Latif, Aneeka
IST 736
HW 5

## Introduction

Living in San Francisco, it's almost as if new restaurants pop up every day. The best way to hear about these restaurants are either through the news or restaurant related apps such as Yelp! that serve as a database and promote restaurants based on user reviews. The average user expects that these reviews to be true renditions of real experiences but unfortunately this is not the case.

Yelp especially has been in hot water over the years due to accusations of untrue positive reviews. Yelp has since taken some action on preventing fake reviews or extremely polarized reviews (Hawkins, 2018[1]) but will not be able to catch all instances of false reviews, for example, in the case of employees boosting their own ratings. This is unfortunate because employees have very clear incentives to want good online postings, they could have promotions depending on their perceived customer service or could be coerced into doing so by managers and employers. Naturally, the employer's incentive in having good reviews, or rather, not having poor reviews is directly correlated with the amount of foot traffic and customers that will enter their establishment; for the employer this is free advertisement.  Users are not looking for poor experiences, they rely on word of mouth and word-of-the-internet to discover must-see eateries. Once a user's read a bad review, it's highly unlikely for them to proceed to dine there.



Figure 1: Poor restaurant review for a bar in Los Angeles          Figure 2: Great restaurant review for a diner in Brooklyn

## Analysis and Models

### About the Data

The dataset is a csv file comprised of 25 yelp reviews. Each row in the file consists of the full user review in text format, review_id, user_id, business_id, star_rating, date, and integer values representing how many users marked a particular review as useful, funny, or cool, respectively. The dataset was reduced significantly to expedite annotations from human raters. The text data was isolated for submission to Amazon Mechanical Turk.

### Crowdsourcing Annotations

[1] Hawkins, Joy; Yelp vs Google: How they deal with fake reviews

Latif, Aneeka
IST 736
HW 5

Amazon Mechanical Turk is a cost effective and time effective method of crowd sourcing annotations. Rather than annotating a data set yourself, you can submit a dataset to AMT, split up the dataset into smaller sets and pay others to do the labeling on your behalf. The only caveat is making sure the annotators provide quality work and aren't spamming or rushing through the set too quickly.

| | review_id | user_id | business_id | stars | useful | funny | cool | text |
|---|---|---|---|---|---|---|---|---|
| 1 | Q1sbwvVQXV27348PgoKj4Q | hG7b0MtEbXx5QzbzE6C_VA | ujmEBvifdJM6h6RLv4wQlg | 1 | 6 | 1 | 0 | Total bill for this horrible service? Over $8Gs. These crooks actually had the ne |
| 2 | GJXCdrlo3ASJOqKeVWPi6Q | yXQM5uF2j56es16SJzNHfg | NZnhc2sEQy3RmzKTZnqtwQ | 5 | 0 | 0 | 0 | I "adore" Travis at the Hard Rock's new Kelly Cardenas Salon! I'm always a fa Travis's greets you with his perfectly green swoosh in his otherwise perfectly st Next Travis started with the flat iron. The way he flipped his wrist to get volume Travis, I will see you every single time I'm out in Vegas. You make me feel bea |
| 3 | 2TzJDVDEuAW6MR5Vuc1ug | n6-Gk65crPZL6UzitqRm3NYw | WTqjgwHi0bSFevF32_DJVw | 5 | 3 | 0 | 0 | I have to say that this office really has it together, they are so organized and fri |
| 4 | yi0R0Ugj_xUx_NekO-_Qig | dacAiZ6fTM6mqwW5uxkskg | ikCg8xy5Jig_NGPx-MSIDA | 5 | 0 | 0 | 0 | Went in for a lunch. Steak sandwich was delicious, and the Caesar salad had a Drink prices were pretty good. The Server, Dawn, was friendly and accommodating. Very happy with her. In summation, a great pub experience. Would go again! |
| 5 | 11a8sVPMUFsaC7_ABRkmtw | ssoyf2_x0EQMed6fgHeMyQ | b1b1eb3uo-w561D0ZfCEiQ | 1 | 7 | 0 | 0 | Today was my second out of three sessions I had paid for. Although my first se |
| 6 | fdiNeiN_hoCxCMy2wTRW9g | w31MKYsNFMrjhWxxAb5wiw | eU_713ac6lTGNO4BegRaww | 4 | 0 | 0 | 0 | I'll be the first to admit that I was not excited about going to La Tavolta. Being a My hubby got the crab tortellini and also loved his. I heard "mmmm this is so g So- Do order the calamari and fried zucchini appetizers. Leave out the mussels. If they have the sea bass special, I highly recommend it. The chicken parm an Do make a reservation but still expect to wait for your food. Go with a large gro |

