

# Reputation offsets trust judgments based on social biases among Airbnb users

Bruno Abrahao<sup>a,1</sup>, Paolo Parigi<sup>b</sup>, Alok Gupta<sup>c</sup>, and Karen S. Cook<sup>a,1</sup>

<sup>a</sup>Department of Sociology, Stanford University, Stanford, CA 94305; <sup>b</sup>Institute for Research in the Social Sciences, Stanford University, Stanford, CA 94305; and <sup>c</sup>Data Science Division, Airbnb, San Francisco, CA 94103

Contributed by Karen S. Cook, July 20, 2017 (sent for review March 15, 2016; reviewed by David Lazer and Chris Snijders)

To provide social exchange on a global level, sharing-economy companies leverage interpersonal trust between their members on a scale unimaginable even a few years ago. A challenge to this mission is the presence of social biases among a large heterogeneous and independent population of users, a factor that hinders the growth of these services. We investigate whether and to what extent a sharing-economy platform can design artificially engineered features, such as reputation systems, to override people's natural tendency to base judgments of trustworthiness on social biases. We focus on the common tendency to trust others who are similar (i.e., homophily) as a source of bias. We test this argument through an online experiment with 8,906 users of Airbnb, a leading hospitality company in the sharing economy. The experiment is based on an interpersonal investment game, in which we vary the characteristics of recipients to study trust through the interplay between homophily and reputation. Our findings show that reputation systems can significantly increase the trust between dissimilar users and that risk aversion has an inverse relationship with trust given high reputation. We also present evidence that our experimental findings are confirmed by analyses of 1 million actual hospitality interactions among users of Airbnb.

online trust | reputation systems | sharing economy | social biases | risk

A new wave of companies, emerging under the banner of the sharing economy (1), is profoundly altering the way we interact and exchange with one another. These Internet-based services are driving a major change in our cultural and technological landscapes and have achieved astounding success, enabling users to share their own personal resources, such as their vehicles, real estate properties, time, or skills. A growing number of individuals trust the sharing economy with a variety of services to satisfy their needs, to generate income, or, more simply, to meet new people. Examples of sharing-economy transactions include hiring a “tasker” from Task Rabbit to run errands, sharing a “couch” with a perfect stranger through CouchSurfing, hiring a “driver” on Uber, or staying in someone's home while traveling using Airbnb.

Users in the sharing economy seek to connect with others engaged in activities on the same platform. Compared with exchanges via traditional e-commerce companies, where transactions are relatively anonymous, the sharing economy exposes us to the more personal character of such interactions. This inevitably prompts attention to the users' sociodemographic characteristics as factors that drive selection.

As a consequence, social biases figure as major hurdles to the growth of sharing-economy services, as they influence users' perceptions of trust and risk. To enable trust between strangers so that everyone can exchange with anyone, beyond cultural and social boundaries, these companies face daunting obstacles in their attempts to minimize these biases.

In this study, we investigate whether and to what extent a sharing-economy platform can design technological features to counteract natural behavioral tendencies that may lead to social biases. This question is of central importance in the social sciences more broadly, but also in the engineering of platforms that aim to enable trust.

Social biases are a result of a number of mechanisms that are difficult to measure. In this work, we make social biases amenable to investigation by focusing on a form of social bias that naturally maps into a quantifiable interpretation and that we expect to be at work in these environments. At the same time, this source of bias is well understood in the social sciences so that we can rely on previous literature, instead of opening up a new dimension of complexity. To this end, we focus on homophily (2–6), the higher likelihood that people trust others who are similar to themselves.

McPherson (4) proposed a theory of how homophily structures modern societies using a construct of social space defined in Blau's theory of preferences (6). Each individual occupies a position in the social space whose coordinates are a function of his or her sociodemographic characteristics. The more features two individuals share in common, the more likely they are to form relationships based on mutual trust.

To operationalize homophily in a structured way, we use Blau's construct of social space to induce and measure the effect of homophily in an experimental setting whose volunteers are active members of the sharing economy. (At the time of writing, the online experiment is accepting participants for demonstration purposes at [stanfordexchange.org](http://stanfordexchange.org).)

Building on this baseline, the heart of our experiment is the measurement of the extent to which another source of information that can be artificially engineered could potentially alter the

## Significance

**We investigate the extent to which artificial features engineered by sharing-economy platforms, such as reputation systems, can be used to override people's tendency to base judgments of trustworthiness on social biases, such as to trust others who are similar (i.e., homophily). To this end, we engaged 8,906 users of Airbnb as volunteers in an online experiment. We demonstrate that homophily based on several demographic characteristics is a relatively weak driver of trust. In fact, having high reputation is enough to counteract homophily. Using Airbnb data, we present evidence that the effects we found experimentally are at work in the actual platform. Lastly, we found an inverse relationship between risk aversion and trust in those with positive reputations.**

Author contributions: B.A., P.P., and K.S.C. designed research; B.A., P.P., and A.G. performed research; B.A. and P.P. contributed new reagents/analytic tools; B.A., P.P., A.G., and K.S.C. analyzed data; and B.A., P.P., and K.S.C. wrote the paper.

