarity across reviews provided by reviewers who provide EIRs. In particular, self-similarity is zero and less than 5 for 75% and 95% of normal users, respectively. However, 95% (and 25%) of users who provide EIRs exhibit a degree of self-similarity more than 1 (more than 10). This suggests that: *Reviewers that provide some EIRs tend to write more similar reviews than normal users.*

Length of Reviews: The overall length of a review and its title could be viewed as measures of its level of details. Figure 13 shows the summary distribution of review length (in terms of characters) across EIRs and normal reviews. We observe that the typical (*i.e.*, mean) length of an EIR (599 characters) is more than three times longer than a normal review (179 characters). Interestingly, the longest normal review (14.8K character) is much longer than the longest EIR (11K character). We observe a similar pattern for the length of reviews considering word count. Furthermore, the title for EIRs are typically 6.6 words long which is two words longer than the title of normal reviews.

Helpfulness: Another important aspect of reviews is how helpful they are to other users. We measure the level of helpfulness of a review by its total number of *helpful votes* (or up-votes). A slightly larger fraction of normal reviews (12.68%) receive up-votes compare to the EIRs (10.87%). Figure 8 shows the Complementary Cumulative Distribution Function (CCDF) of the number of up-votes for EIRs and normal reviews. The main difference is in the tail of these distributions (for reviews with many up-votes) that shows most popular normal reviews receive a significantly larger number of up-votes than EIR ones. In short, EIRs and normal reviews commonly exhibit the same degree of helpfulness except the extreme cases for normal reviews are much more helpful.

The maximum number of likes for normal reviews is 8,404 while for the biased reviews is 289, the minimum for both is 0. The average for normal reviews is around 0.62 while for biased reviews is 0.37 and we have more liked reviews among normal reviews (12.68%) compare to the biased reviews (10.87%).

Star Rating: A critical aspect of a review is the star rating (in the range of 1 to 5 stars) that it assigns to a product. Figure 14 presents the summary distribution of assigned star rating by EIRs and normal reviews. We observe that the assigned rating by EIRs is frequently more positive than normal reviews. More specifically, 95% (75%) of EIRs associated the rating of at least 3 (5) stars while this number drops to 1 (4) for normal reviews.

Reviewer-Review Mapping Per Product: A majority (99.8%) of reviewers in our EIR dataset (DS2) have written only one EIR for each product. We only found 73 users who have written multiple EIRs for at least one product. These reviews add up to the total of 151 EIRs for 32 unique products. None of the users in our user-centric dataset (DS1) writes

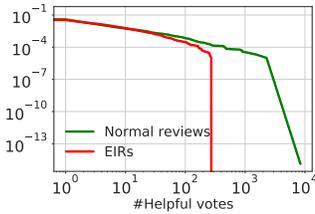


Fig. 8. CCDF of Helpfulness

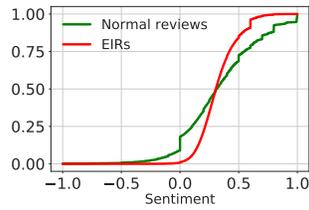


Fig. 9. CDF of Review sentiment

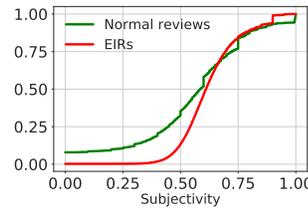


Fig. 10. CDF of Review subjectivity

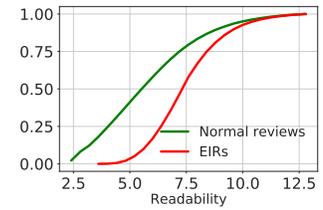


Fig. 11. CDF of Review Readability

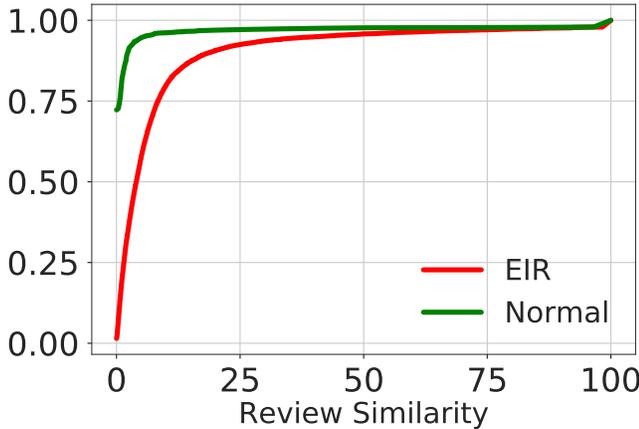


Fig. 12. Self-Similarity of reviews per user

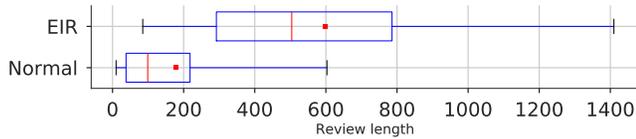


Fig. 13. Distribution of review length

multiple EIRs for a single product. Given the one-to-one relationship between the absolute majority of reviewer-review pairs per product, for the rest of our analysis, we assume each reviewer has only a single review per product and vice versa.