| text |
|---|
| Total bill for this horrible service? Over $8Gs. These crooks actually had the nerve to charge us $69 for 3 pills. I checked online the pills can be had for 19 cents EACH! Avoid Hospital ERs at all costs. |
| I "adore" Travis at the Hard Rock's new Kelly Cardenas Salon! I'm always a fan of a great blowout and no stranger to the chains that offer this service; however, Travis has taken the flawless blowout to a whole new level! Travis's greets you with his perfectly green swoosh in his otherwise perfectly styled black hair and a Vegas-worthy rockstar outfit. Next comes the most relaxing and incredible shampoo -- where you get a full head message th Next Travis started with the flat iron. The way he flipped his wrist to get volume all around without over-doing it and making me look like a Texas pageant girl was admirable. It's also worth noting that he didn't fry my hair -- som Travis, I will see you every single time I'm out in Vegas. You make me feel beauuuutiful! |
| I have to say that this office really has it together, they are so organized and friendly! Dr. J. Phillipp is a great dentist, very friendly and professional. The dental assistants that helped in my procedure were amazing, Jewel and |
| Went in for a lunch. Steak sandwich was delicious, and the Caesar salad had an absolutely delicious dressing, with a perfect amount of dressing, and distributed perfectly across each leaf. I know I'm going on about the salad . Drink prices were pretty good. The Server, Dawn, was friendly and accommodating. Very happy with her. In summation, a great pub experience. Would go again! |
| Today was my second out of three sessions I had paid for. Although my first session went well, I could tell Meredith had a particular enjoyment for her male clients over her female. However, I returned because she did my teeth |
| I'll be the first to admit that I was not excited about going to La Tavolta. Being a food snob, when a group of friends suggested we go for dinner I looked online at the menu and to me there was nothing special and it seemed ove |

Figure 3: original CSV file

Figure 4: sample set with isolated reviews for AMT

The smaller sample set of 25 reviews was uploaded to Amazon Mechanical Turk with payment of $0.05 per annotated task per 5 reviewers. The payment of $0.05 was intended to be small enough that spammers and bad users would not feel incentivized to accept the task but was also the average payment per task at the time of submission. With a total of 125 annotations, the total payment to AMT was around $7.00, including miscellaneous fees. The annotators were limited by the percent of previous tasks that have been accepted/approved at a rate of 75% or greater. Assuming that all annotators are fluent English speakers, no other limitations were deemed necessary for this task. There was also no need for master raters or annotators with some form of advanced knowledge.

Figure 5: Setting up the AMT task

---

[1] Hawkins, Joy; Yelp vs Google: How they deal with fake reviews

Figure 6: Adding a constraint of HIT acceptance > 75% to ensure good annotations

All five annotators were shown every review in the 25 yelp review set and were asked to label the review with a number from 1 to 5 indicating whether the review was Very Positive, Somewhat Positive, Neutral, Somewhat Negative, or Negative.

These labels were chosen based on personal work experience with in annotation. At a previous annotation project at Yahoo, we found that if the labels were too extreme, it would cause unintended bias. On the project's initial launch, we asked annotators to label ads on relevancy with "perfect", "excellent", "good", "bad", or "poor". We found that regardless of age, gender, professional or educational background, annotators refrained from using the "perfect" label. The reasoning was that from an entire dataset, there should only be one or two cases of perfect relevancy. Our results would always show a maximum of two "perfect" labels and a majority of "good" labels. We re-ran that same exact dataset with different labels that have less of a conceptual bias: these words were "relevant", "somewhat relevant", "barely relevant", "not relevant", and "not applicable" to cover all data issues or other unforeseen issues. We found that annotators began to give a much more normalized distribution of labels over the exact same dataset.

---

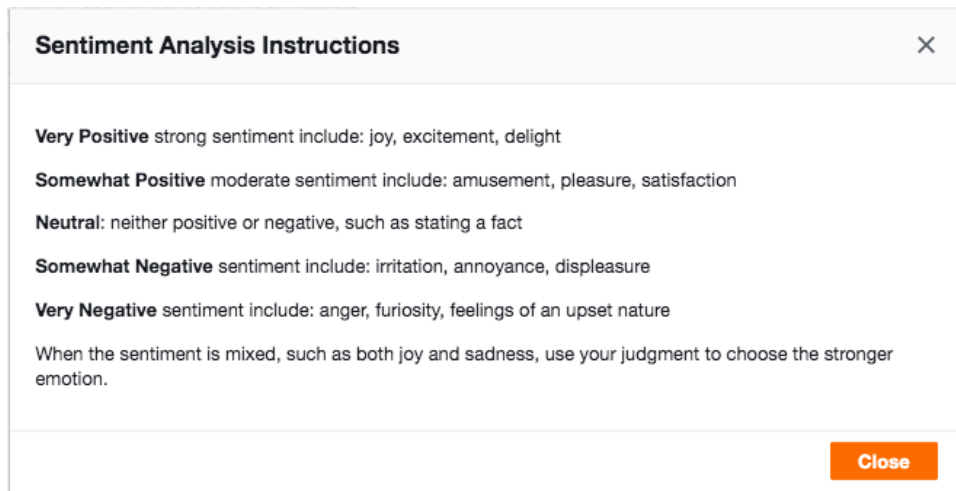[1] Hawkins, Joy; Yelp vs Google: How they deal with fake reviews

**Sentiment Analysis Instructions**                    ✕

**Very Positive** strong sentiment include: joy, excitement, delight

**Somewhat Positive** moderate sentiment include: amusement, pleasure, satisfaction

**Neutral**: neither positive or negative, such as stating a fact

**Somewhat Negative** sentiment include: irritation, annoyance, displeasure

**Very Negative** sentiment include: anger, furiosity, feelings of an upset nature

When the sentiment is mixed, such as both joy and sadness, use your judgment to choose the stronger emotion.

