

Bias and Reciprocity in Online Reviews: Evidence From Field Experiments on Airbnb

(Preliminary - Do not cite without permission)

Andrey Fradkin^{*1}, Elena Grewal^{†2}, David Holtz^{‡2}, and Matthew Pearson^{§2}

¹MIT Sloan School of Management and Airbnb, Inc.

²Airbnb, Inc.

July 15, 2015

Abstract

Reviews and other evaluations are used by consumers to decide what goods to buy and by firms to choose whom to trade with, hire, or promote. However, because potential reviewers are not compensated for submitting reviews and may have reasons to omit relevant information in their reviews, reviews may be biased. We use the setting of Airbnb to study the determinants of reviewing behavior, the extent to which reviews are biased, and whether changes in the design of reputation systems can reduce that bias. We find that reviews on Airbnb are generally positive and informative. 97% of guests privately and anonymously report having positive experiences and 74% of guests submit a five out of five star overall rating. Furthermore, when guests do not recommend a listing, this is reflected in a lower than five star rating over 90% of the time. We use the results from two field experiments intended to reduce bias to show that non-reviewers tend to have worse experiences than reviewers and that strategic reviewing behavior occurred on the site, although the aggregate effect of the strategic behavior was relatively small. We then model three mechanisms causing bias and show that they decrease the rate of reviews of negative experiences by .65 percentage points.

We are grateful to Jon Levin, Liran Einav, Caroline Hoxby, Shane Greenstein, Ramesh Johari, Mike Luca, Chris Dellarocas, John Horton, Chiara Faronato, Jeff Naecker, Fred Panier, and seminar participants at Microsoft, Ebay, and the CODE Conference for comments. Note: The views expressed in this paper are solely the authors' and do not necessarily reflect the views of Airbnb Inc.

*Primary Author: fradkina@nber.org

†Primary Experiment Designer: elena.grewal@airbnb.com

‡dave.holtz@airbnb.com

§matthew.pearson@airbnb.com

1 Introduction

Reviews and other evaluations are used by consumers to decide what goods to buy and by firms to choose whom to trade with, hire, or promote. These reviews are especially important for online marketplaces (e.g. Ebay, Amazon, Airbnb, and Etsy), where economic agents often interact with new trading partners who provide heterogeneous goods and services.¹ However, potential reviewers are not compensated for submitting reviews and may have reasons to omit relevant information when reviewing. Therefore, basic economic theory suggests that accurate reviews constitute a public good and are likely to be under-provided (Avery et al. (1999), Miller et al. (2005)). As a result, the distribution of evaluations for a given agent may not accurately represent the outcomes of that agent’s previous transactions. The presence of this review bias may consequently reduce market efficiency. For example, it may cause agents to engage in suboptimal transactions (Horton (2014), Nosko and Tadelis (2014)) or to leave the platform altogether.

We study the determinants of reviewing behavior, the informational content of online reviews, and the effects of changes in the design of reputation systems. The setting of this paper is Airbnb, a large online marketplace for accommodations. Reputation is thought to be particularly important for transactions on Airbnb because guests and hosts interact in person, often in the primary home of the host.² As in many other marketplaces, reviews are predominantly positive.³ Over 70% of the guests in our sample submit a five out of five star rating for their hosts. These positive ratings may genuinely reflect positive experiences of guests or they may be the result of biases in the reputation system.

We use proprietary data from Airbnb to show that, although bias does exist, this positive review rate typically reflects positive experiences. Our empirical strategy relies on the fact that Airbnb’s review system solicits anonymous recommendations which are never displayed on the site in addition to the publicly displayed review text and ratings. Reviewers should have no reason to misreport their experience in this anonymous recommendation. We find that over 97% of reviewers recommend their counterparty and that when they do not recommend their counterparty, they leave a lower than five star ratings over 90% of the time. Furthermore, public ratings are correlated with customer service calls and return rates of guests, indicating that these ratings are consistent with actual experiences.

Although we find that reviews typically reflect private beliefs, several forms of bias that may still exist in the review system. We describe a theoretical framework for measuring bias and use experiments as well as non-experimental analysis to study its magnitude. In our theoretical framework, there are two conditions under which reviews can be biased.⁴ First, those that review an agent might differ systematically in their experiences from those that

¹There is a large literature studying the effects of reputation scores on market outcomes. Pallais (2014) uses experiments to show that reviews affect demand for workers on Odesk. Cabral and Hortaçsu (2010) use panel data to show that reputation affects exit decisions by firms on Ebay. Luca (2013) shows that Yelp reputation has especially large effects on non-chain restaurants.

³For example, Thomas Friedman wrote the following in the New York Times: “Airbnb’s real innovation — a platform of ‘trust’ — where everyone could not only see everyone else’s identity but also rate them as good, bad or indifferent hosts or guests. This meant everyone using the system would pretty quickly develop a relevant ‘reputation’ visible to everyone else in the system.”

⁴For example, Horton (2014) shows that 91% of ratings on Odesk in 2014 were four or five (out of five) stars.

do not review an agent. Second, reviewers might not reveal their experiences in the public review. The extent of bias is a function of the utility of reviewing and the design of the reputation system. For example, reviewers who care about their own reputation, may be afraid of retaliation by the counterparty and may consequently inflate their reviews. Changes in the reputation system that remove the possibility of retaliation would then make reviews more honest and increase review rates. In our model, the best review system occurs if each agent submits an honest report after a transaction.

We use our first field experiment, described in more detail in [section 4](#), to study selection into reviewing. The experimental treatment offers a \$25 coupon in exchange for a review. The treatment increases review rates by 6.4 percentage points and decreases the share of those reviews that are five stars by 2.1 percentage points. Furthermore, the coupon does not cause guests to change their reviewing style as measured by prior reviewing behavior, suggesting that the effect of the experiment is primarily due to selection.

The second condition for review bias occurs when reviewers do not reveal their experiences in the review. We show that this misrepresentation does occur in the data. For example, 6% of guests who anonymously answered that they would not recommend their host nonetheless submitted a public review with a five star rating. One possible reason for this misrepresentation is strategic behavior on behalf of reviewers. For example, [Cabral and Hortaçsu \(2010\)](#) and [Saeedi et al. \(2015\)](#) show that when Ebay had a two sided review system, over 20% of negative buyer reviews were followed by negative seller reviews, interpreted by the authors as retaliatory. Furthermore, [Bolton et al. \(2012\)](#) provides laboratory evidence that a system in which reviews are hidden until both parties submit a review (“simultaneous reveal”) reduces retaliation and makes markets more efficient.

We document the first experimental test of a simultaneous reveal mechanism in an actual online marketplace and use it to test whether misrepresentation is due to the fact that reviewers are strategic. This test was designed by Airbnb to determine the effect of the change in mechanisms. The treatment increased review rates by guests while decreasing the share of five star reviews by 1.6 percentage points. On the host side, the treatment increases review rates by 7 percentage points but does not affect recommendation rates. We show that strategic motives affected reviewing behavior in the control group by showing that the relationship between first and second ratings changes due to the experiment. Our results differ from the laboratory results in [Bolton et al. \(2012\)](#) in two ways. First, our experimental treatment effect on five star ratings of -1.5 percentage points is smaller than their -7.7 percentage point effect found in a laboratory setting. Second, our simultaneous reveal mechanism experiment increases review rates while the same mechanism caused reductions in review rates in the lab. Following the experiment, Airbnb released the simultaneous review treatment to all users.

Mismatch between public and private reviews occurs even in the simultaneous reveal treatment group. In [section 6](#) we use non-experimental evidence to study several explanations for this mismatch. We find that mismatch between public and private ratings in the cross-section is predicted by property type (entire home or a room in a home) and host type

⁴There is also considerable evidence about fake promotional reviews, which occur when firms post reviews either promoting themselves or disparaging competitors (see [\(Mayzlin et al., 2014\)](#) for a recent contribution). Promotional reviews are likely to be rare in our setting because a transaction is required before a review can be submitted.

(multi-listing host or casual host). We use two distinct identification strategies to show that the coefficients on these characteristics likely represent causal effects.

First, we compare guest reviewing behavior in cases when a given host sometimes rents out her entire place and other times just a room. We find that guests to the private room are more likely to submit a four or five star rating when they do not recommend the listing. Second, we consider cases when a host who was once a casual host became a multi-listing host. We find that the rate of mismatch decreases when the host becomes a multi-listing host.

We hypothesize that these effects occur because buyers and sellers sometimes have a social relationship. For example, guests who rent a room within a property may talk to their hosts in person. Alternatively, guests may feel more empathy for casual hosts rather than multi-listing hosts, who may communicate in a more transactional manner. Social communication can lead reviewers to omit negative comments due to two reasons. First, conversation can cause buyers and sellers to feel empathy towards each other (Andreoni and Rao (2011)). This may cause buyers to assume that any problem that occurs during the trip is inadvertent and not actually the fault of the seller. Second, social interaction may cause buyers to feel an obligation towards sellers because those sellers offered a service and were “nice” (Malmendier and Schmidt, 2012). This obligation can lead buyers to omit negative feedback because it would hurt the seller or because it would be awkward.⁵

Lastly, we conduct a quantitative exercise to measure the magnitude of bias in online reviews. Bias occurs when a negative experience does not result in a negative public review. We show that bias decreases the rate of reviews with negative text and a non-recommendation by just .86 percentage points. This result is due to the fact that most guests respond that they would recommend their host. However, although the overall bias is small, when negative guest experiences do occur, they are not captured in the review text 56% of the time. We find that most of this effect is caused by sorting bias and the fact that not everyone reviews. This suggests that inducing additional reviews and displaying data on non-reviews can increase market efficiency.

Empirical Context and Related Literature:

Our empirical strategy has at least three advantages over the prior literature on bias in online reviews. First, we conduct two large field experiments that vary the incentives of reviewers on Airbnb. This allows us to credibly identify the causal effects of changes to review systems. Second, we use proprietary data which is observed by Airbnb but not by market participants. This gives us two pieces of information, transactions and private review information, which are typically not used by prior studies. We can use this data to study selection into reviewing and differences between the publicly submitted review and the privately reported quality of a person’s experiences. Lastly, Airbnb (along with Uber, Taskrabbit, Postmates, and others) is a part of a new sector, often referred to as the “Sharing Economy”, which facilitates the exchange of services and underutilized assets between buyers and semi-professional sellers.

⁵Airbnb has conducted surveys of guests who did not submit a review asking why they did not submit one. Typical responses include: “Our host made us feel very welcome and the accommodation was very nice so we didn’t want to have any bad feelings”. “I also assume that if they can do anything about it they will, and didn’t want that feedback to mar their reputation!”

There has been relatively little empirical work on this sector.⁶

Other evidence about Airbnb reviews comes from comparisons with hotel reviews. [Zervas et al. \(2015\)](#) compare the distribution of reviews for the same property on both TripAdvisor and Airbnb and shows that ratings on Expedia are lower than those on Airbnb by an average of at least .7 stars. More generally, the rate of five star reviews is 31% on TripAdvisor and 44% on Expedia ([Mayzlin et al. \(2014\)](#)) compared to 75% on Airbnb. This difference in ratings has led some to conclude that two-sided review systems induce bias in ratings. Our analysis suggests that the five star rate on Airbnb would be substantially higher than 44% even if the three forms of bias that we consider are removed.

There are other potential explanations for the observed differences in ratings distributions between platforms. For example, a much lower share of bookers submit a review on Expedia than on Airbnb.⁷ This may lead reviews on Expedia to be negatively biased if only guests with extreme experiences submit reviews. Alternatively, guests on Airbnb and guests of hotels may have different expectations when they book a listing. A particular listing may justifiably receive a five star rating if it delivered the experience that an Airbnb guest was looking for at the transaction price, even if an Expedia guest would not have been satisfied.⁸

Numerous studies have proposed theoretical reasons why bias may occur but most of the evidence on the importance of these theoretical concerns is observational or conducted in a laboratory setting. For example, [Dellarocas and Wood \(2007\)](#) use observational data from Ebay to estimate model of reviewing behavior.⁹ They show that buyers and sellers with mediocre experiences review fewer than 3 percent of the time. Although our experimental results confirm that mediocre users are less likely to review, the selection is less severe. [Nosko and Tadelis \(2014\)](#) show that Ebay’s search algorithms create better matches when they account for review bias using a sellers Effective Positive Percentage (EPP), the ratio of positive reviews to transactions (rather than total reviews). We provide the first causal evidence that buyers who dont review have worse experiences and, by doing so, provide support for using the EPP metric.

Our coupon intervention reduced bias, but, because coupons are expensive and prone to manipulation, this intervention is not scalable. [Li and Xiao \(2014\)](#) propose an alternative way to induce reviews by allowing sellers to offer guaranteed rebates to buyers who leave a review. However, [Cabral and Li \(2014\)](#) show that rebates actually induce reciprocity in buyers and increase the bias in reviews.