V. TEMPORAL ANALYSIS

All of our previous analysis have focused on the overall characteristics of reviews, reviewers, and products over their entire lifetime. Intuitively, product sellers offer various incentives to attract reviewers and obtain incentivized reviews for their specific product. Obtaining these incentivized reviews

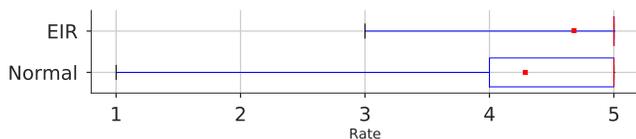


Fig. 14. Distribution of reviews' assigned star rating

over time increases the available information and improves the overall image (*e.g.*, rating) of the product. This, in turn, expands the level of interest among (ordinary) users who may consider to buy the product and provide their own review. Examining the temporal pattern of submitted reviews (by various reviewers) for a product or submitted reviews by a reviewer (for any product) sheds an insightful light in various dynamics among seller products, reviews, and reviewers.

In this section, we tackle two important issues: First, we inspect the “*review profile of sample products*” to study how the temporal pattern of obtained EIRs for a product affect the level of interest among other users. Second, we examine the “*review profile of sample reviewers*” to explore how reviewers get engaged in producing EIRs. To tackle these questions we have inspected temporal patterns for many products and reviewers, and only present a few sample cases to illustrate our key findings better.

In this analysis, we primarily focus on the number of EIRs, non-EIRs (*i.e.*, reviews that are not tagged as EIR by our method) associated with a product (or a reviewer) per day and their (cumulative) average rating.³ across EIRs or non-EIRs that a product receives or a reviewer assigns.

Product Reviews

We consider four different products to examine the temporal correlations between the daily number of EIRs and the level of interest among other users, namely the number of non-EIRs and their ratings, for each product.

Note that a product seller can (loosely) control the arrival rate of EIRs by offering incentives (or promotions) with a particular deadline to a specific set of reviewers. We refer to such an event as a *promotional campaign*. The goal of our analysis is to investigate whether and to what extent such a campaign affects the number of non-EIRs and their rating for individual products. Note that a product seller can (loosely) control the arrival rate of EIRs by offering incentives (or promotions) with a certain deadline to a specific set of reviewers. We refer to such an event as a *promotional campaign*. By specifying a deadline for the incentive or promotion, the seller can also force interested users to write their reviews within a specific

³Amazon appears to rely on some weighted averaging method [5] to calculate the overall rating of a product based on factors such as the recency of a review, its helpfulness and whether it is associated with a verified purchase. Since the details of Amazon’s rating method is unknown, we simply rely on a linear moving average of all ratings to determine the overall rating of each product or reviewer over time.

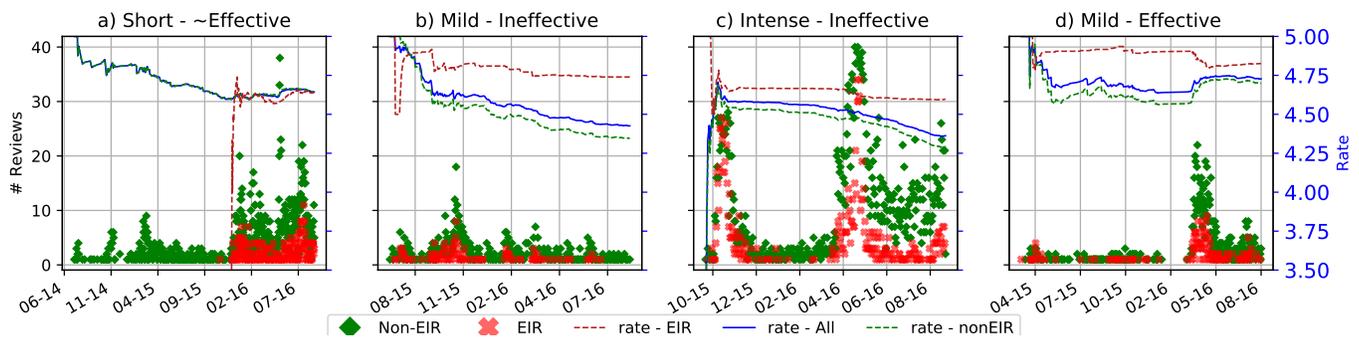


Fig. 15. Temporal Patterns of Reviews for Individual Products

window of time. We simply assume that any measurable, sudden increase in the number of daily EIRs for a product is triggered, by a promotional campaign that is initiated by its seller. The goal of our analysis is to investigate whether and to what extent such a campaign affects the number of non-EIRs and their rating for individual products. Each plot in Figure 15 presents the daily number of EIRs (with a red X), the daily number of non-EIRs (with a green diamond), the cumulative average rating for all non-EIR (with dotted green) and EIR (with dotted red lines) for a single product. Each plot also shows the cumulative rating of all reviews with a solid blue line. Three rating lines on each plot are based on the right Y-axis showing the star rating (1 to 5 scale).