**Close**

Figure 7: This image displays the description of each label chosen for the annotations



**What sentiment does this text convey?**

Total bill for this horrible service? Over $8Gs. These crooks actually had the nerve to charge us $69 for 3 pills. I checked online the pills can be had for 19 cents EACH! Avoid Hospital ERs at all costs.

**Select an option**

| | |
|---|---|
| Very Positive | 1 |
| Somewhat Positive | 2 |
| Neutral | 3 |
| Somewhat Negative | 4 |
| Very Negative | 5 |

Submit

Figure 8:  This image displays a sample task/HIT an annotator would see. Each HIT includes a question to answer, a text snippet, and a label to choose for each task.

The dataset was submitted to Amazon Mechanical Turk with a time limit of 3 hours for each annotator; annotators were able to complete the dataset in a single day. Each of 5 annotators annotated 25 tasks, which resulted in 125 ratings overall. Amazon Mechanical Turk compiled all of these ratings by task ID and provided a downloadable csv file for further analysis.

---

[1] Hawkins, Joy; Yelp vs Google: How they deal with fake reviews

Figure 9: status bar for AMT task



Figure 10: Snippet of results

## Cohen's Kappa

Cohen's Kappa is a measurement of inter-rater agreement; this measures how likely it is for two raters to come to the same decision based on true agreement versus whether two raters came to the decision by chance. Kappa values closer to one indicate strong intentional agreement and kappa values closer to 0 indicate agreement due to chance. The Kappa coefficient is thought to be a much better test than raw agreement because factoring in 'chance' can be thought of as a form of normalization. Cohen's Kappa requires data be in a categorical form, so the sentiment labels were not converted to their numeric counterparts.

## Results

The pairwise Kappa scores were calculated for each pair of annotators. Annotators 2 and 3 seem to have moderately strong agreement with their labels. Annotators 1 and 2 also have moderately strong agreement. The pairs of Annotators 2 and 4 and annotators 3 and 4 have weak agreement but arguably have more intentional agreement than agreement by chance. Annotator 5 has the least agreement with any other annotator which can imply that this annotator may either have not understood the task, or, was randomly choosing labels to finish the task and get paid quicker.

The overall average agreement for all the annotators is 0.236. Surprisingly, despite annotator 5 having the least agreement with the rest of the annotators, removing annotator 5's labels will not significantly increase the average kappa score. The average kappa score without including annotator 5 is only 0.239.

| | Annotator 1 | Annotator 2 | Annotator 3 | Annotator 4 | Annotator 5 |
|---|---|---|---|---|---|
| Annotator 1 | | 0.382 | 0.170 | 0.064 | 0.174 |

[1] Hawkins, Joy; Yelp vs Google: How they deal with fake reviews

| | | | | | |
|---|---|---|---|---|---|
| Annotator 2 | 0.382 | | 0.460 | 0.314 | 0.155 |
| Annotator 3 | 0.170 | 0.460 | | 0.284 | 0.145 |
| Annotator 4 | 0.064 | 0.314 | 0.284 | | 0.215 |
| Annotator 5 | 0.174 | 0.155 | 0.145 | 0.215 | |

Figure 11: confusion matrix of Cohen's Kappa inter-rater agreement on a scale between 0 and 1

Average kappa score = 0.236

Average kappa (with annotator 5 removed) = .239

**Conclusion**

While it may be quick and very convenient to crowd source annotations, it is not without risk. Annotators need to be audited to make sure that labels are correct and valid. One method of checking agreement between annotators is to use Cohen's Kappa which will calculate the percent of true agreement while also taking into consideration the probability of agreement due to chance.

Human annotators fare better than machine learning algorithms in that they are able to process nuances more easily; for example: a model may understand the phrase "this is awfully good" as a negative statement whereas a human annotator can more easily overlook the negation to understand the positive intent. Human annotation can be costly, however, so the most optimal method for model training may be a split between algorithmic prediction and human annotation for validation.

Poor labeling will ultimately affect both the restaurant owners and potential customers. Assuming that companies like Yelp use this form of annotation to rank/remove false reviews, if the annotation data is not correct, the algorithm will not properly flag all the negative reviews and may also incorrectly flag positive reviews for removal. This will result in a poor portrayal of the restaurant and can discourage future clientele from patronizing the restaurant.

[1] Hawkins, Joy; Yelp vs Google: How they deal with fake reviews