There are other potential problems with review systems which we do not study. Reviews may be too coarse if many types experiences are considered by guests to be worthy of five stars. Another potential problem is that reviewers may react in response to existing reviews (e.g. [Moe and Schweidel \(2011\)](#) and [Nagle and Riedl \(2014\)](#)). Because reviewers

⁶Although see recent contributions by [Fradkin \(2014\)](#) about Airbnb and [Cullen and Farronato \(2015\)](#) about Taskrabbit.

⁷A rough estimate of review rates on Expedia can be derived as follows. Expedia had approximately \$30 billion in bookings in 2012 and approximately 1 million reviews (<http://content26.com/blog/expedias-emily-pearce-user-reviews-rule-the-roost/>). If trips have an average price of \$1000 then the review rates on Expedia are around 3%. In comparison, review rates on Airbnb are over 70%.

⁸Below, we list three other reasons why the distribution of reviews on Airbnb and hotel review sites may differ. One, the price a given listing charges on the two sites may be different. Two, TripAdvisor in particular is prone to fake reviews which tend to deflate overall ratings ([Mayzlin et al. \(2014\)](#)). Three, low rated listings may be filtered out of the site at different rates on the two sites.

on Airbnb typically enter the review flow through an email or notification, they are unlikely to be reading prior reviews when choosing to submit a review. Lastly, even in an unbiased review system, cognitive constraints may prevent agent from using all of the available review information to make decisions.

There are several parallels between social influence in reviewing behavior and social influence in giving experiments. [Bohnet and Frey \(1999\)](#) use laboratory experiment to show that giving decreases with social distance and ([Sally \(1995\)](#)) shows that giving increases with non-binding communication. Anonymity is another important factor in giving behavior. For example, [Hoffman et al. \(1994\)](#) and [Hoffman et al. \(1996\)](#) find that giving decreases with more anonymity and increases with language suggesting sharing. Since transactions on Airbnb are frequently in person, involve social communication, and are branded as sharing, they represent a real world analogue to the above experiments.

Similarly, [Malmendier et al. \(2014\)](#), [Lazear et al. \(2012\)](#), and [DellaVigna et al. \(2012\)](#) find that when given the choice, many subjects opt-out of giving games. When subjects that opt-out are induced to participate through monetary incentives, they give less than subjects that opt-in even without a payment. We find the same effect with regards to reviews — when those that opt-out of reviewing are paid to review, they leave lower ratings. Our results are therefore consistent with models in which leaving a positive review is an act of giving from the reviewer to the reviewed.

2 Setting and Descriptive Statistics

Airbnb describes itself as a trusted community marketplace for people to list, discover, and book unique accommodations around the world. Since 2008, Airbnb has accommodated over 30 million guests and has listed over one million listings. Airbnb has created a market for a previously rare transaction: the rental of an apartment or part of an apartment in a city for a short term stay by a stranger.

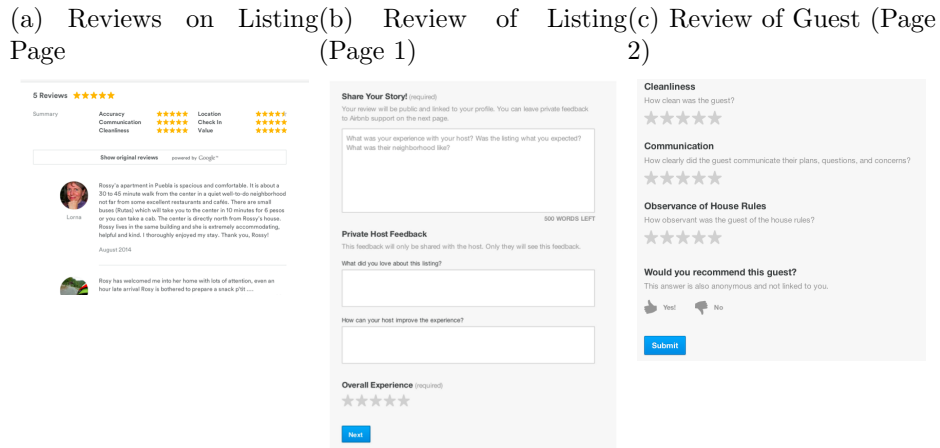
In every Airbnb transaction that occurs, there are two parties - the “Host”, to whom the listing belongs, and the “Guest”, who has booked the listing. After the guest checks out of the listing, there is a period of time (throughout this paper either 14 or 30 days) during which both the guest and host can review each other. Both the guest and host are prompted to review via e-mail the day after checkout. The host and guest also see reminders to review their transaction partner if they log onto the Airbnb website or open the Airbnb app. A reminder is automatically sent by email if a person has not reviewed within a given time period that depends on the overall review period or if the counter-party has left a review.

Airbnb’s prompt for reviews of listings consists of 3 pages asking public, private, and anonymous questions (shown in Figure 1). Guests are first asked to leave feedback consisting of publicly shown text, a one to five star rating,¹⁰ and private comments to the host. The next page asks guests to rate the host in six specific categories: accuracy of the listing compared to the guest’s expectations, the communicativeness of the host, the cleanliness of the listing, the location listing, the value of the listing, and the quality of the amenities provided by the listing. Rounded averages of the overall score and the sub-scores are displayed on each

⁹See [Dai et al. \(2012\)](#) for an interesting structural model which tries to infer restaurant quality and the determinants of reviewing behavior using the sequence of observed Yelp reviews.

listing’s page once there are at least 3 reviews. Importantly, the second page also contains an anonymous question that asks whether the guest would recommend staying in the listing being reviewed.

Figure 1: Review flow on the website



The host is asked whether they would recommend the guest (yes/no), and to rate the guest in three specific categories: the communicativeness of the guest, the cleanliness of the guest, and how well the guest respected the house rules set forth by the host. The answers to these questions are not displayed anywhere on the website. Hosts also submit written reviews that will be publicly visible on the guest’s profile page. Fradkin (2014) shows that, conditional on observable characteristics, reviewed guests experience lower rejection rates by potential hosts. Finally, the host can provide private text feedback about the quality of their hosting experience to the guest and to Airbnb.

2.1 Descriptive Statistics

In this section, we describe the characteristics of reviews on Airbnb. We use data for 59981 trips between May 10, 2014 and June 12, 2014, which are in the control group of the simultaneous reveal experiment.¹¹ The summary statistics for these trips are shown in Table 1. Turning first to review rates, 67% of trips result in a guest review and 72% result in a host review. Furthermore, reviews are typically submitted within several days of the checkout, with hosts taking an average of 2.7 days to leave a review and guests taking an average of 3.3 days. Hosts review at higher rates and review first more often for two reasons. First, because hosts receive inquiries from other guests, they check the Airbnb website more frequently than guests. Second, because hosts use the platform more frequently than guests and rely on Airbnb to earn money, they have more to gain than guests from inducing a positive guest review.

We first consider guest reviews of hosts. 97% of guests who submit a review for a listing, recommend that listing in an anonymous question prompt. This suggests that most guests

¹⁰In the mobile app, the stars are labeled (in ascending order) “terrible”, “not great”, “average”, “great”,

Table 1: Summary Statistics

Reviewer	Reviews	Five Star	Recommends	Overall Rating	Text Classified Positive	First Reviewer	Time to Review (Days)
Guest	0.671	0.741	0.975	4.675	0.779	0.350	4.284
Host	0.715	-	0.989	-	-	0.491	3.667

do have a good experience. Figure 2 shows the distribution of star ratings for submitted reviews. Guests submit a five star overall rating 74% of the time and a four star rating 20% of the time.¹²

Figure 2: Distribution of Guest Overall Ratings of Listings

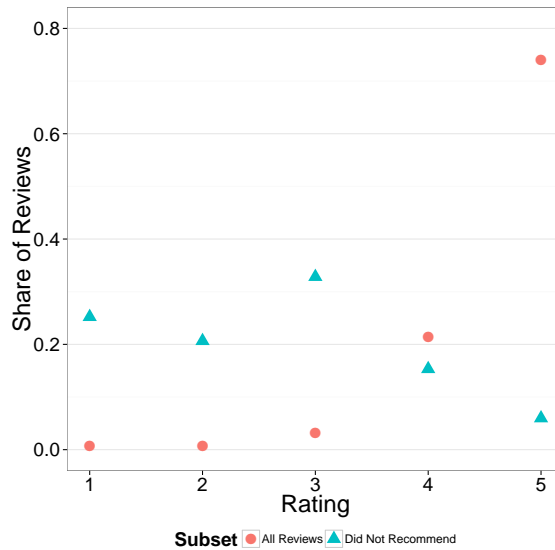


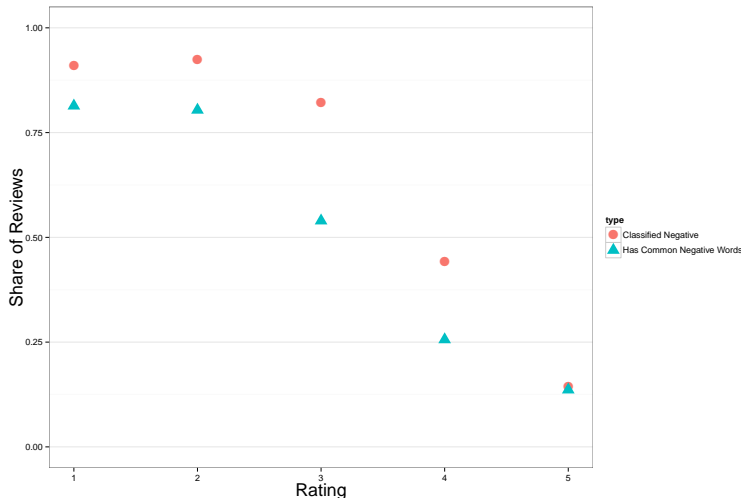
Figure 2 displays the star ratings conditional on whether a guest recommended the listing. As expected, the distribution of ratings for guests who do not recommend is lower than the distribution of ratings for those that do recommend. However, in over 20% of cases where the guest does not recommend the host, the guest submits a four or five star rating. Therefore, guests sometimes misrepresent the quality of their experiences in star ratings. This misrepresentation can occur purposefully or because the guests do not understand the review prompt. Although we have no way to determine whether reviewing mistakes occur, the fact that fewer than 5% of reviewers recommend a listing when they submit a lower than four star rating suggests that guests typically understand the review prompt.¹³

and “fantastic”. The stars are not labeled on the browser during most of the sample period.

¹¹The experiments are randomized at a host level. Only the first trip for each host is included because the experimental treatment can affect the probability of having a subsequent trip. To the extent that better listings are more likely to receive subsequent bookings, these summary statistics understate the true rates of positive reviews in the website.

¹²There is no spike in the distribution for 1 star reviews, as seen on retail sites like Amazon.com. This is likely due to the fact that review rates are much lower for retail websites than for Airbnb.

Figure 3: Prevalence of Negative Text Conditional on Rating



“Classified Negative” refers to the classification by the regularized logistic regression based on the textual features of a review. “Has Common Negative Words” is a binary indicator for whether the review contains a word or phrase that occurs in at least 1% of non-recommended reviews and occurs at least 3 times as frequently in guest reviews with non-recommendations as in guest reviews with five star ratings.

The text of a review is the most public aspect of the information collected by the review system because the text of a review is permanently associated with an individual reviewer. There is also evidence that review text influences consumer decisions even when star ratings are present (Archak et al. (2011)). Review text can contain a variety of information about the quality of a transaction and the characteristics of a product. In this paper, we focus on the sentiment of the text, e.g. whether the text contains only positive information or whether it includes negative phrases and qualifications. We use two approaches to measure sentiment. The first and preferred strategy uses regularized logistic regression, a common technique in machine learning, to classify the review text based on the words and phrases that appear in the text. This approach is described in greater detail in Appendix A.

The most important choice in this procedure is what data to use to “train” (estimate) the model. Our training sample for guest reviews of hosts consists of reviews with five stars, which are labeled as “positive”, and reviews with one or two stars, which are labeled as “negative”. Our training sample for host reviews of guests uses either a non-recommendation or a sub-rating that is lower than 4 stars as a negative label and a recommendation with all sub-ratings as five stars as a positive label. After training the models, we apply them to predict the sentiment in the set of reviews we study for this paper. Both the guest and host review samples are taken from the period before the experiments so that the model training is not affected by the experimental we study. As an alternative classification strategy, we code whether a review had at least one negative word or phrase. A word or phrase is considered negative if it appears three times as frequently in reviews with negative recommendations as reviews with five star ratings and recommendations. The word or phrase must also meet a minimum frequency threshold.

Phrases that commonly show up in negative reviews by guests concern cleanliness, smell,

unsuitable furniture, noise, and sentiment (see [Figure A1](#) for specific examples). In [Figure 3](#) we show the share of reviews with negative text conditional on the rating. Over 90% of 1 and 2 star reviews are classified as negative and these reviews contain the most common negative phrases at over 75% of the time. Three star reviews have text that is classified as negative over 75% of the time. Therefore, we find that guests who are willing to leave negative ratings are also typically willing to leave negative text.