Short & Moderately Effective Campaigns: Figure 15-a shows a product that has been consistently receiving a few daily non-EIR (and not a single EIR) reviews over a roughly two year period. Its average product rating rather consistently drops during 2015. A persistent daily rate of EIR suddenly starts in early 2016 and continues for a few months indicating a likely promotional campaign. The campaign triggers a significant increase in the number of non-EIRs. Interestingly, the average rating of EIR rapidly converges to the average rating of non-EIR (and the overall rating) and not only prevents further dropping but also rather improves the overall rating of this product. This appears to be a short-term (over a few months) and moderately effective promotional campaign by the seller.

Multiple Mild but Ineffective Campaigns: Figure 15-b presents another product that consistently receives non-EIRs over a one year period. We can also observe ON and OFF periods of EIRs that did not seem to seriously engage other users with this product (*i.e.*, no major increase in the daily rate of non-EIRs). The assigned rating by EIRs is relatively constant, and their gap with the rating of non-EIRs (and overall rating) rapidly grows. Clearly, these multiple mild campaigns are not effective in raising the ratings of the product.

Multiple Intense but Ineffective Campaigns: Figure 15-c shows a product that has been consistently receiving both EIR and non-EIRs over a year-long period. However, there are two (and possibly three) distinct windows of time (each one is a few weeks long) with pronounced peaks in the number of daily

EIRs which suggests two intense campaigns. Interestingly, the first campaign only generates short-term interest among ordinary users (shown as a short-term increase in the daily number of non-EIRs) while the second campaign triggers more non-EIRs. The average rating of EIR is clearly above non-EIRs. However, the average rating of non-EIRs (and even EIRs) continues to drop over time despite the increased level of attention by other users after the second campaign. Therefore, these multiple intense campaigns were not able to improve the overall rating of this product.

Multiple Mild and Effective Campaigns: Figure 15-d shows a product with a low and persistent daily EIR and non-EIR over a one-year period. We then observe a couple of months with absolutely no reviews that suggest the unavailability of the product. This is followed by a more active campaign of EIRs over a month that continues at a lower rate. This last campaign seems to significantly increase the level of interest among the regular users as well as their rating for this product. In particular, the average rating by non-EIRs was relatively stable and clearly below the rating by EIRs until the last campaign. Interestingly, the last campaign decreases the overall rating by EIRs while enhances the overall rating by non-EIRs. Therefore, we consider this as an effective campaign.

These examples collectively demonstrate that while a seller could loosely control the duration and intensity of its promotional campaign for a product, its impact on the level of engagement by other users could be affected by many other factors (*e.g.*, quality of reviews, strategies of competitors, quality of the product, and product rank on different search queries on Amazon) and thus widely vary across different products.

User Reviews

We now focus on the written EIRs and non-EIRs by individual users over time. Similar to the temporal patterns of product reviews, we show the number of daily EIRs (with a red X), non-EIRs (with a green circle). We also show average assigned rating by EIRs (with red dotted line) and non-EIRs (with green dotted line) of the reviewer over time. The three plots in Figure 16 present the temporal pattern of all reviews (for any product) and their rating for three different reviewers.

Persistent EIR Writer: Figure 16-a shows a user who provided a single review in 2013 and was inactive for more than a year. Starting April 2014, she has been providing a couple of EIRs and/or non-EIRs a day for 2.5 years, and then her activity significantly dropped. Her average rating for EIRs and non-EIRs are very similar. It seems that this reviewer becomes active in Amazon to provide EIRs.

Active EIR Writer: Figure 16-b shows a user who has been actively writing non-EIRs over 16 years since 2001, and her level of activity has gradually increased. Interestingly, she started posting EIRs from 2015 for two years and then stopped. These two years are perfectly aligned with the period in which EIRs have become rapidly popular in Amazon (as we showed in figure 3). Furthermore, the overall assigned rating by this reviewer in non-EIRs was relatively stable over time which was slightly lower than her assigned rating in EIR reviews. This reviewer is a perfect example of a serious Amazon reviewer who takes advantage of offered incentives by sellers for writing EIRs.

Casual EIR Writer: Figure 16-c shows the temporal pattern of review submission by a user who has been in the system since 2013. However, he became moderately active in the middle of 2015 and provided some EIRs and mostly non-EIRs in the past two years. The number of his EIRs are limited and mostly written over a one year period. It is rather surprising that his rating in EIRs gradually grew over time and was always slightly lower than his ratings for non-EIRs. Far from normal behavior, he has written 49 non-EIRs in one day in 2016 (the green dot above the rating lines). Overall, he appears to be a moderate reviewer who casually writes EIRs.

