

NBER WORKING PAPER SERIES

SMILES IN PROFILES:
IMPROVING FAIRNESS AND EFFICIENCY USING ESTIMATES OF
USER PREFERENCES IN ONLINE MARKETPLACES

Susan Athey
Dean Karlan
Emil Palikot
Yuan Yuan

Working Paper 30633
<http://www.nber.org/papers/w30633>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
November 2022, Revised March 2023

We thank Kiva Microfunds for generously sharing data and discussing the research questions. We thank Herman Donner and Kristine Koutout at Stanford Graduate School of Business and Allison Koenecke at Cornell University for comments and suggestions. The Golub Capital Social Impact Lab at Stanford Graduate School of Business provided funding for this research. This research has been subject to review and approval by Research Compliance Office at Stanford University, protocol number IRB-62442 and registered at AEA RCT registry with the number 0010030. The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

At least one co-author has disclosed additional relationships of potential relevance for this research. Further information is available online at <http://www.nber.org/papers/w30633>

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2022 by Susan Athey, Dean Karlan, Emil Palikot, and Yuan Yuan. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Smiles in Profiles: Improving Fairness and Efficiency Using Estimates of User Preferences
in Online Marketplaces

Susan Athey, Dean Karlan, Emil Palikot, and Yuan Yuan

NBER Working Paper No. 30633

November 2022, Revised March 2023

JEL No. D0,D40,J0,J02,O1

ABSTRACT

Online platforms often face the challenge of being both fair (i.e., non-discriminatory) and efficient (i.e., maximizing revenue). Using computer vision algorithms and observational data from a micro-lending marketplace, we find that the choices that online borrowers make when creating online profiles impact both of these objectives. We further support this finding with a web-based randomized survey experiment. In the experiment, we create profile images using Generative Adversarial Networks that differ in a specific feature and estimate the impact of the feature on lender demand. We then evaluate counterfactual platform policies based on the changeable profile features, and identify approaches that can ameliorate the fairness-efficiency tension.

Susan Athey
Graduate School of Business
Stanford University
655 Knight Way
Stanford, CA 94305
and NBER
athey@stanford.edu

Emil Palikot
306 Fair Oaks st
San Francisco, CA 94110
emil.palikot@gmail.com

Yuan Yuan
5000 Forbes Avenue
Pittsburgh, PA 15213
yuany3@andrew.cmu.edu

Dean Karlan
Kellogg Global Hub
Northwestern University
2211 Campus Drive
Evanston, IL 60208
and CEPR
and also NBER
dean.karlan@gmail.com

A randomized controlled trials registry entry is available at:
<https://www.socialscienceregistry.org/trials/10030>

1 Introduction

Personal profile images play an important role in the success of many online platforms (Ert et al., 2016). At the same time, profile images that reveal users' characteristics might enable discrimination, and can lead to severe inequities in outcomes (Edelman and Luca, 2014). Using data from Kiva, an online micro-lending platform, we examine how features of profile images affect the fairness and efficiency goals of a platform. We then use simulation exercises to examine how platforms can intervene to promote fairness and efficiency.

To build intuition we consider a simple example, sellers on an online marketplace differ in two dimensions: characteristics that are fixed when they create their online profiles - *types*, and characteristics they choose at that moment - *style*. Suppose that *type* can be high or low and *style* is smiling or not smiling, and that buyers prefer sellers of high *type* and sellers with smiling profiles. If *types* and *style* choices are uncorrelated, we have two distinct sources of inequity: high *type* sellers outperform low *type* sellers and sellers with smiling profiles outperform those who are not smiling in their profile photo. When *type* and *style* are positively correlated, the two inequities compound, and when they are negatively correlated, they mitigate one another.

We analyze the *type* and *style* features of the online profiles on a non-profit micro-lending platform, Kiva. On Kiva, individual lenders make loans to borrowers, selecting from a curated catalog of borrowing campaigns.¹ In principle, Kiva balances two objectives: fairness and efficiency. One way to think about fairness for Kiva is that it gives different *types* of borrowers equal access to capital. An important component of efficiency is volume, specifically flow of capital from individual lenders (mostly in the United States) to borrowers in low-income countries.