With regards to four star reviews, the results are mixed. Guests only leave negatively classified text 45% of the time. Therefore, the review frequently does not contain information about why the guest left a four star rating. Lastly, even when guests leave a five star rating, they leave negative text approximately 13% of the time. This is due to two reasons. First, even when the experience is not perfect, the listing may be worthy of a five star rating. Guests in that case may nonetheless explain any shortcomings of the listing in the review text. Second, our classifier has some measurement error and this may explain why some of these reviews were classified as negative.

Host reviews of guests are almost always positive. Over 99% of hosts responded that they would recommend a guest. Furthermore, only 14% of reviews by hosts have a category rating that is lower than five stars. These high ratings are present even though the prompt states: “This answer is also anonymous and not linked to you.” We view this as evidence that most guests do not inconvenience their hosts beyond what is expected. In the rare cases when negative reviews by hosts do occur, they contain phrases concerning sentiment, personal communication, money, cleanliness, and damage (see [Figure A2](#) for examples).

3 Theoretical Framework for Review Bias and Its Effects

In this section we describe a simple model of review bias and how reviewing behavior affects market efficiency. Suppose there is a marketplace that brings together buyers and sellers. There are two types of sellers, a high type, H, and a low type, L. The low type sellers always generate a worse experience than the high type sellers. Each seller stays in the market for 2 periods and each period a mass of .5 sellers enter the market, with a probability, μ of being a high type. Sellers choose a price, $p \geq 0$ and their marginal cost is 0. Sellers do not know their type in the first period.

On the demand side, there are $K > 1$ identical buyers each period. Each buyer receives utility u_h if she transacts with a high type and u_l if she transacts with a low type. Furthermore, buyers have a reservation utility $\underline{u} > u_L$ and $\underline{u} \leq \frac{(1-\mu)u_l + .5\mu u_h}{1-.5\mu}$. These assumptions ensure that buyers would not want to transact with low quality sellers but would want to transact with non-reviewed sellers. Lastly, after the transaction, the buyer can review the seller. Buyers can see the reviews of a seller but not the total amount of prior transactions.

After a transaction, buyers can choose whether and how to review sellers. Each buyer, i , has the following utility function for reviewing sellers:

$$\begin{aligned} \kappa_{ih} &= \alpha_i + \beta_i \\ \kappa_{il} &= \max(\alpha_i + \beta_i - \gamma, \beta_i) \end{aligned} \tag{1}$$

where h and l refer to experiences with type H and L sellers respectively. β_i refers to the utility of a review and is potentially influenced by the cost of time, financial incentives to review, and the fear of retaliation from a negative review. α_i refers to the utility of being positive in a review and is influenced by social and strategic reciprocity and the overall preference of individuals to be positive. γ is the disutility from being dishonest. In the case of an interaction with a low quality seller, buyers have to make a choice between misrepresenting their experience, telling the truth, or not reviewing at all.

Observation 1: Both types of sellers can have either no review or a positive review.

Whenever $\alpha_i > 0$, some amount of sorting by type occurs in reviews. This is an implication of the fact that not everyone reviews and that not everyone tells the truth about a low quality experience. One argument against the generality of this observation is that if there were more periods, than all low sellers would eventually get a negative review and would be identified as low quality. In practice, review ratings are rounded to the nearest half a star and sometimes even good sellers get bad reviews. Therefore, buyers still face situations where multiple seller types have the same rating.

Observation 2: The platform knows more information than the buyers do not know about the likely quality of a seller.

Since high type sellers are more likely to be reviewed, a non-review is predictive of the quality of a seller. The platform sees non-reviews while buyers do not and can use that information. Second, platforms often observe private signals associated with a given transaction or review and can use that information to identify the quality of a seller. In our setting, Airbnb can see guests' anonymous recommendations and customer service calls.

Let r_p be equal to the probability that a buyer who transacts with an H seller leaves a positive review, let r_{lp} be the probability that a buyer who transacts with an L seller leaves a positive review, and let r_{lu} be the probability that a buyer who transacts with an L seller leaves a negative review. These probabilities are functions of the utility of review parameters $(\alpha_i, \beta_i, \gamma)$.

All sellers without a negative review transact because buyers' expected utility from the transaction is higher than their reservation utility. The welfare gains from having the marketplace (excluding the disutility from reviewing) are:

$$Welfare = \mu u_h + (1 - \mu)(1 - .5r_{lu})u_l - (1 - .5(1 - \mu)r_{lu})\underline{u} \quad (2)$$

Now suppose that everyone reviewed and did so honestly. This corresponds to $r_p = 1$ and $r_{lu} = 1$. The difference in welfare between the scenario where everyone reviews honestly and the status quo is $.5(1 - \mu)(\underline{u} - u_L)(1 - r_{lu})$.

Therefore, the gain from having a better review system is a function both of the prevalence of bad actors $(1 - \mu)$, the probability guests submit honest negative reports about transactions with low quality listings, r_{lu} , and the disutility from transacting with a low quality seller compared to the utility from the outside option. This analysis justifies our focus on cases where a negative experience is not reported either due to a lack of review, which occurs with

probability $1 - r_{ll} - r_{lp}$, or a misreported review, which occurs with probability r_{lp} . However, because most guests have a positive experience, $1 - \mu$ is close to 0. Therefore, an imperfect review system only causes large welfare losses on the platform when the utility from negative experiences, u_l , is very low relative to the outside option.

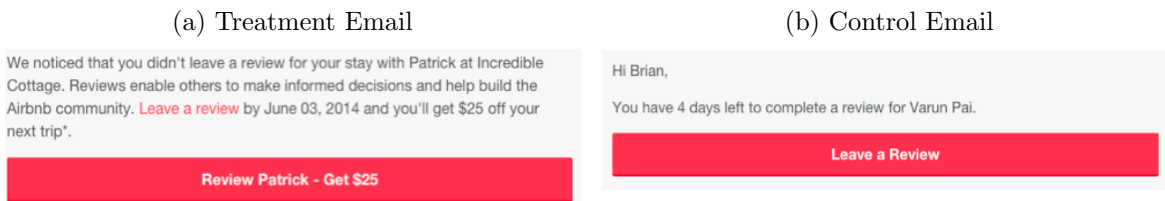
Next, consider the effects of changing the parameters related to reviewing. Increasing β_i can induce additional buyers to review but it does not change their decision to report honestly or dishonestly. Therefore, the welfare gains from increasing β_i come from inducing buyers who were previously not reviewing to honestly review. In our results, this parameter change corresponds to offering a coupon for buyers to review.

Increasing α_i induces additional buyers to review positively and induces truth-tellers to misreport. The welfare change from increasing α_i comes from inducing dishonest reports because only r_{ll} matters for welfare. Tying this to our later empirical results, increasing socially induced reciprocity corresponds to an increase in α_i . Therefore, while socially induced reciprocity increases review rates for high quality sellers, it also increases misreporting for low quality sellers and therefore reduces market efficiency in this model.

4 Sorting Into Reviewing and the Incentivized Review Experiment

In this section we describe an experiment intended to induce additional reviews and use it to study the magnitude of sorting into reviewing based on the quality of experience by reviewers. In the experiment, which was conducted between April and July of 2014, all trips to non-reviewed listings for which the guest did not leave a review within 9 days were assigned to either a treatment group or a control group, each assigned with a 50% probability at a host level. Guests in the treatment group received an email offering a \$25 Airbnb coupon while guests in the control group received a normal reminder email (shown in Figure 4).

Figure 4: Incentivized Review Experiment Emails



The treatment affected the probability of a review and consequently the probability of additional bookings for a listing. This resulted in more trips in the experimental sample to listings in the control group than listings in the treatment group. Therefore, we limit the analysis to the first trip to a listing in the experiment.¹⁴ Appendix C demonstrates that the randomization for this experiment is valid.

Table 2 displays the review related summary statistics of the treatment and control groups in this experiment. First, note that the 23% review rate in the control group is smaller than

¹⁴We plan to investigate what occurred in subsequent trips in follow-up work.

Table 2: Summary Statistics: Incentivized Review Experiment

	Control		Treatment	
	Guest	Host	Guest	Host
Reviews	0.257	0.626	0.426	0.632
Five Star	0.687	-	0.606	-
Recommends	0.773	0.985	0.737	0.986
High Likelihood to Recommend Airbnb	0.731	-	0.708	-
Overall Rating	4.599	-	4.488	-
All Sub-Ratings Five Star	0.458	0.795	0.389	0.805
Responds to Review	0.021	0.051	0.019	0.040
Private Feedback	0.432	0.273	0.439	0.275
Feedback to Airbnb	0.102	0.089	0.117	0.089
Median Review Length (Characters)	345	126	301	128
Negative Sentiment Given Not-Recommend	0.882	0.849	0.930	0.673
Text Classified Positive	0.757	-	0.688	-
Median Private Feedback Length (Characters)	131	95	126	96
First Reviewer	0.072	0.599	0.168	0.570
Time to Review (Days)	18.420	5.864	13.709	5.715
Time Between Reviews (Hours)	292.393	-	215.487	-
Num. Obs.	15470	15470	15759	15759

the overall review rate (67%). The lower review rate is due to the fact that those guests who do not review within 9 days are less likely to leave a review than the average guest. The treatment increases the review rate in this sample by 70% and decreases the share of five star reviews by 11%. The left panel of figure 5 displays the distribution of overall star ratings in the treatment versus the control. The treatment increases the number of ratings in each star rating category. It also shifts the distribution of overall ratings, increasing the relative share of 3 and 4 star ratings compared to the control. The non-public responses of guests are also lower in the treatment, with a 2 percentage point decrease in the recommendation and likelihood to recommend Airbnb rates.

The effect of this experiment on the review ratings might be caused by one of several mechanisms. In Appendix D we show that the effect is not driven by guest characteristics, guest leniency in reviewing, listing characteristics, and fear of retaliation in a host review.

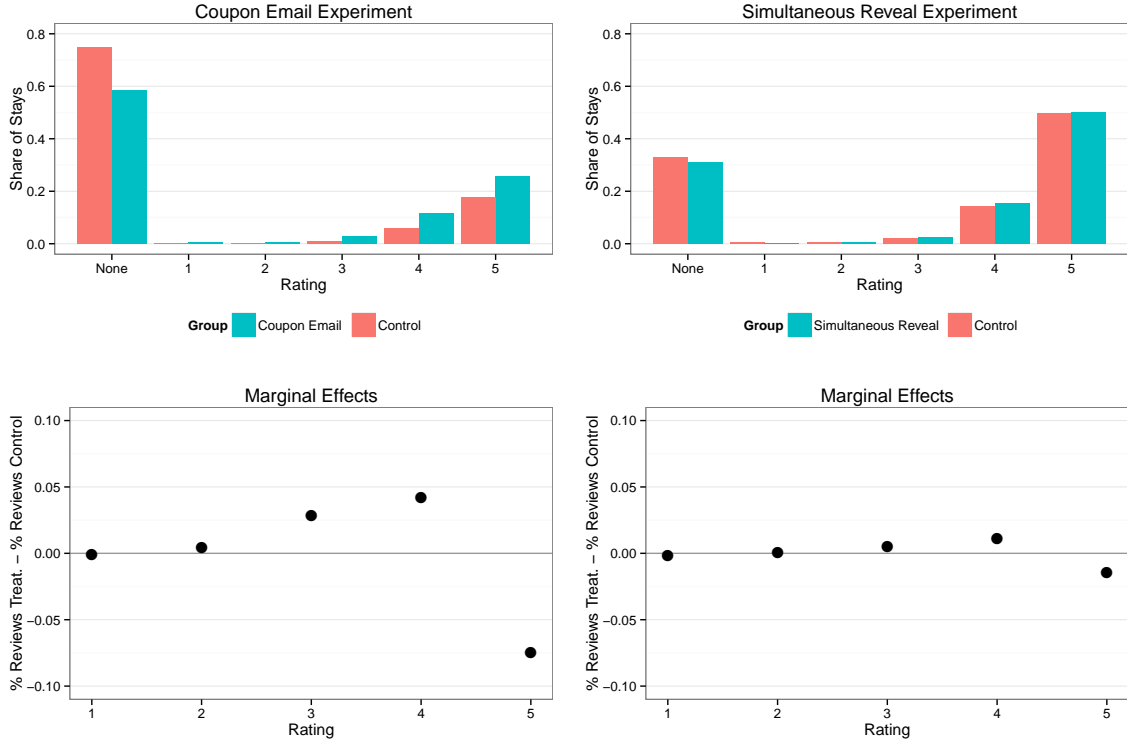
Because only those guests who had not left a review within 9 days are eligible to be in the experiment, the estimated treatment effects do not represent changes to the overall distribution of ratings for non-reviewed listings. We use the following equation to adjust the experimental treatment effects to represent the overall effect on ratings for listings with 0 reviews.

$$e_m = \frac{s_{\leq 9} r_{m, \leq 9} + (s_{ctr} + t_{rev})(r_{m, ctr} + t_m)}{s_{\leq 9} + s_{ctr} + t_{rev}} - \frac{s_{\leq 9} r_{m, \leq 9} + s_{ctr} r_{m, ctr}}{s_{\leq 9} + s_{ctr}} \quad (3)$$

where e_m is the adjusted treatment effect for metric m , s refers to the share of trips in each group, t refers to the experimental treatment effect, and r_m refers to the mean value of a review metric, m . “ ≤ 9 ” refers to the sample of trips where the guest reviews within 9 days, “ ctr ” refers to the control group, and t_{rev} refers to the treatment effect of the experiment on review rates.