In summary, our user-level temporal analysis of EIRs and non-EIRs indicates that: *Reviewers exhibit different temporal patterns in producing EIRs. However, users are more active while incentives are offered.*

VI. DETECTING OTHER INCENTIVIZED REVIEWS

So far in this paper, we primarily focused on EIRs for our analysis since we can reliably detect and label them as incentivized reviews. However, in practice, there might exist a whole spectrum of explicitly or implicitly incentivized reviews besides EIRs. An intriguing question is *whether all these incentivized reviews (regardless of their implicit and explicit nature) share some common features that can be leveraged to detect them in an automated fashion?* To tackle this question, we consider a number of machine learning and neural network classification methods that are trained using a combination of basic and text features of the reviews.

Pre-processing Reviews: We use 100K random EIRs (from the DS3 dataset) and the same number of normal reviews as our labeled data. First, we remove the sentence that indicates the explicit incentive of a reviewer from each EIR before using the EIRs in this analysis so that these sentences do not serve as a dominant and prominent explicit feature. Second, we consider the following pre-processing of text of reviews to examine their exclusive or combined effect on the accuracy of various detection methods: (i) Converting all characters to

lower-case, 2) Using the stem of each word in the review (e.g., “wait is the stem for words “waiting”, “waits”, “waited”). (iii) Using only alphabet characters in a review, and (iv) Removing all the stop-words.

Classification Methods: We examine a number of classification methods including *Multi-Layer Perceptron (MLP)*, *SVM*, *GaussianProcess*, *DecisionTree*, *RandomForest*, *AdaBoost* Classifiers, *Bi-grams* and *Tri-grams* (with and without $tf-idf^4$), and *character-based bi- and tri-grams*.

Each classifier is trained and tested in three scenarios with a different combination of review features follows: (i) *Basic Features:* Using nine basic features of reviews, length, sentiment, subjectivity, and readability of review text, star-rating and helpfulness of reviews, and length, sentiment, and subjectivity of title. (ii) *Text Features:* Using extracted text features of the character-based Tri-grams (limited to $2^{*}10$ text features) of the reviews, (iii) *All Features Basic + Text* (character Using a combination of basic and text features.

TABLE II
THE BASIC REVIEW FEATURES, USED FOR CLASSIFICATION

| | | |
|-------------------|----------------------|--------------------|
| Text Length | Star-rating (1 to 5) | Text Sentiment |
| Title Sentiment | Helpfulness | Title Subjectivity |
| Text Subjectivity | Readability | Title Length |

Individual methods are evaluated in 5 and 10-fold cross-validation as well as 70/30 test and training split manner. We only present the result for the MLP method using pre-processed reviews after removing all stop words and replacing all remaining words with their stem part as this combination exhibits the highest level of accuracy. The results for all other cases are available in our related online technical report [8].

We found MLPC considerably better regarding memory usage, computation time, and accuracy on a 50-50% combination of EIR and normal reviews in the training set. We use 90% of data for training and testing and 10% of data for hyperparameter tuning using the *grid-search* in SciKitLearn library. We tune the values of *alpha* and also decide on the number of hidden layers and nodes in each layer. We trained the MLP classifier using default parameters⁵, except for *alpha* (the L2 penalty regularization term) and *hidden_layer_size* that we set to 0.1 and (50,30,10), respectively.

Table III presents the average accuracy, recall, precision, F1-score, Precision-Recall Area Under Curve (P-R AUC), and the Receiver Operating Characteristic (ROC) AUC for MLP Classifier (MLPC) overall runs and their standard deviation on each feature set. These results indicate that even without the explicit acknowledgment sentence in EIRs, a classifier can accurately detect EIRs from normal reviews using basic or text feature. The accuracy further improves if we combine both sets of features.

We examine the ability of a classifier for detecting EIRs in other categories. To this end, we divide EIRs and normal