Using observational data on funding outcomes from Kiva, we start by documenting substantial inequities in average daily funds collected by different borrowing campaigns listed on Kiva. Next, we detect features of images using an off-the-shelf machine learning algorithm and select *style* characteristics: features that are changeable when users create their profiles. Comparing the predictive performance of models trained with and without *style* features, we show that jointly *style* features are predictive of funding outcomes. We showcase several specific features that have a large and statistically significant impact on funding outcomes, both unconditionally and after adjusting for other observable features of the borrower. In particular, we show that *smile* is associated with better, and

¹Technically, the loan is made to a microcredit institution and earmarked to the specific borrower.

body-shot with worse funding outcomes. We find that *style* features are not predictive of the probability of repaying the loan.² Finally, we analyze the correlation between *types* and *styles*. A borrower *type* is a collection of characteristics that are fixed at the time the borrower is creating an online profile, such as race, gender, or country of origin. We document correlations between *type* and *style* characteristics and show that the desirable *style* features are generally more prevalent among borrowers with *types* associated with better funding outcomes. For example, high-performing borrower types, such as women, are more likely to have profiles with *smile* and less likely to have *body-shot* profile images. This evidence indicates that the distribution of *style* features exacerbates existing inequities in a way that is unfair to the borrowers in that it is not justified by different repayment probabilities.

Estimates of the impact of *style* features on outcomes from observational data rely on the assumption of unconfoundedness, which is not directly testable. In this setting, it would require that for a particular feature such as *smile*, no aspects of the photographs or profile descriptions correlated with *smile*, other than those that are adjusted for, matter for funding outcomes. Even though we use a state-of-the-art feature detection algorithm, it is plausible that we do not capture all information that lenders discern from profile images. To address this issue, we provide evidence from an experiment with recruited subjects on the Prolific.co platform. Subjects choose between borrowing campaigns featuring fabricated profile images. We use images that are generated with Generative Adversarial Networks and can be thought of as pairs of images that are identical except for a feature that we specify. We analyze two features, *smile* and *body-shot*. The estimates that we obtain about the effect of these features on preferences are consistent with our estimates from the observational data from the Kiva platform.

The evidence of the positive correlation between the *types* associated with high funding outcomes and desirable *style* features indicates that inequities due to intrinsic borrowers' characteristics are exacerbated due to profile *style* choices. However, this also means that fairer outcomes could be achieved by implementing a policy that encourages low-*type* borrowers to change their profile images. Section 6 focuses on the impact of different platform policies on fairness and efficiency.

We consider platform policies that change the conditional distribution of *style* features of borrowers' profiles and that vary probabilities of including borrowers in the lenders' choice set based on borrowers' characteristics. We calibrate a model of lenders' demand for borrowers using estimates from

²This evidence indicates that the inequity due to *style* features is not justified by different repayment rates. However, *style* features might also be informative about the developmental impact of the borrowing campaign. We do not adjust the estimates for the expected impact because we do not have a good measure of it; this is a limitation of this work.

the recruited experiment. We find that policies that alter the distribution of desirable *style* features such that they become less correlated with *types* improve both fairness and efficiency. Specifically, we show that a policy of *style* curation, which encourages borrowers to have profiles with *smile* and to avoid *body-shots*, improves fairness, as measured by the Gini coefficient or the market share of the bottom tercile of borrowing campaigns. It also leads to a higher number of transactions.³ In contrast, a policy that increases the prominence of campaigns with *smile* and without *body-shot*, for example, by ranking them higher on the search page, leads to less fair outcomes, although it boosts efficiency. This is because promoting the selected features increases the prominence of high *type* borrowers. Note that if a platform trains a recommendation system based on funding data, and the recommendation system accounts for image features, it is likely that the recommendation system would indeed increase the prominence of profiles with *style* features that are attractive to users, so this policy captures the expected outcome if a platform implements a recommendation system that incorporates image features.