Table 3 displays the baseline treatment effects (Column 1) and adjusted treatment effects (Column 2) for this experiment using the sample of trips that were also in the treatment of the subsequent experiment (this sample is chosen for comparability of results). The 17 percentage point treatment effect on review rates in the experiment drops to a 6.4 percentage

Figure 5: Distribution of Ratings - Experiments



point effect when scaled. Because of this scaling, the effect of the experiment is smaller on the overall distribution of reviews than on the distribution of reviews in the experiment. Another reason why there is a difference between columns (1) and (2) is that guests who review after 9 days tend to give lower ratings on average. Therefore, even if the experiment did not change the composition of reviews among those that did not review within 9 days, it would still have an effect on the distribution of ratings by inducing more of these guests to review. In total, the experiment decreases the overall share of five star ratings by 2.1 percentage points and the share of reviews with recommendations by .5 percentage points.

The effects discussed above do not capture the full bias due to sorting because the experiment induced only 6.4% of guests to review, leaving 27% of trips without guest reviews. Our data cannot tell us about the experiences of those non-reviewers. On the one hand, those who do review in the treatment may have even worse experiences than reviewers because they did not bother to take the coupon. Alternatively, those who did not review may have simply been too busy to review or may have cared about a \$25 coupon for Airbnb, especially if they do not plan on using Airbnb again. In column (5) of [Table 3](#) we show the imputed selection effect if non-reviewers had the same behavior as reviewers in the treatment group of the experiment. In this case, there would be a 5.8 percentage point lower five star review rate for previously non-reviewed listings on Airbnb. We view this as conservative estimate of total sorting bias because those not induced by the coupon are likely to have had even worse experiences than those who did take up the coupon.

Although we have documented sorting bias, this bias may not matter much for market

Table 3: Magnitudes of Experimental Treatment Effects

Experiment:	Coupon	Coupon	Sim. Reveal	Sim. Reveal	Coupon
Sample:	Experimental Sample	No Prior Reviews	All Listings	No Prior Reviews	No Prior Reviews
Adjustment:		Effect on Distribution			If Everyone Reviewed
Specification:	(1)	(2)	(3)	(4)	(5)
Reviewed	0.166***	0.064	0.018***	0.008	0.323
Five Star	-0.128***	-0.024	-0.015***	-0.010*	-0.060
Recommends	-0.012	-0.004	-0.001	-0.001	-0.011
Neg. Sentiment	0.135***	0.027	0.020***	0.028***	0.068

Columns (1), (3), and (4) display treatment effects in a linear probability model where the dependent variable is listed in the first column. Column (2) adjusts the treatment effects in column (1) to account for the fact that only guests who had not reviewed within 9 days were eligible for the coupon experiment. Therefore, the treatment effect in column (2) can be interpreted as the effect of the coupon experiment on average outcomes for all trips to non-reviewed listings. Controls for trip and reviewer characteristics include: number of guests, nights, checkout date, guest origin, listing country, and guest experience. The regressions predicting five star reviews, recommendations, and sentiment are all conditional on a review being submitted. “Negative sentiment” is an indicator variable for whether the review text contains one of the phrases identified as negative. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$ (Estimates in Column (2) do not have associated standard errors.)

efficiency if it is constant across listings. In that case, consumers can rationally adjust their expectations that reviews are inflated. However, if sorting differs between listings with similar ratings then even knowledgeable guests may mistakenly book a low quality listing because that listing has especially biased reviews. We demonstrate that there is heterogeneity in bias by comparing the distribution of the difference between the two quality measures: the share of five star reviews out of all reviews (which does not take into account sorting) and the more reliable Effective Positive Percentage (EPP) proposed by [Nosko and Tadelis \(2014\)](#).¹⁵ We use a sample of listings where the average star rating is greater than 4.75, so that the overall star ratings are rounded to 5 for all listings. Although all of these listings have similar average star ratings, their five star rate minus EPP varies greatly, with an interquartile range for the difference of 16% - 27%. (See [Figure A3](#) for a histogram). In [Appendix B](#) we also show that EPP at the time of bookings predicts future ratings. Therefore, we conclude that sorting bias varies across listings and therefore affects market efficiency.

5 The Simultaneous Reveal Experiment

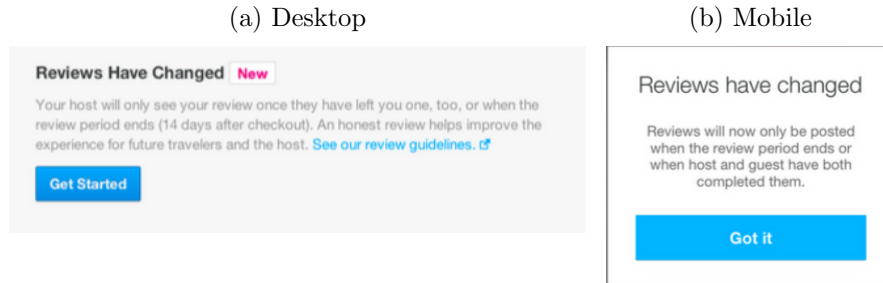
In this section we study the effects of a change in Airbnb’s review system intended to remove strategic retaliation and reciprocation of reviews. Prior to May 8, 2014, both guests and hosts had 30 days after the checkout date to review each other and any submitted review was immediately posted to the website. This allowed for the possibility of the second reviewer to retaliate or reciprocate the first review. To the extent that this retaliation or reciprocation did not accurately reflect the quality of a reviewer’s trip, it made the review system less informative.

The second experiment precludes this strategic reciprocity by changing the timing with which reviews are publicly revealed on Airbnb. Starting on May 8, 2014, Airbnb ran an experiment in which one third of hosts were assigned to a treatment in which reviews were hidden until either both guest and host submitted a review or 14 days had expired (shown

¹⁵EPP is measured in this paper by the share of five star reviews out of all trips.

in Figure 6). Another third of hosts were assigned to a control group where reviews were revealed as soon as they were submitted and there were also 14 days to review.

Figure 6: Simultaneous Reveal Notification



For this analysis we limit the data we use to the first trip to every listing that was in the experiment. We exclude subsequent trips because the treatment may affect re-booking rates which would make the experiment unbalanced. Appendix C documents the validity of our experimental design. Table 4 shows the summary statistics for the treatment and control groups in the “simultaneous reveal” experiment. The treatment increases review rates for guests by 2 percentage points and for hosts by 7 percentage points. The rate of five star reviews by guests decreases by 1.6 percentage points, while the recommendation rate decreases by .4 percentage points. Furthermore, the drop in positive text by guests of 1.5 percentage points mirrors the drop in five star reviews. This suggests that the fear of retaliation had a similar effect on both the averaged ratings and the text in which the reviewer was identifiable. The treatment induced a 6.4 percentage points higher rate of guest suggestions to hosts. (see Table AVI). This increase was present even when conditioning on guest recommendations and star ratings. The relatively larger increase in private feedback suggests that without the fear of retaliation, guests felt they could speak more freely to the hosts about problems with the listing. However, many guests still did not wish to leave negative feedback because they did not want to hurt the hosts, regardless of the possibility of retaliation.

Columns (3) and (4) of Table 3 display the experimental treatment effects on guest reviews when controlling for trip and guest characteristics. Column (3) uses the entire experimental sample while column (4) shows estimates from a sample of previously non-reviewed listings. Of note is that although the experiment has statistically significant effects on reviewing behavior, they are generally smaller than the effects of the coupon. This is evident when comparing column (4) to column (5) of the table. The results from the coupon experiment suggest that eliminating sorting by inducing everyone to review would decrease the rate of five star reviews by 5.8 percentage points, whereas removing strategic motivations only has a 1 percentage point effect. Therefore, sorting is more important for determining the distribution of star ratings than strategic factors.

Turning to hosts, the rate of reviews increases by 7 percentage points, demonstrating that hosts were aware of the experiment and were induced to review. Furthermore, the rate of recommendations by hosts did not significantly change, suggesting that the recommendation is not affected by strategic motivates. However, the text of the submitted reviews does

Table 4: Summary Statistics: Simultaneous Reveal Experiment

	<u>Guest</u>		<u>Host</u>	
	Control	Treatment	Control	Treatment
Reviews	0.671	0.690	0.715	0.787
Five Star	0.741	0.726	-	-
Recommends	0.975	0.974	0.989	0.990
High Likelihood to Recommend Airbnb	0.765	0.759	-	-
Overall Rating	4.675	4.661	-	-
All Sub-Ratings Five Star	0.500	0.485	0.854	0.840
Responds to Review	0.025	0.066	0.067	0.097
Private Feedback	0.496	0.567	0.318	0.317
Feedback to Airbnb	0.106	0.109	0.068	0.072
Median Review Length (Characters)	327	333	147	148
Negative Sentiment Given Not-Recommend	0.861	0.866	0.714	0.743
Text Classified Positive	0.779	0.764	-	-
Median Private Feedback Length (Characters)	131	129	101	88
First Reviewer	0.350	0.340	0.491	0.518
Time to Review (Days)	4.284	3.897	3.667	3.430
Time Between Reviews (Hours)	63.680	47.478	-	-
Num. Obs.	60743	61018	60743	61018

change. The rate of negative sentiment conditional on a non-recommend (calculated using the methodology described in [section 2](#)) increases from 71% to 74%. This suggests that the experiment did have the intended effect of allowing people to be more honest in their public feedback.

5.1 Evidence for Retaliation and Reciprocity

In this section, we use experimental variation to quantify the importance of strategic reviewing on Airbnb. We first test for responses by the second reviewer to the first review. In the control group of the experiment, second reviewers see the first review and can respond accordingly. In the treatment group, second reviewers cannot respond to the content of the first review. If the experiment has an effect, the first review text should have no effect on the second review, conditional on the host’s recommendation. Our specification to test this is:

$$y_{gl} = \alpha_0 t_l + \alpha_1 FNR_{gl} + \alpha_2 FNS_{gl} + \alpha_3 t_l * FNR_{gl} + \alpha_4 t_l * FNS_{gl} + \beta' X_{gl} + \epsilon_{gl} \quad (4)$$

where y_{gl} is a negative review outcome, t_l is an indicator for whether the listing is in the treatment group, FNR_{gl} is an indicator for whether the first reviewer did not recommend, FNS_{gl} is an indicator for whether the first review text contained negative sentiment, and X_{gl} are guest, trip and listing controls.

If guests reciprocate positive first reviews, then the guests in the treatment should leave less positive reviews after a positive review by a host. This response corresponds to α_0 being positive. Second, α_1 should be positive if there is positive correlation between guest and host experiences. Third, if there is retaliation against negative host reviews, α_2 , because negative first review text induces negative second reviews. Moving to the interactions, α_2 should be approximately equal $-\alpha_4$ because second reviewers in the treatment can no longer see the first review. Lastly, we expect that α_3 , the interaction of the non-recommendation

with the treatment to be close to 0. The reason is that second reviewers do not see the recommendation regardless of the experimental assignment.¹⁶

Table 5 displays estimates of Equation 4 for cases when the guest reviews second. Columns (1) - (3) show the estimates for guest non-recommendations, low ratings, and negative sentiment respectively. Turning first to the estimates of α_0 , the effect is a precisely estimated 0 for non-recommendations and positive for the other metrics. This demonstrates that guests do not change their non-public feedback in response to positive host reviews. However, when the first review is positive, guests are more likely to respond with positive ratings and text. Next, we consider the effect of a host’s review having negative sentiment conditional on a non-recommendation. Across the three outcome variables, the coefficients on host negative sentiment range between .55 and .71. Furthermore, the interaction with the treatment is of the opposite sign and of similar magnitude. Therefore, the correlation between first and second negative reviews is driven by retaliation. Therefore, we conclude that guests both retaliate and reciprocate host reviews.