⁴Term Frequency Inverse Document Frequency

⁵The full list of parameters are available on: <https://goo.gl/TmWueZ>

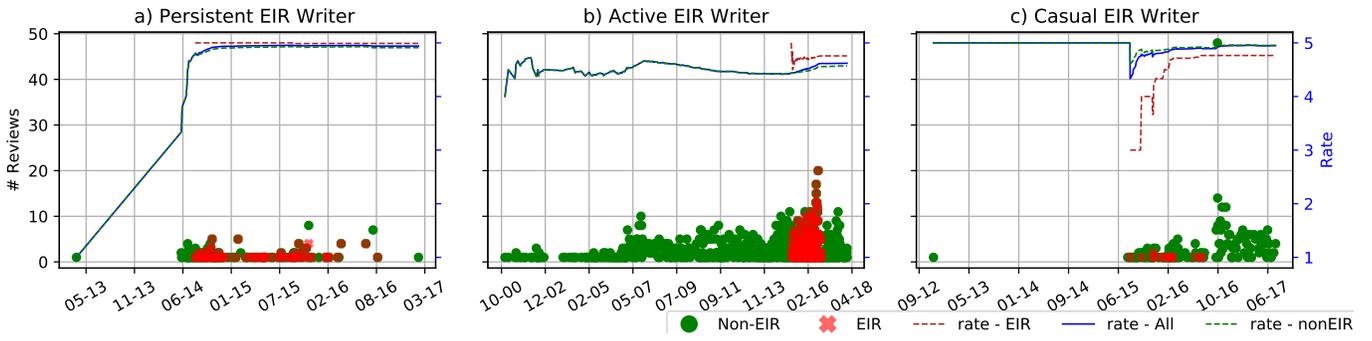


Fig. 16. Temporal Patterns of Reviews for Individual Reviewers

reviews into two groups based on the category of their corresponding product (*i.e.*, Electronics and Health). We train two classifiers, called *C-Health* and *C-Elect.*, where each one only uses EIRs and normal reviews (with a combination of basic and text features) associated with products in one category. Finally, we test each classifier on reviews from the other category to assess their accuracy in detecting EIR and normal reviews. The last two rows of Table III present the accuracy of MLPC for this cross-category detection. These results show that the accuracy of cross-category detection of EIRs (for these two categories) is still sufficiently high ($\leq 80\%$). Interestingly, the classifier that is trained with Health reviews exhibits a higher accuracy in detecting Electronics reviews.

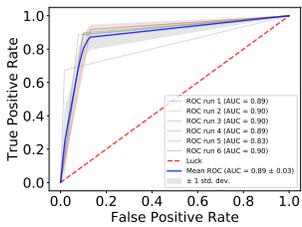


Fig. 17. ROC Area Under Curve for MLP classifier on basic and text features

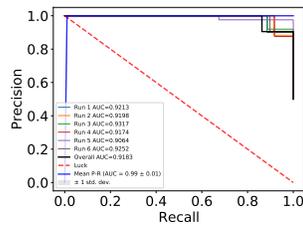


Fig. 18. Precision-Recall Area Under Curve for MLP classifier on basic and text features

Next, we investigate the ability of our trained classifier using the basic and text-based features in detecting other incentivized reviews, namely implicitly incentivized review (IIR) and other explicitly incentivized reviews that do not contain the identified regular expressions and thus our method has not detected. We randomly select 100,000 reviews (during 2016) from the DS1 dataset that are neither EIR nor normal reviews and then use the trained classifier to determine whether any of these *unseen* reviews are classified as incentivized or normal. After removing reviews with less than three words in the text, we kept 98,594 reviews. The classifier flags 20,892 (21.19%) of these reviews as incentivized. Our manual inspection of the content of these reviews revealed that they can be broadly divided into two groups as follows:

Other Explicitly Incentivized Reviews: 3,799 (18%) of reviews labeled as incentivized contain other explicit patterns indicating the incentivized nature of their reviews that we had

not considered, *e.g.*, “I had the opportunity to get it for my review”, “received with a promotion rate”.

Implicitly Incentivized Reviews (IIRs): We note that the absence of any explicit disclosure of incentives in the remaining reviews does not imply that they are not incentivized. We hypothesize that some of these are implicitly incentivized reviews (IIR). To verify this hypothesis across all the remaining flagged reviews, we rely on the pairwise relationship between review-product and review-reviewer and check any of these reviews against the following two conditions: (*i*) whether a review is associated with a product that had received at least one other EIR, or (*ii*) whether a review is provided by a user who has submitted at least one other EIR. We observe that 296 reviews are affiliated with both EIR reviewers and EIR products (*i.e.*, meet both conditions) while 8544 (41%) of them are only affiliated with EIR products and 63 reviews are only affiliated with EIR reviewers. Intuitively, meeting both conditions offers a stronger evidence that a review could be (implicitly) incentivized. Our manual inspection of reviews in these three groups confirmed this intuition. While reviews that met both conditions contain some indication of incentive (*e.g.*, *for my honest result, promotional price*), reviews related only to products only contained moderate hints (*e.g.*, *I have to thank this seller*).

VII. RELATED WORK

Detection and analyzing of spam reviews started in 2008 by labeling the (near) duplicate reviews as spam and using supervised learning techniques to detect spam reviews [9]. Since then, different aspects of online reviews have been investigated such as behavioral abnormalities of reviewers [13] and review quality and helpfulness [16], [11], [14]. Studies on spam detection have deployed a diverse set of techniques. Early studies relied on unexpected class association rules [10] and standard word and part of speech n-gram features with supervised learning [17] that are later improved by using more diverse feature sets [12]. *FraudEagle* [2] was proposed as a scalable and unsupervised framework that formulates opinion fraud as a network classification problem on a signed network of software product reviews of an app store. These studies also relied on different strategies, such as Amazon Mechanical

Turk [17] or manual labeling [12] to create a labeled dataset for their analysis.