We showcase a specific dimension of *type*-based inequity: the gender gap in favor of campaigns with *female* profiles.⁴ We show that campaigns with *male* profiles collect 32% fewer funds per day. We corroborate this finding in the recruited experiment where we find that subjects are 31% more likely to choose a *female* profile.⁵ The distribution of selected *style* features exacerbates the gap: 77% of borrowers that our algorithm classified as *females* have profiles with *smile* and 22% of them have *body-shot* profile images. In contrast, 33% of *male* borrowers have profiles with *smile* and 26% with *body-shot*. In the counterfactual simulations, we show that a profile-*style* recommendation policy can substantially narrow the gender gap, while the policy of increasing the prominence of profiles with selected *style* features exacerbates the disparity.

The evidence we present demonstrates how in two-sided markets where users have preferences for certain features in profile images, the correlation between *types* and *style* choices can matter for fairness and efficiency. Thus, platforms faced with balancing the two objectives need to account for this correlation before implementing policies based on profile images.

³We compare the outcomes under counterfactual policies to a *fair* benchmark in which the distribution of outcomes is unaffected by *style* choices; when the outcomes under a counterfactual policy are closer to the benchmark, we argue that the policy improves fairness.

⁴Throughout, we use *male* and *female* to denote the feature detection algorithm's prediction of the gender of the person in the image.

⁵Lenders might prefer campaigns with female profiles for a variety of reasons. For example, there is ample evidence that female entrepreneurs that obtain funding through microfinance generally use the funds effectively (D'Espallier et al., 2011; Aggarwal et al., 2015). Also, lenders might want to compensate for discrimination against women in traditional entrepreneurial finance (Alesina et al., 2013).

The paper is organized as follows: Section 2 presents related literature. Section 3 describes how micro-lending platforms operate and provide institutional details about Kiva. Section 4 presents the observational data and its analysis. Section 5 describes the design of the experiment and its results. Section 6 focuses on counterfactual simulations, and Section 7 concludes.

2 Literature review

There is a rich literature on the role of images in shaping choices online. In the context of Airbnb, Ert et al. (2016) show that personal photos increase the sense of personal contact and improve users' perception of the service. Many papers focus on the impact of *type* features. Edelman et al. (2017) and Ge et al. (2016) provide evidence from field experiments that demographic characteristics revealed in images impact the choices of users on hospitality and ride-sharing platforms. In the context of online lending Theseira (2009), Pope and Sydnor (2011), and Younkin and Kuppuswamy (2018) show that loan applications with pictures of black borrowers are less likely to be funded. Jenq et al. (2015) document that lenders on online peer-to-peer lending platforms favor more attractive and light-skinned borrowers, and Ravina (2019) documents an impact of the physical beauty of the borrower. Park et al. (2019) use an online lab experiment to show that the interaction of borrowers' perceived gender and beauty affects lending decisions. In Kiva's context, Galak et al. (2011) show experimentally that lenders have a preference for borrowers of the same gender.

Other papers focus separately on the impact of *style* on outcomes. Duarte et al. (2012) show that borrowers who appear more trustworthy are more likely to have their loans funded.⁶ Septianto and Paramita (2021), in a recruited experiment, document that profiles with happy images receive more donations. Pham and Septianto (2019) and Jordan et al. (2019) show that smiling increases the attractiveness of profiles in the charitable giving context. In an analogous setting to this paper, Ai et al. (2016) use a field experiment to document that lenders are more likely to join teams from similar locations. We contribute to this literature by using Generative Adversarial Networks to provide causal evidence of the impact of selected profile features on outcomes. We introduce a distinction between features of images that are intrinsic to borrowers (*type*) and characteristics that can be modified (*style*) and show that features in both of these categories impact outcomes. Troncoso and Luo (2022), in a context of a freelancing platform, show that both the features, which we classify as *type* as well as those that we consider as *style* impact whether the applicant is considered as a good fit for a job.

⁶Trustworthiness is rated by human-raters based on profile images. Krumbhuber et al. (2007) argue that trustworthiness is related to the dynamics of facial expressions.