Table 5: Retaliation and Induced Reciprocity - Guest

	Does Not Recommend (1)	Overall Rating < 5 (2)	Negative Text (3)
Treatment	0.002 (0.002)	0.032*** (0.006)	0.036*** (0.006)
Host Negative Sentiment	0.674*** (0.130)	0.705*** (0.124)	0.555*** (0.143)
Host Does Not Recommend	0.133 (0.094)	0.047 (0.109)	0.235* (0.131)
Treatment * Host Negative Sentiment	-0.625*** (0.160)	-0.694*** (0.177)	-0.630*** (0.195)
Treatment * Host Does Not Recommend	-0.014 (0.120)	0.252* (0.151)	0.086 (0.172)
Guest, Trip, and Listing Char. Observations	Yes 17,995	Yes 17,995	Yes 17,995

5.2 Evidence for Fear of Retaliation and Strategically Induced Reciprocity

We now investigate whether first reviewers strategically choose review content to induce positive reviews and to avoid retaliation. Strategic actors have an incentive to omit negative feedback from reviews and to wait until the other person has left a review before leaving

¹⁶There are two complications to the above predictions. First, the experiment not only changes incentives but also changes the composition and ordering of host and guest reviews. If, for example, trips with bad outcomes were more likely to have the host review first in the treatment, then the predictions of the above paragraph may not hold exactly. Second, because we measure sentiment with error, the coefficients on the interaction of the treatment with non-recommendations may capture some effects of retaliation.

a negative review. Because the simultaneous reveal treatment removes this incentive, we expect a higher share of first reviewers to have negative experiences and to leave negative feedback, conditional on having a negative experience. We test for these effects using the following specification:

$$y_{gl} = \alpha_0 t_l + \alpha_1 FNR_{gl} + \alpha_2 FNR_{gl} * t_l + \epsilon_{gl} \quad (5)$$

where y_{gl} is a negative review outcome, t_l is an indicator for whether the listing is in the treatment group and FNR_{gl} is an indicator for whether the reviewer did not anonymously recommend the counter-party. We expect α_0 to be positive because first reviews should be more honest in the treatment. α_1 should be positive because reviewers who do not recommend typically leave negative public reviews. Lastly, α_2 should be positive because reviewers should be more honest in the treatment.

Table 6 displays estimates of Equation 5 for first reviews by hosts. Column (1) shows that hosts are 2.7 percentage points more likely to review first in the treatment. This demonstrates that hosts change their timing of reviews to a greater extent than guests, presumably because they don't have an incentive to threaten retaliation against negative guest reviews. Columns (2) and (3) display the main specification, where y_{gl} is an indicator for the presence of negative sentiment in the host's review text. There is a .2 percentage point increase in the overall rate of negative text in first host reviews. Column (3) shows that this effect is concentrated amongst hosts that do not recommend the guest. The treatment causes hosts to include negative review text an additional 12 percentage points when they do not recommend the guest.

These results demonstrate that hosts are aware of strategic considerations and omit negative feedback from public reviews even if they have a negative experience. Furthermore, the effect is concentrated on hosts who do not recommend their guests. The lack of an effect of the treatment on those hosts who do recommend their guests, suggests that hosts review honestly in those cases. Appendix E discusses the analogous results for guests reviewing first.

Table 6: Fear of Retaliation - Host

	Reviews First (1)	Neg. Sentiment (First)	
		(2)	(3)
Treatment	0.027*** (0.003)	0.002* (0.001)	-0.0005 (0.001)
Does Not Recommend			0.613*** (0.044)
Treatment * Does Not Recommend			0.117** (0.055)
Guest, Trip, and Listing Char. Observations	Yes 121,380	Yes 31,547	Yes 31,547

6 Misreporting and Socially Induced Reciprocity

Reviewers leave conflicting private and public feedback even when there is no possibility of retaliation. In the simultaneous reveal treatment, guests who do not recommend a listing fail to leave negative text 33% of the time and leave four or five star ratings 20% of the time. Similarly, hosts do not leave negative text in 29% of cases when they do not recommend the guest. In this section, we link this misreporting in public reviews to the type of relationship between the guest and host.

Stays on Airbnb frequently involve a social component. Guests typically communicate with hosts about the availability of the room and the details of the check-in. Guests and hosts also often socialize while the stay is happening. This social interaction can occur when hosts and guests are sharing the same living room or kitchen. Other times, the host might offer to show the guest around town or the guest might ask for advice from the host. Lastly, the type of communication that occurs may differ between hosts who are professionals managing multiple listings and hosts who only rent out their own place.

Internal Airbnb surveys of guests who did not leave a review suggest that the social aspect of Airbnb affects reviewing behavior. Guests often mention that it feels awkward to leave a negative review after interacting with a host. For example, one guest said: “I liked the host so felt bad telling him more of the issues.” Second, guests frequently mention that they don’t want the host to feel bad. One respondent said: “I often don’t tell the host about bad experiences because I just don’t want to hurt their feelings”. Third, guests don’t want to hurt the host’s reputation. A typical response is: “My hosts were all lovely people and I know they will do their best to fix the problems, so I didn’t want to ruin their reputations.” Lastly, guests sometimes doubt their own judgment of the experience. For example, one guest claimed that “I think my expectations were too high”. T

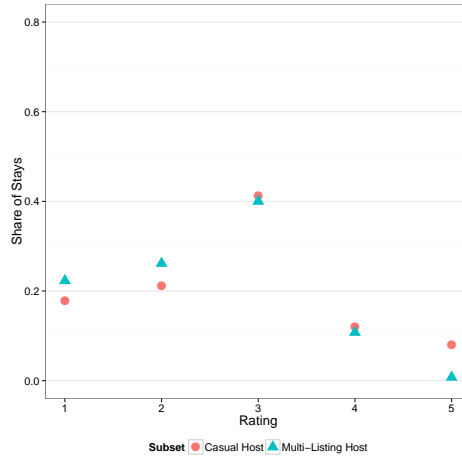
We do not directly observe whether social interaction occurs, but we do observe variables correlated with the degree of social interaction between guest and host. Our first proxy for the degree of social interaction is whether the trip was to a private room within a home or to an entire property. Stays in a private room are more likely to result in social interaction with the host because of shared space. Our second proxy for social interaction is whether the host is a multi-listing host (defined as a host with more than 3 listings). Multi-listing hosts are less likely to interact with guests because they are busy managing other properties and because they typically do not reside in the properties they manage.

Figure 7 plots the distribution of guest ratings conditional on not recommending the host as a function of property type. Guests staying with casual hosts are over 5% more likely to submit a five star overall rating than guests staying with multi-listing managers. That is, even though all guests in the sample would not recommend the listing they stayed at, those staying with multi-listing hosts were more likely to voice that opinion in a review rating. The baseline regression specification we use to study this effect is as follows:

$$y_{glt} = \alpha_0 PR_l + \alpha_1 MLH_l + \alpha_2 R_{glt} + \alpha_3 MLH_l * NR_{glt} + \alpha_4 R_{glt} * NR_{glt} + \beta' X_{gl} + \gamma_g + \epsilon_{gl} \quad (6)$$

where y_{glt} is a negative review by guest g for listing l at time t , PR_l is an indicator for whether the listing is a private room, MLH_l is an indicator for whether the host is a multi-listing host, R_{gl} is a vector of rating indicators, X_{gl} are guest and trip characteristics, and γ_g is a guest fixed effect. If socially induced reciprocity occurs, then α_3 should be negative

Figure 7: Ratings When Guest Does Not Recommend - Simultaneous Reveal



because guests to private rooms should leave less negative feedback and α_4 to be positive because multi-listing hosts induce less reciprocity in guests.

Of course reviews across listing types may differ for reasons other than the degree of social interaction. Different listing or host types may have different qualities and this could cause the rates of misrepresentation to of experiences to change. To control for these factors, we use two forms of variation in the data. First, sometimes a particular property is rented out completely while at other times just a room in that property is rented out. Other than the size of the room, the price, and the degree of social interaction, there should be minimal difference in the quality of the two listings. We add address-specific fixed effects in some specifications to isolate the effect of staying in a private room. Similarly, stays with a multi-listing host may differ for a variety of reasons unrelated to socially induced reciprocity. Therefore, we use variation within listing to study the effects of multi-listing hosts. This variation exists because hosts sometimes start as casual hosts but then expand their operations over time.

Table 7 displays the results of regressions predicting whether a review rating had more than 3 stars. Columns (1) contains a specification with a variety of controls while column (2) adds guest fixed effects. In both specifications, entire properties are 1 percentage point less likely to receive high rating, but the effect goes away if a guest recommends a listing. Similarly, multi-listing hosts are 4.5% less likely to receive high rating when they are not recommended by the guest. Column (3) adds listing fixed effects, using variation in host status over time to identify the effect of a multi-listing host. In this case, reviews of multi-listing hosts are 3.2 percentage points less likely to receive high rating if the guest does not recommend.

Table 8 contains the specifications with address fixed effects. Column (1) shows a regression in which the entire property indicator is not interacted with the recommendation. In this case, there is not difference on average between reviews of entire properties and private rooms at the same location. However, when interactions are added in columns (2) and (3), there is 4.6 percentage point decrease in the probability of high ratings for entire properties

relative to private rooms when there is a non-recommendation. We also conducted the same exercise when the outcome variable was negative sentiment in review text and the results were similar. This evidence confirms that guest's willingness to be honest in reviews is a function of the degree of social interaction they had with the host. Furthermore, these estimates of socially induced reciprocity are likely to be underestimates because even stays at entire properties with multi-listing hosts still sometimes have social interactions.

Table 7: Socially Induced Reciprocity - Star Rating

	Rating > 3		
	(1)	(2)	(3)
Entire Property	-0.011*** (0.001)	-0.013*** (0.002)	
Listing Reviews	-0.0001*** (0.00000)	-0.0001*** (0.00001)	-0.00002 (0.00002)
Checkout Date	-0.000*** (0.000)	-0.000*** (0.000)	-0.000*** (0.000)
Nights	0.0003*** (0.00002)	0.0003*** (0.00003)	0.0001*** (0.00005)
Guests	-0.003*** (0.0001)	-0.001*** (0.0002)	-0.002*** (0.0003)
Customer Support	-0.026*** (0.001)	-0.025*** (0.001)	-0.020*** (0.001)
Total Bookings by Guest	0.0004*** (0.00003)	-0.0002*** (0.0001)	-0.0002** (0.0001)
Price	0.0001*** (0.00000)	0.0001*** (0.00000)	-0.00003*** (0.00001)
Effective Positive Percentage	0.055*** (0.001)	0.055*** (0.001)	-0.009*** (0.001)
No Trips	0.003 (0.008)	0.007 (0.010)	0.028 (0.020)
Person Capacity	-0.001*** (0.0001)	-0.001*** (0.0001)	-0.0001 (0.0004)
Multi-Listing Host	-0.045*** (0.001)	-0.045*** (0.002)	-0.032*** (0.003)
Recommended	0.758*** (0.001)	0.742*** (0.001)	0.691*** (0.002)
Multi-Listing * Recommended	0.031*** (0.001)	0.030*** (0.002)	0.030*** (0.002)
Entire Prop. * Recommended	0.014*** (0.001)	0.015*** (0.002)	0.010*** (0.002)
Guest FE	No	Yes	Yes
Market FE	Yes	Yes	No
Listing FE	No	No	Yes
Observations	2,274,159	2,274,159	2,274,159

The outcome in the above regression is whether the guest's star rating is greater than 3. The estimation is done on all trips between 2012 and 2014 for a 50% sample of guests. *p<0.10, ** p<0.05, *** p<0.01

Table 8: Socially Induced Reciprocity - Address Fixed Effects

	Rating > 3		
	(1)	(2)	(3)
Entire Property	0.0005 (0.002)	-0.046*** (0.005)	-0.046*** (0.007)
Listing Reviews	0.0001** (0.00003)	0.00005 (0.00003)	0.00001 (0.0001)
Checkout Date	-0.000*** (0.000)	-0.000*** (0.000)	-0.000** (0.000)
Nights	0.0001 (0.0001)	0.0001 (0.0001)	0.0001 (0.0001)
Guests	0.001 (0.0005)	-0.0005 (0.0004)	-0.0003 (0.001)
Customer Support	-0.075*** (0.002)	-0.023*** (0.002)	-0.022*** (0.003)
Log(Guest Bookings)	-0.002*** (0.001)	0.002*** (0.0005)	-0.001 (0.001)
Log(Price Per Night)	-0.019*** (0.002)	-0.008*** (0.002)	-0.008*** (0.002)
High LTR			0.037*** (0.002)
Recommends		0.726*** (0.003)	0.734*** (0.005)
Entire Prop. * Recommends		0.050*** (0.005)	0.040*** (0.006)
Entire Prop. * High LTR			0.011*** (0.003)
Address FE	YES	YES	YES
Observations	232,899	205,085	112,783

The outcome in the above regression is whether the guest's star rating is greater than 3. The sample used is the set of trips to addresses that had multiple listing types, of which one had more than 1 bedroom, which took place between 2012 and 2014. "High LTR" occurs when the guest's likelihood to recommend is greater than 8 (out of 10). *p<0.10, ** p<0.05, *** p<0.01

7 How Large is the Bias?

8 Measuring the Size of Bias

Our analysis has shown that submitted reviews on Airbnb exhibit bias from sorting, strategic reciprocity, and socially induced reciprocity. In this section, we describe a methodology for using experimental estimates to measure bias and quantify the relative importance of the mechanisms documented in this paper.