The effect of incentives on reviewers and quality of reviews are studied by Qiao et al. [19]. They showed that external incentives might implicitly shift an individual's decision-making context from a pro-social environment to an incentive-based environment. Wang et al. [22] modeled the impact of bonus rewards, sponsorship disclosure, and choice freedom on the quality of paid reviews. In a qualitative study, Petrescu et al. [18] examined the motivations behind incentivized reviews as well as the relationship between incentivized reviews and the satisfaction ratings assigned by consumers to a product. They showed that the level of user engagement depends on a cost-benefit analysis. Burtch et al. [6] focused on social norms instead of financial incentives. By informing individuals about the volume of reviews authored by peers, they test the impact of financial, social norms, and a combination of both incentives in motivating reviewers. The study by Xie [23] unveiled the underground market for app promotion and statistically analyzed the promotion incentives, characteristics of promoted apps and suspicious reviewers in multiple app review services.

To the best of our knowledge, none of the prior studies have systematically examined the prevalence of EIRs, their basic characteristics, and their influence on the level of interest among other users to a product based on large-scale quantitative measurements in a major e-commerce platform.

VIII. CONCLUSION

In this paper, we presented a detailed characterization of Explicitly Incentivized Reviews (EIRs) in two popular categories of Amazon products. We presented a technique to detect EIRs, collected a few datasets from Amazon and identified a large number of EIRs in Amazon along with their associated product and reviewer information. Using this information, we compared and contrasted various features of EIRs with reasonably normal reviews. We showed that EIRs exhibit different features compared to normal reviews and discussed the implications of these differences. We then zoomed into the temporal pattern of submitted EIR reviews for a few specific products and submitted reviews by a few specific reviewers. These temporal dynamics demonstrated whether/how promotional campaigns by a seller could affect the level of interest by other users and how reviewers could get engaged in providing EIRs. Finally, we illustrated that machine learning techniques can identify EIRs from normal reviews with a high level of accuracy. Moreover, such techniques can accurately identify other explicitly and implicitly incentivized reviews. We leverage affiliation of reviews with reviewers and products to infer their incentivized nature.

Some of our future plans are as follows: We plan to iteratively improve the performance of classifiers by incorporating other explicit patterns. Furthermore, we deploy probabilistic techniques to infer the likelihood that a review is incentivized based on its affiliation with other products and reviewers.

Finally, we explore whether the incentivized reviews have disappeared entirely from Amazon or become more implicit.

REFERENCES

- [1] aboutAmazon.com. Update customer review, <https://goo.gl/fiVa8j>, 2016.
- [2] L. Akoglu, R. Chandy, and C. Faloutsos. Opinion fraud detection in online reviews by network effects. *ICWSM*, 2013.
- [3] Amazon.com. Community guidelines, <https://www.amazon.com/gp/help/customer/display.html?nodeid=14279631>, 2018.
- [4] Amazon.com. About amazon verified purchase reviews, <https://goo.gl/aNcPCR>, 2018.
- [5] T. Bishop. <https://www.geekwire.com/2015/amazon-changes-its-influential-formula-for-calculating-product-ratings/>, 2015.
- [6] G. Burtch, Y. Hong, R. Bapna, and V. Griskevicius. Stimulating online reviews by combining financial incentives and social norms. *Management Science*, 2017.
- [7] R. Gunning. The technique of clear writing. *McGraw-Hill, NY*, 1952.
- [8] S. Jamshidi, <https://goo.gl/hR8kFV>, 2018.
- [9] N. Jindal and B. Liu. Opinion spam and analysis. In *ACM International Conference on Web Search and Data Mining*, 2008.
- [10] N. Jindal, B. Liu, and E.-P. Lim. Finding unusual review patterns using unexpected rules. In *Proceedings of the 19th ACM international conference on Information and knowledge management*, 2010.
- [11] S.-M. Kim, P. Pantel, T. Chklovski, and M. Pennacchiotti. Automatically assessing review helpfulness. In *Proceedings of the ACL Conference on empirical methods in natural language processing*, 2006.
- [12] F. Li, M. Huang, Y. Yang, and X. Zhu. Learning to identify review spam. In *Proceedings of IJCAI*, 2011.
- [13] P. Lim, V. Nguyen, N. Jindal, B. Liu, and H. Lauw. Detecting product review spammers using rating behaviors. In *Proceedings of ACM international conference on Info. and knowledge management*, 2010.
- [14] J. Liu, Y. Cao, C.-Y. Lin, Y. Huang, and M. Zhou. Low-quality product review detection in opinion summarization. In *Proceedings of the Joint Conference on EMNLP-CoNLL*, 2007.
- [15] L. Michael and J. Otterbacher. Write like i write: Herding in the language of online reviews. In *ICWSM*, 2014.
- [16] S. Mudambi. What makes a helpful online review? a study of customer reviews on amazon.com. *MIS Quarterly*, 34:185–200, 2010.
- [17] M. Ott, Y. Choi, C. Cardie, and J. T. Hancock. Finding deceptive opinion spam by any stretch of the imagination. In *Proceedings of the ACL Human Language Technologies*, 2011.
- [18] M. Petrescu, K. OLeary, D. Goldring, and S. B. Mrad. Incentivized reviews: Promising the moon for a few stars. *Journal of Retailing and Consumer Services*, 2017.
- [19] D. Qiao, S.-Y. Lee, A. Whinston, and Q. Wei. Incentive provision and pro-social behaviors. In *Proceedings of the Hawaii International Conference on System Sciences*, 2017.
- [20] ReviewMeta.com. Analysis of 7 million amazon reviews, <https://goo.gl/CPzHpB>, 2016.
- [21] K. Shyong, D. Frankowski, J. Riedl, et al. Do you trust your recommendations? an exploration of security and privacy issues in recommender systems. *Emerging Trends in ICS*, 2006.
- [22] J. Wang, A. Ghose, and P. Ipeirotis. Bonus, disclosure, and choice: What motivates the creation of high-quality paid reviews? In *Proceedings of the International Conference on Information Systems*, 2012.
- [23] Z. Xie and S. Zhu. Appwatcher: Unveiling the underground market of trading mobile app reviews. In *Proceedings of the ACM Conference on Security & Privacy in Wireless and Mobile Networks*, 2015.