Reviewers: D.L., Northeastern University; and C.S., Eindhoven University.

Conflict of interest statement: A.G. is a data scientist at Airbnb who performed the experiments that rely on the company's private data. Paolo Parigi began working at Uber after the research design, experiment execution, data analysis, and writing of the study were completed.

Freely available online through the PNAS open access option.

Data deposition: The data necessary to replicate our experimental results are available through the Stanford Exchange Project at [stanfordexchange.org/project](http://stanfordexchange.org/project).

<sup>1</sup>To whom correspondence may be addressed. Email: [abrahao@cs.stanford.edu](mailto:abrahao@cs.stanford.edu) or [kcook@stanford.edu](mailto:kcook@stanford.edu).

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1604234114/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1604234114/-DCSupplemental).



is a postinteraction subjective evaluation of an alter. It consists of the assignment of zero to five stars, where the number of stars is proportional to the degree of positiveness. The ratings a member receives are averaged over all of their raters, rounded to the half unit, and presented in the member's profile on the platform. Similarly, an interaction grants the two parties the opportunity to mutually provide free-form written reviews. Due to the difficulty of manipulating textual contents of reviews experimentally, we restricted our attention to the number of reviews a user received.

We manipulated these two dimensions in a structured way to study their effects on trust. Among the five profiles participants saw on the screen, four had reputation features with similar values, chosen independently at random for each participant's session, which we refer to as the baseline reputation. These were the profiles at social distances  $d = 0, 1, 2$  and one of the profiles at  $d = 4$ . The other generated profile at distance  $d = 4$  had one of the reputation features randomly selected to be switched to either a better or a worse value than baseline (see [Game Design Details](#) for how we manipulated the numerical values of reputation). For convenience, we refer to the profile that has a different reputation feature than the baseline as being at distance  $d = 5$ .

We randomly assigned users to two possible worlds. In world 1, the profile at  $d = 5$  not only had the largest distance from the participant, but also a weaker reputation than all other profiles (the baseline reputation). In this case, reputation did not compete with the tendency toward homophily. In world 2, the profile at  $d = 5$  had a better reputation than the baseline reputation. This induced a tension between placing trust in the most distant profile with a better reputation or in the other profiles closer to the participant in social space. [Fig. S1](#) shows a partial view of the screen users see in the experiment, and [Fig. S2](#) shows a diagram that exemplifies the structure of a user's session.

We gave participants a single "wallet" with 100 credits, which they could keep or invest in receivers in whatever way they chose. Therefore, participants could gain or lose credits through their investments. Because this was a one-time game, it was easy to show that the Nash equilibrium was not to invest any amount, since the dominant strategy for receivers was not to return any amount. (Nevertheless, we observed such rational behavior only in rare instances.)

It is argued that risk is a component of trust in general, and some definitions of trust include risk (8). Even though previous research has attempted to relate trust and risk, the empirical evi-

dence of the connection between risk attitudes and trust has been weak (17). Moreover, research that has addressed this question has been limited to laboratory experiments or small datasets.