Our first measure of bias, B_{avg} , is the difference between average experience and the reported experience. This metric represents how off the reputation system is on average in representing the experience of users. Our second measure of bias, B_{neg} , is the share of those with negative experiences who reported negatively. This rate quantifies how many bad guests or hosts are “caught”. To the extent that a bad agent imposes a negative externality on other agents (Nosko and Tadelis (2014)), the platform may especially care about catching these bad agents in the review system. Furthermore, the welfare losses from imperfect reputation systems in section 3 are a direct function of B_{neg} .

8.1 Empirical Analogues of Bias Measures

Suppose that each trip results in a positive experience with probability, g , and a negative experience (denoted n) with probability, $1 - g$. Then an unbiased review system would have a share, g , of positive ratings. Furthermore, suppose that there are only two types of reviews, positive (s_g) and negative. Then the share of submitted ratings that are positive is:

$$\bar{s} = \frac{gPr(r|g)Pr(s_g|g,r) + (1-g)Pr(r|n)Pr(s_g|n,r)}{Pr(r)} \quad (7)$$

where r is an indicator for whether a review was submitted. The difference between the average actual experience and the average submitted review is:

$$B_{avg} = (1-g)\frac{Pr(r|n)Pr(s_g|n,r)}{Pr(r)} - g\left(1 - \frac{Pr(r|g)Pr(s_g|g,r)}{Pr(r)}\right) \quad (8)$$

Where the first term is the share of reviewers with bad experiences who report positively and the second term is the share of all guests with positive experiences who report negatively. Note, these two forms of bias push the average in opposite directions. So looking at average ratings understates the amount of misreporting.

Our second measure of bias is the share of negative experiences not-reported by reviewers:

$$B_{neg} = 1 - \frac{N_{n|n}}{N_{all}(1-g)} \quad (9)$$

where $N_{n|n}$ is the number of negative reports given the reviewer has a negative experience and N_{all} is the total number of trips.

In order to operationalize these metrics, we assume that guests honestly recommend when they leave a review (because the recommendation is anonymous). To calibrate the empirical analogue to g , we need to make assumptions about the degree of selection into

reviewing. Because the recommendation rate for guests in the incentivized review experiment was lower than in the control, $Pr(r|g) \neq Pr(r|b)$. Therefore, we cannot simply use the rates of recommendations in the data to back out g . Instead, we calibrate g by using the recommendation rates from the incentivized review experiment, which eliminate some of the effect of selection into reviewing. However, because the coupon experiment was only conducted for listings with 0 reviews, we must extrapolate to the sample of all reviews. To do so, we assume that the relative bias due to sorting for listings with 0 reviews is the same as the bias due to sorting for the overall sample. We then reweigh the baseline rate of recommendation for listings with 0 reviews by the relative rates of recommendations in the overall sample.

$$\hat{g} = s_{0,ir,sr} \frac{s_{all,sr}}{s_{0,c,sr}} \quad (10)$$

where $s_{0,ir,sr}$ is the share of recommendations in the incentivized review (ir) and simultaneous reveal (sr) treatments, $s_{0,c,sr}$ is the share of recommendations in the ir control and sr treatment, and $s_{all,sr}$ is the share of positive reviews in the entire sr treatment.

For \hat{g} to be an unbiased estimate of the rate of good experiences, we need to make two assumptions about reviews. First, the rate of positive experiences for those that do not review in the coupon experiment must be equal to the rate of positive experiences in the overall sample. We view this assumption as conservative, given that those not induced to review by the Airbnb coupon are likely to have even worse experiences on average than those that did review. Second, the relative rate of bias due to sorting must be the same across listings with different amounts of reviews. In the absence of experimental variation, we cannot confirm or reject this proposition. Lastly, we need to calibrate the review probabilities and mis-reporting rates conditional on leaving a review. We describe how to do so in the next section.

8.2 The Size of Bias

We measure bias for guest reviews of listings in five scenarios, each with progressively less bias. Scenario 1 represents the baseline scenario in the control group in the simultaneous reveal experiment. In this case all three biases (sorting, strategic, and social) operate. Scenario 2 corresponds to the treatment group of the simultaneous reveal experiment (note, there are effects on both ratings and review rates). In both scenarios, we calculate measures of bias by making simple transformations of the moments in the data. $Pr(\widehat{s_g|n}, r)$ is equal to the empirical rate of positive text without a recommendation. $\hat{g} = 3.0\%$ is our best estimate of the true rate of negative experiences in the data and $Pr(\widehat{r|n}) = \frac{Pr(\widehat{n|r}) * P(\widehat{r})}{(1-\hat{g})}$. Scenario 3 represent the bias if there was no socially induced reciprocity in the reviewing process. To calculate the review rates in this scenario, we let $Pr(\widehat{s_g|n}, r)$ equal to the adjusted rate of positive text for stays with multi-listing hosts in entire properties. Scenario 4 removes sorting bias from reviews. This corresponds to the rate of non-recommendation if the share of reviewers with negative experiences was equal to the share of guests with negative experiences. The no-sorting calculation keeps the overall review rate equal to the review rate in the simultaneous reveal treatment. Lastly, scenario 5 computes the two measures of bias if everyone submits reviews.

Table 9 displays both measures of bias in each of the five scenarios.¹⁷ We first turn to the case when all biases are present (row 1). In this scenario, positive reviews occur .78% more of the time than positive experiences. Furthermore, 55% of negative experiences are not reported in text. Removing strategic considerations barely changes these rates of bias. Therefore, we conclude that strategic motivations have little effect on the rate at which negative experiences are reported in text by guests.

Table 9: Size of Bias
(Guest does not recommend listing but omits negative text.)

Counterfactual:	Measure of Bias:	
	B_{avg} Average	B_{neg} % Negative Missing
Baseline	0.78	55.60
Simultaneous Reveal	0.72	52.49
Simultaneous Reveal + No Social Reciprocity	0.49	47.64
Simultaneous Reveal + No Social Reciprocity + No Sorting	0.13	34.15
Above + Everyone Reviews	0.13	4.55

The above table displays three measures of bias under five scenarios. The mis-reporting used to calculate bias occurs when the guest does not recommend the listing but omits any negative review text. B_{avg} is the difference between the average rate of negative sentiment for reviews (where the guest does not recommend), and the overall rate of trips where the guest has a negative experience. B_{neg} is share of all stays where a negative experience was not reported.

Row 3 shows the bias in the case where social reciprocity is removed as a motivation for reviews. The overall bias and percent missing fall by .23 percentage points and 5.2 percentage points respectively. This shows that socially induced reciprocity is a bigger source of bias is more important than strategic bias in the setting of Airbnb. However, both strategic and social reciprocity account for a relatively small portion of the bias in the system.

In row 4, we remove sorting bias. There is a fall of .36 percentage points in average bias and 13 percentage points in the share of negative experiences missing. Sorting is a more important source of bias than strategic and socially induced reciprocity in both cases. Lastly, in Row 5, we report what our measures of bias would be if every guest submitted a review conditional on removing the aforementioned biases. In this case, B_{avg} does not change in this scenario because the rate of misreporting does not change. However, B_{neg} falls by an additional 30 percentage points due to the fact that even without sorting into reviewing, some non-reviewers would have negative experiences which would not be reported. Lastly, there is a residual 4.55% of negative experiences that would still go unreported. This misreporting can correspond to two scenarios: measurement error or residual socially induced reciprocity that occurs even when guests stay at the properties of multi-listing hosts.

The overall level of bias on this site is small. There is a less than 1% difference between the average rate of negative experiences and the average rate of negative reviews. Furthermore, while both strategic and social reciprocity exist, they result in relatively small increases in bias. On the other hand, sorting into reviewing represents the biggest source of bias on the platform. However, the low rate of average bias on Airbnb is driven by the fact that most guests have good experiences. When guests do have negative experiences, they do not report them 52.5% of the time. The main factor in non-reporting is that approximately 30% of guests do not submit a review.

¹⁷See [Table AVII](#) for a measure of bias using five star ratings conditional on a non-recommendation.

9 Discussion

We have shown that the high rate of five star ratings on the site is caused by the fact that most users genuinely have a positive experience. However, even though most users have a positive experience, bias in review still occurs. Our analysis documents review bias due to sorting, strategic reciprocity, and socially induced reciprocity. These biases leads to cases when negative experiences are not reported in review text on the website. If the three biases were eliminated, then an additional 20% of negative experiences would be documented in reviews on the website.

Based on our results, the most important bias to tackle in review systems is sorting into reviewing. There are several potential interventions that might reduce sorting bias in review systems. First, marketplaces can change the way in which reviews are prompted and displayed in order to increase review rates. For example, the simultaneous reveal experiment described in this paper increased review rates and consequently reduced the rate of sorting into reviewing. Other potential interventions include making reviews mandatory (as on Uber), strategically offering coupons for reviews, or making the review easier to submit. Second, online marketplaces can display ratings that adjust for bias in the review system. For example, the effective positive percentage could be shown on a listing page in addition to the standard ratings. Alternatively, listing pages could be augmented with data on other signals of customer experience, such as customer support calls. Lastly, as in [Nosko and Tadelis \(2014\)](#), the platform can use its private information regarding the likely quality of a listing to design a search ranking algorithm.

Our theoretical model suggests that a key variable determining the costs of imperfect reputation systems is the rate of positive experiences on a platform. This begs the question of why the rate of positive experiences is so high on the Airbnb platform. One potential answer is that most bad actors or listings are caught by the Airbnb’s trust and safety efforts. These efforts include verifying the identities of guests and hosts, tracking and preemptively eliminating scams, encouraging detailed profiles, and subsidizing high resolution photos. Alternatively, the law of large numbers may ensure that any low quality listings are eventually negatively reviewed and consequently never booked again. Lastly, because Airbnb is typically a lower priced alternative to hotels, guests may have lower expectations regarding the quality of their experience. Therefore, even when an Airbnb stay is not flawless, guests may nonetheless be satisfied by the value they’ve received. In future work, we plan to study the importance of these mechanisms for equilibrium outcomes such as transaction volume, prices, and welfare.

References

- Andreoni, James, and Justin M. Rao.** 2011. “The power of asking: How communication affects selfishness, empathy, and altruism.” *Journal of Public Economics*, 95(7-8): 513–520.
- Archak, Nikolay, Anindya Ghose, and Panagiotis G. Ipeirotis.** 2011. “Deriving the Pricing Power of Product Features by Mining Consumer Reviews.” *Management Science*, 57(8): 1485–1509.
- Avery, Christopher, Paul Resnick, and Richard Zeckhauser.** 1999. “The Market for Evaluations.” *American Economic Review*, 89(3): 564–584.
- Bohnet, Iris, and Bruno S Frey.** 1999. “The sound of silence in prisoner’s dilemma and dictator games.” *Journal of Economic Behavior & Organization*, 38(1): 43–57.