TABLE III
THE EVALUATION OF MLP CLASSIFIER IN DETECTING EIRS.

| | Acc. | Rec. | Prec. | F1-score | P-R AUC | AUC |
|--|-----------|-----------|-----------|-----------------------|-----------|-----------|
| Basic | 0.84±0.03 | 0.81±0.01 | 0.78±0.03 | 0.81±0.01 (0.82,0.81) | 0.86±0.01 | 0.81±0.01 |
| Text | 0.88±0.02 | 0.89±0.0 | 0.89±0.02 | 0.89±0.0 (0.89,0.89) | 0.91±0.0 | 0.89±0.0 |
| Basic & Text | 0.92±0.01 | 0.89±0.01 | 0.86±0.02 | 0.89±0.01 (0.9,0.89) | 0.93±0.0 | 0.89±0.01 |
| ==== Category based Training ===== | | | | | | |
| C-Elect. | 0.8 | 0.8 | 0.79 | 0.8 | 0.85 | 0.8 |
| C-Health | 0.87 | 0.86 | 0.84 | 0.86 | 0.9 | 0.86 |
| ==== Only Text Features ===== | | | | | | |
| Tri-grams-char-JustAlpha-noStop MLP | 0.88±0.02 | 0.89±0.0 | 0.89±0.02 | 0.89±0.0 (0.89,0.89) | 0.91±0.0 | 0.89±0.0 |
| Tri-grams-char-LOW-noStop MLP | 0.89±0.02 | 0.88±0.01 | 0.88±0.02 | 0.88±0.0 (0.88,0.88) | 0.91±0.01 | 0.88±0.01 |
| Tri-grams-char-noStop MLP | 0.89±0.01 | 0.89±0.0 | 0.87±0.02 | 0.89±0.01 (0.89,0.89) | 0.91±0.0 | 0.89±0.0 |
| Tri-grams-char-LOW MLP | 0.89±0.02 | 0.89±0.0 | 0.87±0.03 | 0.89±0.0 (0.89,0.89) | 0.91±0.0 | 0.89±0.0 |
| Tri-grams-char-JustAlpha MLP | 0.89±0.01 | 0.89±0.0 | 0.88±0.02 | 0.89±0.0 (0.9,0.89) | 0.91±0.0 | 0.89±0.0 |
| Tri-grams-char-noStop MLP | 0.89±0.01 | 0.89±0.0 | 0.88±0.01 | 0.89±0.0 (0.89,0.89) | 0.91±0.0 | 0.89±0.0 |
| ==== EIRs to normal ratio in Training set ===== | | | | | | |
| MLPC-all-JustAlpha-noStop-0.1 | 0.94±0.01 | 0.93±0.0 | 0.99±0.01 | 0.96±0.0 (0.44,0.96) | 0.97±0.01 | 0.66±0.06 |
| MLPC-all-JustAlpha-noStop-0.2 | 0.94±0.03 | 0.9±0.01 | 0.94±0.05 | 0.94±0.01 (0.68,0.94) | 0.96±0.01 | 0.81±0.07 |
| MLPC-all-JustAlpha-noStop-0.3 | 0.92±0.03 | 0.89±0.01 | 0.95±0.03 | 0.93±0.01 (0.74,0.93) | 0.95±0.01 | 0.83±0.05 |
| MLPC-all-JustAlpha-noStop-0.4 | 0.9±0.02 | 0.89±0.01 | 0.95±0.03 | 0.93±0.01 (0.8,0.93) | 0.94±0.01 | 0.85±0.03 |
| MLPC-all-JustAlpha-noStop-0.5 | 0.92±0.01 | 0.89±0.01 | 0.86±0.02 | 0.89±0.01 (0.9,0.89) | 0.93±0.0 | 0.89±0.01 |
| MLPC-all-JustAlpha-noStop-0.6 | 0.91±0.03 | 0.88±0.01 | 0.9±0.04 | 0.9±0.01 (0.84,0.9) | 0.94±0.01 | 0.87±0.02 |
| MLPC-all-JustAlpha-noStop-0.7 | 0.92±0.02 | 0.9±0.0 | 0.91±0.02 | 0.91±0.0 (0.87,0.91) | 0.94±0.0 | 0.89±0.01 |