Because we argue that platforms can implement policies that balance fairness and efficiency by exploiting the correlation between desirable *style* features and borrowers *types*, our paper relates to the fairness-efficiency trade-off literature. There is ample empirical evidence that the implementation of more efficient algorithms can exacerbate inequities (Lepri et al., 2018; Williams et al., 2018; Zhang et al., 2021; Zhang and Yang, 2021). Only a few papers compare algorithms based on their impact on fairness and efficiency; Rhue and Clark (2020) simulate a marketplace and counterfactually adjust algorithmic decision thresholds to highlight the tension between fairness and core business outcomes in an online crowdfunding platform. In the context of criminal sentencing, Kasy and Abebe (2021) present a theoretical model and calibrated simulations to show how various algorithms impact the race gap. We contribute to this literature by studying a new class of policies based on features in profile images. We show that the impact of such policies depends on the correlation between the changeable features of images and the fixed characteristics of borrowers depicted in the images. We also demonstrate this point in a counterfactual simulation of various platform policies based on features in images.

We showcase how Generative Adversarial Networks (GANs) can be used to estimate preferences for specific features in images. The pipeline that we propose is particularly suitable for audit studies. Audit studies have been commonly used to study biases in how economic agents make decisions (Mullainathan et al., 2012; Kline et al., 2021; Salminen et al., 2022). Our method applies GANs to address the confounding problem by generating images that differ only in the selected dimension; additionally, GAN-generated images are realistic, which allows us to study the effect of the feature in a life-like setting. Fong and Luttmer (2009) varied racial information in images of Hurricane Katrina victims by showing images of black or white hurricane victims. They adjusted for other dimensions in which the images differ by reducing image quality and controlling for observed characteristics. Ash et al. (2022) used the interaction between the text of a news article and the associated image to measure the extent of racial stereotypes. Flores-Macías and Zarkin (2022), in a conjoint experiment, studied the effect of military uniform, gender, and skin color on the perception of the effectiveness of law enforcement. They used Photoshop to modify the images used in the experiment. GANs improve upon these methods by varying only the selected feature and producing high-quality realistic images; they can also be scaled to a larger set of images at a low cost. Ludwig and Mullainathan (2022) used GANs to morph images in a direction that shifts the choices of decision-makers and then asked subjects to name the changed feature. We use GANs to achieve a different objective. We start by identifying features; we prioritize for study those that appear to have a causal effect in observational

data. Then we modify a selected feature in the images using GANs, and finally estimate the impact of the feature on the choices in a recruited experiment. In contrast, Ludwig and Mullainathan (2022) highlights how GANs can be used to identify impactful features.

3 Empirical context

Microcredit, sometimes called microlending, typically refers to small, uncollateralized loans to low-income households at terms more favorable than otherwise available, often on the premise of supporting micro-enterprise development.⁷

Over the past two decades, online microcredit platforms have broadened opportunities for participation by expanding access to individual lenders. Sun et al. (2019) argue that online microcredit platforms, such as Kiva, help establish personal connections between lenders and borrowers, and encourage participation by simplifying the discovery and lending processes.

While microcredit platforms have enabled the participation of new lenders, concerns about fairness have increased. Past research using field and lab experiments has shown that online peer-to-peer platforms yield considerable inequities across race, gender, and physical attributes of borrowers.⁸

A number of microcredit platforms are non-profit organizations, including Kiva. Among non-profit organizations, the tension between fairness and efficiency is particularly nuanced. On the one hand, improving efficiency makes it possible to expand lending services to more borrowers, who are otherwise excluded from access to credit; on the other hand, it changes the dynamics of fund access and might create inequality among different types of borrowers.