- Bolton, Gary, Ben Greiner, and Axel Ockenfels.** 2012. "Engineering Trust: Reciprocity in the Production of Reputation Information." *Management Science*, 59(2): 265–285.
- Cabral, Luís, and Ali Hortaçsu.** 2010. "The Dynamics of Seller Reputation: Evidence from Ebay*." *The Journal of Industrial Economics*, 58(1): 54–78.
- Cabral, Luis M. B., and Lingfang (Ivy) Li.** 2014. "A Dollar for Your Thoughts: Feedback-Conditional Rebates on Ebay." Social Science Research Network SSRN Scholarly Paper ID 2133812, Rochester, NY.
- Cullen, Zoe, and Chiara Farronato.** 2015. "Outsourcing Tasks Online: Matching Supply and Demand on Peer-to-Peer Internet Platforms."
- Dai, Weijia, Ginger Jin, Jungmin Lee, and Michael Luca.** 2012. "Optimal Aggregation of Consumer Ratings: An Application to Yelp.com."
- Dellarocas, Chrysanthos, and Charles A. Wood.** 2007. "The Sound of Silence in Online Feedback: Estimating Trading Risks in the Presence of Reporting Bias." *Management Science*, 54(3): 460–476.
- DellaVigna, Stefano, John A. List, and Ulrike Malmendier.** 2012. "Testing for Altruism and Social Pressure in Charitable Giving." *The Quarterly Journal of Economics*, 127(1): 1–56.
- Fradkin, Andrey.** 2014. "Search Frictions and the Design of Online Marketplaces."
- Hoffman, Elizabeth, Kevin McCabe, and Vernon L. Smith.** 1996. "Social Distance and Other-Regarding Behavior in Dictator Games." *American Economic Review*, 86(3): 653–60.
- Hoffman, Elizabeth, Kevin McCabe, Keith Shachat, and Vernon Smith.** 1994. "Preferences, Property Rights, and Anonymity in Bargaining Games." *Games and Economic Behavior*, 7(3): 346–380.
- Horton, John J.** 2014. "Reputation Inflation in Online Markets."
- Lazear, Edward P., Ulrike Malmendier, and Roberto A. Weber.** 2012. "Sorting in Experiments with Application to Social Preferences." *American Economic Journal: Applied Economics*, 4(1): 136–163.
- Li, Lingfang (Ivy), and Erte Xiao.** 2014. "Money Talks: Rebate Mechanisms in Reputation System Design." *Management Science*, 60(8): 2054–2072.
- Luca, Michael.** 2013. "Reviews, Reputation, and Revenue: The Case of Yelp.com." *HBS Working Knowledge*.
- Malmendier, Ulrike, and Klaus Schmidt.** 2012. "You Owe Me." National Bureau of Economic Research Working Paper 18543.
- Malmendier, Ulrike, Vera L. te Velde, and Roberto A. Weber.** 2014. "Rethinking Reciprocity." *Annual Review of Economics*, 6(1): 849–874.
- Mayzlin, Dina, Yaniv Dover, and Judith Chevalier.** 2014. "Promotional Reviews: An Empirical Investigation of Online Review Manipulation." *American Economic Review*, 104(8): 2421–2455.
- Miller, Nolan, Paul Resnick, and Richard Zeckhauser.** 2005. "Eliciting Informative Feedback: The Peer-Prediction Method." *Management Science*, 51(9): 1359–1373.
- Moe, Wendy W., and David A. Schweidel.** 2011. "Online Product Opinions: Incidence, Evaluation, and Evolution." *Marketing Science*, 31(3): 372–386.
- Nagle, Frank, and Christoph Riedl.** 2014. "Online Word of Mouth and Product Quality Disagreement." Social Science Research Network SSRN Scholarly Paper ID 2259055, Rochester, NY.
- Nosko, Chris, and Steven Tadelis.** 2014. "The Limits of Reputation in Platform Markets: An Empirical Analysis and Field Experiment."
- Pallais, Amanda.** 2014. "Inefficient Hiring in Entry-Level Labor Markets." *American Economic Review*, 104(11): 3565–99.
- Saeedi, Maryam, Zequian Shen, and Neel Sundaesan.** 2015. "The Value of Feedback: An Analysis of Reputation System."
- Sally, David.** 1995. "Conversation and Cooperation in Social Dilemmas A Meta-Analysis of Experiments from 1958 to 1992." *Rationality and Society*, 7(1): 58–92.
- Zervas, Georgios, Davide Proserpio, and John Byers.** 2015. "A First Look at Online Reputation on Airbnb, Where Every Stay is Above Average." Social Science Research Network SSRN Scholarly Paper ID 2554500, Rochester, NY.

A Classifying Text Sentiment

In this section we describe the procedure used to classify review text. In order to train a classifier, we need “ground truth” labeled examples of both positive and negative reviews. We select a sample of reviews that are highly likely to be either positive or negative based on the ratings that guests submitted. Reviews by guests that we use as positive examples for guests and hosts are ones that have five star ratings. Reviews by guests that are examples of negative reviews are ones with a 1 or 2 star rating. Reviews by hosts that are examples of negative reviews are ones which have either a non-recommendation or a sub-rating lower than 4 stars. Foreign language reviews were excluded from the sample.

We use reviews between January 2013 and March 2014. Because positive reviews are much more common than negative reviews, the classification problem would be unbalanced if we used the entire sample. Therefore, we randomly select 100,000 examples for both positive and negative reviews. Once we obtain these samples, we remove special characters in the text such as punctuation and we remove common “stop words” such as “a” and “that”.¹⁸ Each review is transformed into a vector for which each column represents the presence of a word or phrase (up to 3 words), where only words that occur at least 300 times are included. We tested various thresholds and regularization to determine that this configuration.

We evaluate model accuracy in several ways. First, we look at the confusion matrix describing model predictions on a 20% hold out sample. For guest reviews of listings, 19% of reviews with low ratings were classified as positive and 9% of reviews with high ratings were classified as negative. The relatively high rate of false positives reflects not only predictive error but the fact that some guests misreport their true negative experiences. We also evaluate model accuracy by doing a 10-fold cross-validation. The mean out of sample accuracy for our preferred model is 87%. The top five positive features are “amazing apartment”, “fantastic apartment”, “of shower”, “excellent stay”, and “exceeded”. The top five negative features were “rude”, “ruined”, “not clean”, “not very clean”, and “christmas”.

B Predictors of Review Rates

Table AI displays the results of a linear probability regression that predicts whether a guest reviews as a function of guest, listing, and trip characteristics. Column 2 adds market city of listing fixed effects in addition to the other variables. If worse experiences result in lower review rates, then worse listings should be less likely to receive a review. The regression shows that listings with lower ratings and lower historical review rates per trip have a lower chance of being reviewed. For example, a listing with an average review rating of four stars is .68 percentage points less likely to be reviewed than a listing with an average rating of five stars. Furthermore, trips where the guest calls customer service are associated with an 11% lower review rate.

Guest characteristics also influence the probability that a review is submitted. New guests and guests who found Airbnb through online marketing are less likely to leave reviews after a trip. This might be due to one of several explanations. First, experienced users who found Airbnb through their friends may be more committed to the Airbnb ecosystem and might feel more of an obligation to review. On the other hand, new users and users acquired through online marketing might have less of an expectation to use Airbnb again. Furthermore, these users might have worse experiences on average, either because they picked a bad listing due to inexperience or because they had flawed expectations about using Airbnb.

C Experimental Validity

This section documents that both experimental designs in this paper are valid. Table AII displays the balance of observable characteristics in the experiments. Turning first to the incentivized review experiment, the rate of assignment to the treatment in the data is not statistically different from 50%. Furthermore, there is no statistically significant difference in guest characteristics (experience, origin, tenure) and host characteristics (experience, origin, room type). Therefore, the experimental design is valid.

Similarly, there is no statistically significant difference in characteristics between the treatment and control guest in the for the simultaneous reveal experiment,. However, there is .3% difference between the number of observations in the treatment and control groups. This difference has a p-value of .073, making it barely significant according to commonly used decision rules. We do not know why this result occurs. We do not view this difference as a problem because we find balance on all observables and the overall difference in observations is tiny.

D Robustness Checks for Incentivized Review Experiment

In this section we test for alternative explanations for the effects of the incentivized review experiment. Column (1) of Table AIII displays the baseline treatment effect of the experiment without any control. Column (2) adds in control for guest origin, experience, and trip characteristics. The treatment effects in columns (1) and (2) are approximately equal (-7.5 percentage points), therefore the treatment is not operating by inducing different types of guests to review.

Column (3) shows estimates for a sample of experienced guests and adds controls for the historical judiciousness of a guest when leaving reviews. The guest judiciousness variable measures the extent to which the guest has previously submitted

¹⁸These words are commonly removed in natural language applications because they are thought to contain minimal information.

lower ratings. It is equal to the negative guest-specific fixed effect in a regression of ratings on guest and listing fixed effects.¹⁹ As expected, the coefficient on the guest judiciousness term is negative, with pickier guests leaving lower ratings. However, adding this control and limiting the sample to experienced guests does not diminish the effect of the experiment on ratings. Furthermore, the interaction between the treatment and guest judiciousness is not significant. Therefore, the rating behavior of these guests, conditional on submitting a review, does not change due to the presence of a coupon. In column (4), we test whether more negative reviews are driven by listing composition. Adding controls for listing type, location, price, and number of non-reviewed stays increases the treatment effect to 6.9 percentage points. We conclude that the coupon works mainly by inducing those with worse experiences to submit reviews.

E Additional Results on Strategic Reciprocity

In this appendix we discuss results regarding strategic reciprocity for hosts who review second and guests who review first. Table AIV displays estimates for two outcomes: whether the host does not recommend and whether the host uses negative sentiment. For all specifications, the coefficient on the treatment is small and insignificant. Therefore, there is no evidence of induced reciprocity by positive guest reviews. However, there is evidence of retaliation in all specifications. Specifications (1) and (2) show that a low rating (< 4 stars) by a guest in the control is associated with a 27 percentage points lower recommendation rate and a 32 percentage points lower negative sentiment rate (defined across all host reviews regardless of the host's recommendation). The interaction with the treatment reduces the size of this effect almost completely. In specifications (3) and (4), we look at three types of initial guest feedback: recommendations, ratings, and negative sentiment conditional on not recommending the host. The predominant effect on host behavior across these three variables is the guest text. Guests' negative text increases hosts' use of negative text by 30 percentage points, while the coefficients corresponding to guests' ratings are relatively lower across specifications. This larger response to text is expected because text is always seen by the host whereas the rating is averaged across all prior guests and rounded. Therefore, hosts may not be able to observe and retaliate against a low rating that is submitted by a guest.

Table AV displays the results for fear of retaliation when guests review first. Column (1) shows that there is no difference in whether guests recommend in the treatment and control. Columns (2) and (3) display the effects of the treatment on the likelihood that guests leave a low rating and negative sentiment in their reviews of hosts. There is an overall increase in lower rated reviews by .4 percentage points and an increase in negative sentiment of 1.1 percentage points. Furthermore, column (4) shows that the effect of the treatment does not vary by the quality of the trip, as measured by recommendation rates and ratings. We interpret this small effect as follows. Although guests may fear retaliation, they may have other reasons to omit negative feedback. For example, guests may feel awkward about leaving negative review text or they may not want to hurt the reputation of the host.

One piece of evidence supporting this theory comes from the effect of the treatment on private feedback. Guests have the ability to leave suggestions for a host to improve the listings. Private feedback cannot hurt the host, but it may still trigger retaliation. Table AVI displays the effect of the treatment on whether a guest leaves a suggestion. Column (1) shows that the overall effect of the treatment is 6.3 percentage points, suggesting that guests are indeed motivated by fear of retaliation. Columns (2) and (3) test whether this effect is driven by particular types of trips by interacting the treatment indicator with indicators for guests' recommendations and ratings. The effect of the treatment is especially large for guests that recommend the host. Therefore, the treatment allows guests who have good, but not great, experiences to offer suggestions to the host without a fear of retaliation. In the next section we further explore behavioral reasons for reviewing behavior.

¹⁹The estimation sample for the fixed effects regressions is the year before the start of the experiment.

F Additional Tables

Table AI: Determinants of Guest Reviews

	Reviewed	
Five Star Rate	0.105*** (0.008)	0.106*** (0.008)
Past Booker	0.059*** (0.004)	0.059*** (0.004)
No Reviews	0.026** (0.013)	0.025* (0.013)
No Trips	0.096*** (0.012)	0.098*** (0.012)
Num. Trips	-0.0004*** (0.0001)	-0.0004*** (0.0001)
Customer Service	-0.174*** (0.020)	-0.167*** (0.020)
Entire Property	0.004 (0.005)	0.005 (0.005)
Multi-Listing Host	-0.095*** (0.007)	-0.084*** (0.007)
Log Price per Night	-0.011*** (0.003)	-0.012*** (0.003)
Trip Characteristics	Yes	Yes
Market FE:	No	Yes
Observations	60,579	60,579

Note: *p<0.1; **p<0.05; ***p<0.01

These regressions predict whether a guest submits a review conditional on the observed characteristics of the listing and trip. Only observations in the control group of the simultaneous reveal experiment are used for this estimation.