| MLPC-all-JustAlpha-noStop-0.8 | 0.93±0.01 | 0.89±0.01 | 0.86±0.03 | 0.89±0.01 (0.88,0.89) | 0.93±0.0 | 0.89±0.01 |
| MLPC-all-JustAlpha-noStop-0.9 | 0.92±0.01 | 0.89±0.01 | 0.86±0.03 | 0.89±0.01 (0.89,0.89) | 0.93±0.0 | 0.89±0.01 |
| ==== All Features different text techniques ===== | | | | | | |
| MLPC-all-LOW | 0.85±0.17 | 0.82±0.15 | 0.88±0.06 | 0.84±0.08 (0.74,0.84) | 0.86±0.07 | 0.82±0.14 |
| MLPC-all-JustAlpha | 0.89±0.03 | 0.89±0.01 | 0.88±0.03 | 0.88±0.01 (0.89,0.88) | 0.91±0.01 | 0.89±0.01 |
| MLPC-all-noStop | 0.91±0.02 | 0.89±0.01 | 0.87±0.04 | 0.89±0.01 (0.89,0.89) | 0.92±0.01 | 0.89±0.01 |
| MLPC-all-STEM | 0.94±0.01 | 0.88±0.01 | 0.81±0.04 | 0.87±0.02 (0.89,0.87) | 0.92±0.0 | 0.88±0.01 |
| MLPC-all-LOW-JustAlpha | 0.91±0.03 | 0.89±0.01 | 0.86±0.05 | 0.88±0.02 (0.89,0.88) | 0.92±0.01 | 0.89±0.01 |
| MLPC-all-LOW-noStop | 0.93±0.03 | 0.89±0.01 | 0.84±0.04 | 0.88±0.01 (0.89,0.88) | 0.93±0.0 | 0.89±0.01 |
| MLPC-all-LOW-STEM | 0.92±0.02 | 0.89±0.01 | 0.86±0.04 | 0.89±0.01 (0.89,0.89) | 0.92±0.01 | 0.89±0.01 |
| MLPC-all-LOW-JustAlpha-noStop | 0.89±0.01 | 0.89±0.0 | 0.9±0.02 | 0.89±0.0 (0.89,0.89) | 0.92±0.0 | 0.89±0.0 |
| MLPC-all-LOW-JustAlpha-STEM | 0.91±0.02 | 0.89±0.0 | 0.87±0.02 | 0.89±0.0 (0.89,0.89) | 0.92±0.0 | 0.89±0.0 |
| MLPC-all-LOW-noStop-STEM | 0.91±0.04 | 0.89±0.01 | 0.86±0.05 | 0.88±0.01 (0.89,0.88) | 0.92±0.01 | 0.89±0.01 |
| MLPC-all-JustAlpha-noStop | 0.92±0.01 | 0.89±0.01 | 0.86±0.02 | 0.89±0.01 (0.9,0.89) | 0.93±0.0 | 0.89±0.01 |
| MLPC-all-JustAlpha-STEM | 0.92±0.02 | 0.89±0.01 | 0.86±0.03 | 0.89±0.01 (0.89,0.89) | 0.92±0.01 | 0.89±0.01 |
| MLPC-all-JustAlpha-noStop-STEM | 0.9±0.02 | 0.89±0.01 | 0.89±0.02 | 0.89±0.0 (0.89,0.89) | 0.92±0.01 | 0.89±0.01 |
| MLPC-all-noStop-STEM | 0.92±0.04 | 0.88±0.01 | 0.84±0.07 | 0.88±0.02 (0.89,0.88) | 0.92±0.01 | 0.88±0.01 |
| MLPC-all-LOW-JustAlpha-noStop-STEM | 0.92±0.03 | 0.88±0.01 | 0.85±0.05 | 0.88±0.02 (0.89,0.88) | 0.92±0.01 | 0.88±0.01 |
| regression | 0.74±0.0 | 0.75±0.0 | 0.77±0.0 | 0.75±0.0 (0.74,0.75) | 0.81±0.0 | 0.75±0.0 |
| regression-all | 0.88±0.01 | 0.87±0.01 | 0.86±0.01 | 0.87±0.01 (0.87,0.87) | 0.9±0.01 | 0.87±0.01 |
| SVM | 0.77±0.0 | 0.75±0.0 | 0.73±0.0 | 0.75±0.0 (0.76,0.75) | 0.82±0.0 | 0.75±0.0 |