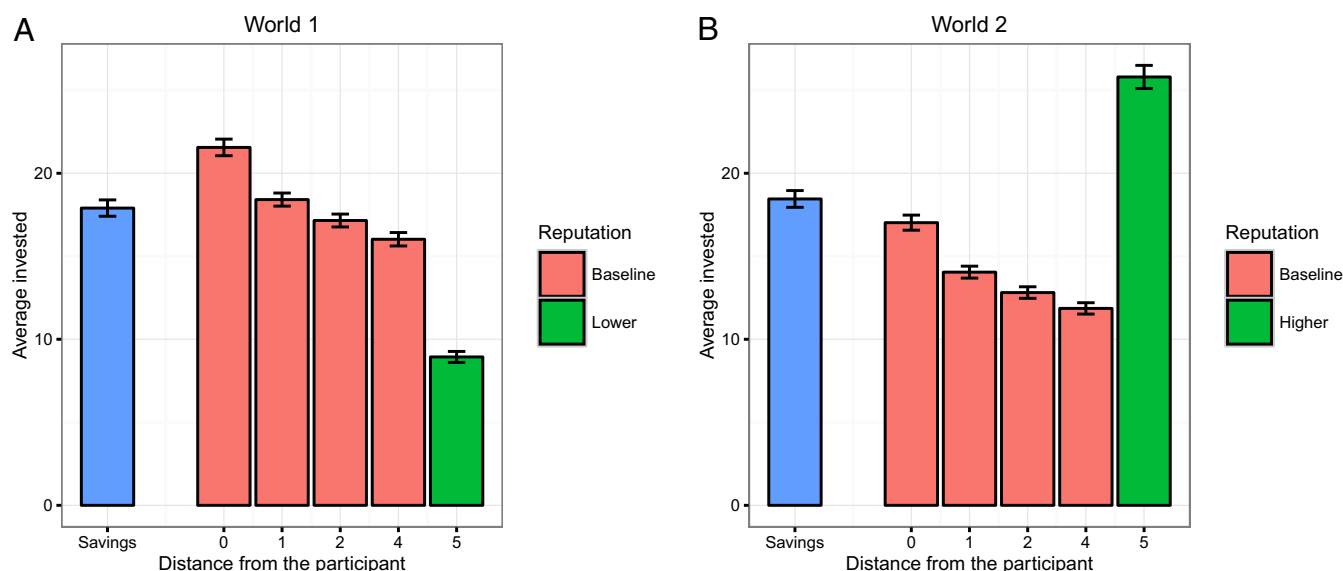
Given the opportunity to study this question using a large population, we introduced a risk-assessment question before the investment game. We worded the question as: "A lottery ticket costs 100 (USD) and people win with 50% chance. How much should the prize be for you to choose to buy a ticket?" Players could enter any numerical value, which corresponded to the minimum reward that would make the participant take the risk of buying a ticket. The prize value 200 (USD) had the expected value of net gain equal to zero (after paying off the ticket) and corresponded to the minimum rational value. Thus, values  $>200$  (USD) measured risk aversion proportional to their magnitude. In [Risk Assessment Question](#), we summarize the distribution of answers ([Table S1](#)) and argue that our measure captures risk behaviors in accordance with previous research ([Table S5](#)) (19, 20).

### Multilevel-Multivariate Analysis

We had five measurements (investments) on each observational unit (participant). As a result, the five investments were correlated, which we accounted for by nesting investments within subjects in a multilevel model. We fitted the model using a multivariate regression with 10 independent variables, one for each investment in the combination  $(d, w)$  of profile distance  $d: \{0, 1, 2, 4, 5\}$  and world  $w: \{1, 2\}$ . The investments a participant made had different sources of mutual correlation. For instance, the sum of the investments had to be at most 100 credits. We accounted for these by computing the model fit with an unconstrained covariance structure that learned from the data the correlations and independent variances across measurements (21).

As a first-order approximation, we fitted the empty model (i.e., without explanatory variables) with 10 intercepts. The five intercepts for each world corresponded to the average distribution of investments among the five profiles across all participants (complete pooling). [Fig. 1](#) shows a plot of the mean estimates, together with the mean number of credits saved, for worlds 1 and 2. [Table S6](#), model 1 shows the numerical estimates from the model fit.

We were mainly interested in the additive effect of the number of different coordinates between two individuals' feature vectors, or their Hamming distance. However, any real-world sociodemographic feature inevitably produces heterogeneous effects on trust (e.g., gender may affect investments more than marital



**Fig. 1.** Empty model estimates of average investment in profile at distance  $d$  and average savings. (A) In world 1, the second profile at distance  $d = 4$  (here identified as  $d = 5$ ) has a worse reputation than baseline. (B) In world 2, the profile at distance  $d = 5$  has a better reputation than the baseline.

whereas positive values increase them. Our main goal was to show that the heterogeneity of the features did not significantly alter the main effects we observed on average investments as a function of  $d$  in the empty multivariate model.

## Results

Fig. 14 shows that homophily dominated investment decisions. That is, the farther away the profile was on the demographic dimensions from the participants, the lower the investment they received, on average. Furthermore, the profile at  $d = 5$  with worse reputation received less investment on average than the equivalent alternative with respect to social distance (i.e., the profile at  $d = 4$ ). Quite strikingly, Fig. 1B shows that reputation builds trust beyond homophily. The average investment in the profile at  $d = 5$ , possessing the best reputation, was significantly higher than the average invested in all of the closest profiles. Note that despite the strong influence of the reputation system in world 2, the magnitude of the investments in the profiles with baseline reputation was still driven by homophily.

The explanatory variables exhibited variance beyond that explained by social distance, which implies that there are differences in investment behavior by demographic group and their interactions. However, as we argue next, the changes in the average investments (model intercepts) that these effects produced in the multivariate model were not strong enough to significantly alter the conclusions regarding homophily and reputation that we previously derived from the empty model.

**Homophily Is at Work.** The covariate “profile distance” was by far the dominant one with respect to variance explained (F value 5668.8,  $P < 0.001$ ). This was followed by the number of reviews with a much smaller F value (26.1,  $P < 0.001$ ).

The dashed lines in Fig. 2 have the values  $\pm 1.37$  and correspond, in the most conservative way, to the smallest difference in average investment between two profiles with baseline reputation, minus two standard errors. That is, a coefficient that