Table AII: Experimental Validity Check

Variable	Experiment	Difference	Mean Treatment	Mean Control	P-Value	Stars
Experienced Guest	Simultaneous Reveal	-0.001	0.557	0.558	0.702	
US Guest	Simultaneous Reveal	-0.001	0.282	0.283	0.761	
Prev. Host Bookings	Simultaneous Reveal	-0.162	14.875	15.037	0.272	
US Host	Simultaneous Reveal	0.001	0.263	0.262	0.801	
Multi-Listing Host	Simultaneous Reveal	0.001	0.082	0.081	0.369	
Entire Property	Simultaneous Reveal	-0.001	0.671	0.671	0.824	
Reviewed Listing	Simultaneous Reveal	-0.003	0.764	0.767	0.167	
Observations	Simultaneous Reveal	0.001			0.431	
Experienced Guest	Incentivized Review	-0.010	0.498	0.508	0.066	*
US Guest	Incentivized Review	0.001	0.228	0.227	0.859	
Prev. Host Bookings	Incentivized Review	-0.008	0.135	0.143	0.134	
US Host	Incentivized Review	0.0002	0.199	0.199	0.973	
Multi-Listing Host	Incentivized Review	0.002	0.169	0.167	0.678	
Entire Property	Incentivized Review	0.002	0.683	0.681	0.645	
Host Reviews Within 7 Days	Incentivized Review	-0.009	0.736	0.745	0.147	
Observations	Incentivized Review	0.005			0.102	

This table displays the averages of variables in the treatment and control groups, as well as the statistical significance of the difference in averages between treatment and control. Note, the sample averages for the two experiments differ because only guests to non-reviewed listings who had not reviewed within 9 days were eligible for the incentivized review experiment. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Table AIII: Effect of Coupon Treatment on Five Star Ratings

	(1)	(2)	(3)	(4)	(5)
Treatment	-0.082*** (0.010)	-0.081*** (0.009)	-0.104** (0.048)	-0.077*** (0.009)	-0.088*** (0.017)
Guest Lenient			0.168*** (0.057)		
Treatment * Guest Lenient			0.044 (0.074)		
Host Rev. First					0.073*** (0.017)
Treatment * Host Rev. First					0.032 (0.021)
Guest Characteristics	No	Yes	Yes	Yes	Yes
Listing Characteristics	No	No	No	Yes	Yes
Observations	10,623	10,623	615	10,623	10,623

The table displays results of a regression predicting whether a guest submitted a five star rating in their review. “Treatment” refers to an email that offers the guest a coupon to leave a review. “Guest Judiciousness” is a guest specific fixed effect that measure a guest’s propensity to leave negative reviews. Judiciousness is estimated on the set of all reviews in the year proceeding the experiment. Guest controls include whether the guest is a host, region of origin, age, gender, nights of trip, number of guests, and checkout date. Listing controls include whether the host is multi-listing host, price, room type of the listing, and listing region. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Table AIV: Retaliation and Induced Reciprocity - Host

	Does Not Recommend (1)	Negative Sentiment (2)	Does Not Recommend (3)	Negative Sentiment (4)
Treatment	-0.002 (0.001)	0.006 (0.008)	-0.001 (0.001)	0.007 (0.009)
Guest Low Rating	0.239*** (0.028)	0.315*** (0.048)	0.104*** (0.031)	0.158*** (0.058)
Guest Review Negative Words			0.349*** (0.093)	0.199* (0.119)
Guest Does Not Recommend			0.083 (0.064)	0.162* (0.092)
Treatment * Low Rating	-0.176 (0.031)	-0.257*** (0.058)	-0.048 (0.037)	-0.117 (0.075)
Treatment * Review Negative Words			-0.260*** (0.098)	-0.180 (0.149)
Treatment * Does Not Recommend			-0.123* (0.065)	-0.149 (0.118)
Guest, Trip, and Listing Char. Observations	Yes 14,376	Yes 8,188	Yes 10,682	Yes 7,518

The above regressions are estimated for the sample where the guest reviews first. Controls include whether there was a contact to customer support, the user type (new vs experienced, organic vs acquired by marketing), nights and number of guests of trip, whether the guest was a host, the age of the guest, the gender of the guest, whether the host is a property manager, the average review rating of the host, and the Effective Positive Percentage of the host. "Treatment" refers to the simultaneous reveal experiment. *p<0.10, ** p<0.05, *** p<0.01

Table AV: Fear of Retaliation - Guest

	Reviews First (1)	< 5 Rating (First) (2)	Neg. Sentiment (First) (3)	Neg. Sentiment (First) (4)
Treatment	0.0005 (0.002)	0.002 (0.004)	0.007 (0.004)	0.009* (0.005)
< 5 Rating				0.157*** (0.009)
Not Recommend		0.658*** (0.007)		0.438*** (0.022)
Treatment * < 5 Rating				-0.008 (0.012)
Treatment * Not Recommend				-0.024 (0.031)
Guest, Trip, and Listing Char. Observations	Yes 38,023	Yes 38,021	Yes 31,350	Yes 29,370

The regressions in columns (2) - (4) are estimated only for cases when the guest reviews first. “Treatment” refers to the simultaneous reveal experiment. Controls include whether there was a contact to customer support, the user type (new vs experienced, organic vs acquired by marketing), nights and number of guests of trip, whether the guest was a host, the age of the guest, the gender of the guest, whether the host is a property manager, the average review rating of the host, and the Effective Positive Percentage of the host. *p<0.10, ** p<0.05, *** p<0.01

Table AVI: Determinants of Private Feedback Increase

	Guest Left Private Suggestion for Host		
	(1)	(2)	(3)
Treatment	0.064*** (0.003)	0.046*** (0.004)	0.052*** (0.007)
Customer Support	0.075*** (0.019)	0.082*** (0.019)	0.079*** (0.019)
Guest Recommends		0.047*** (0.003)	0.052*** (0.003)
Five Star Review			-0.074*** (0.005)
Recommends * Treatment		0.022*** (0.004)	0.023*** (0.004)
Five Star * Treatment			-0.012* (0.007)
Guest, Trip, and Listing Char. Observations	Yes 82,623	Yes 82,623	Yes 82,623

“Treatment” refers to the simultaneous reveal experiment. “Customer Support” refers to a guest initiated customer service complaint. Controls include the user type (new vs experienced, organic vs acquired by marketing), nights and number of guests of trip, whether the guest was a host, the age of the guest, the gender of the guest, whether the host is a property manager, and the five star review rate of the host. *p<0.10, ** p<0.05, *** p<0.01

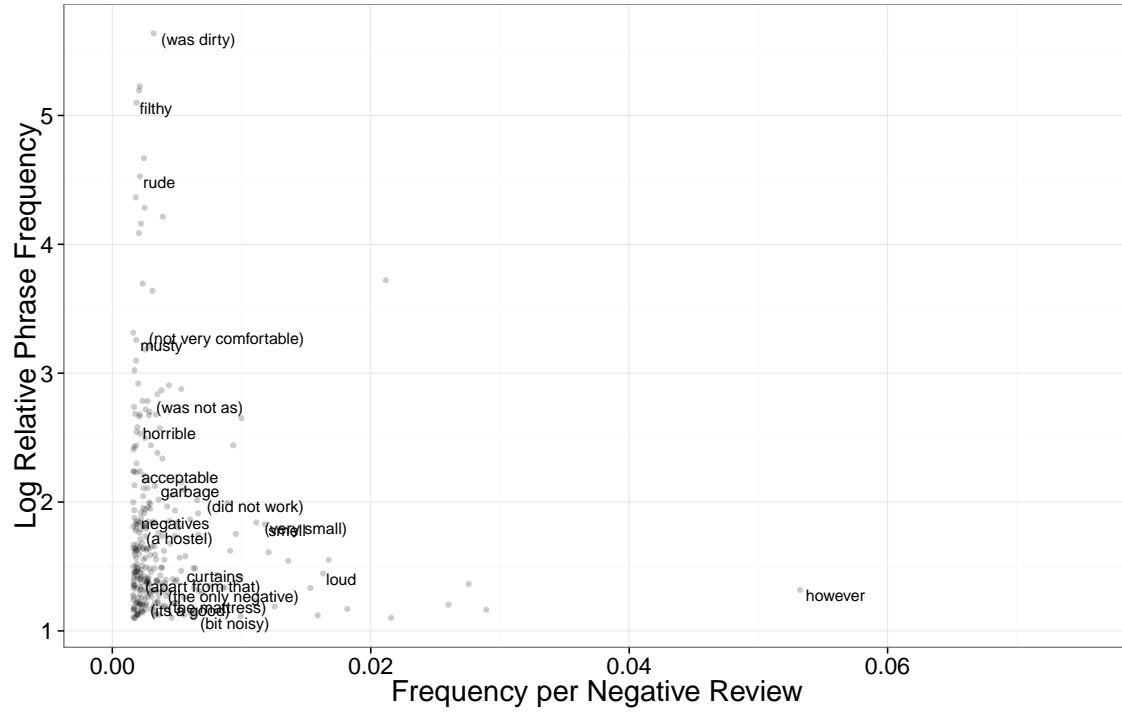
Table AVII: Size of Bias
(Guest does not recommend listing but submits five star rating.)

Counterfactual:	Measure of Bias:		
	B_{avg} Average	B_{mis} % Misreported	B_{neg} % Negative Missing
Baseline	0.58	0.15	51.50
Simultaneous Reveal	0.55	0.18	48.99
Simultaneous Reveal + No Social Reciprocity	0.40	0.03	45.76
Simultaneous Reveal + No Social Reciprocity + No Sorting	0.03	0.03	31.79
Above + Everyone Reviews	0.03	0.03	1.12

The above table displays three measures of bias under five scenarios. The mis-reporting used to calculate bias occurs when the guest does not recommend the listing but omits any negative review text. B_{avg} is the difference between the average rate of negative sentiment for reviews (where the guest does not recommend), and the overall rate of trips where the guest has a negative experience. B_{mis} is the share of all reviews that are mis-reported and B_{neg} is share of all stays where a negative experience was not reported.

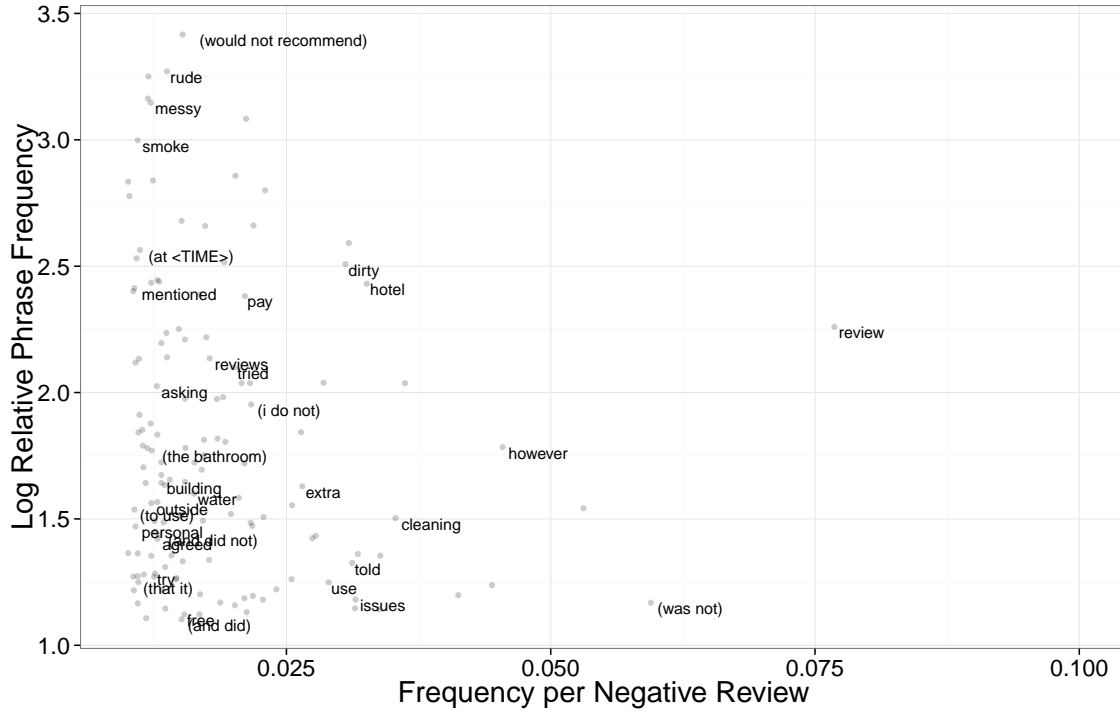
G Additional Figures

Figure A1: Distribution of negative phrases in guest reviews of listings.



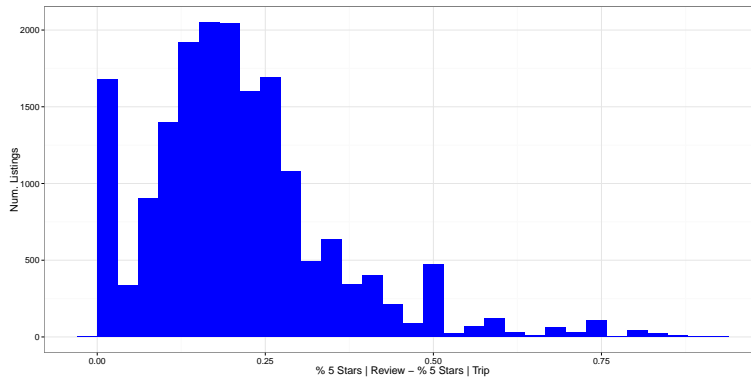
“Relative phrase frequency” refers to the ratio with which the phrase occurs in reviews with a rating of less than five stars.

Figure A2: Distribution of negative phrases in host reviews of guests.



“Relative phrase frequency” refers to the ratio with which the phrase occurs in reviews with a non-recommendation.

Figure A3: Histogram of Difference in Ratings per Listing



The sample used for this figure is composed of highly rated listings (> 4.75 average overall star rating) with at least 3 reviews. This sample is chosen because Airbnb only displays star ratings after 3 reviews are submitted and rounds the star rating the nearest .five stars. Therefore, the listings in this sample seem the same to guests on the overall star rating dimension.