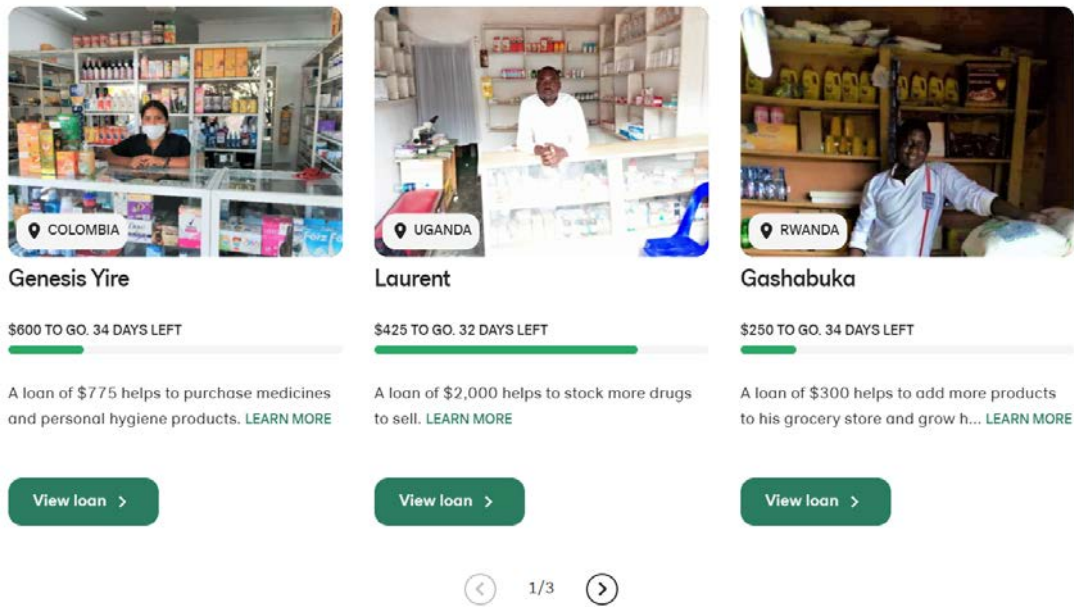
Kiva. Serving borrowers in more than 80 countries, Kiva is one of the most prominent online, non-profit, peer-to-peer microcredit platforms. Since its founding in 2005, Kiva has issued over 1.6 million loans funded by over 2 million lenders, totaling 1.7 billion U.S. Dollars. Kiva is an online marketplace in which borrowers have their own profile pages, with pictures, that prospective lenders can browse to select the borrowing campaigns they want to invest in. Kiva collaborates with local microcredit agencies in vetting, curating, and promoting borrowers.

A potential lender starts the search on the category page. Figure 1 shows an example. Lenders can obtain more information by clicking on "View loan," where they learn more about the loan objective

⁷See Karlan and Morduch (2009) for an academic overview of microcredit, and see Banerjee et al. (2015) for a summary of six randomized controlled trials of microcredit.

⁸Fong and Luttmer (2009) document that lenders prefer borrowers of the same race; Landry et al. (2006) show the role of gender and Park et al. (2019) of perceived beauty.

Figure 1: Kiva category page.



Note: Screenshot from [kiva.org](https://www.kiva.org) collected on 3/3/2022.

and geographical location.

Images play a prominent role in lenders' discovery of borrowers and they help borrowers to tell their stories. For lenders they are a key factor when deciding whether or not to invest (Park et al. (2019)). Borrowers provide their own photographic images, and these vary in quality, content and composition. Some show mostly the borrower, while others focus on their business; facial expressions of borrowers, e.g., serious or smiling, and technical aspects such as quality of lighting or resolution, tend to vary too. To help the borrowers with this important component of their application, Kiva recommends, inter alia, that photos should be high resolution with a horizontal orientation, and they should include the business owner and the business in the background.⁹

4 Analysis of observational data

Our framework for studying fairness and efficiency requires that we establish that lenders have preferences for specific profile *style* features and that they are differentially distributed across borrowers' *types*. To do that we use data on the historical performance of borrowing campaigns on Kiva.

⁹See <https://www.kivaushub.org/profile-photo>

4.1 Kiva data

We construct *Kiva data* by combining three datasets: a publicly available dataset with loan characteristics and lending outcomes, data on features in images associated with the borrowing campaigns, obtained using the methodology described in Appendix A, and a dataset on repayments that Kiva generously shared with us.¹⁰

The publicly available data on characteristics and outcomes of borrowing campaigns spans April 2006 to May 2020 and contains over half a million observations. The data includes key characteristics of each borrowing campaign such as sector, name of activity, country, funding goals, and currency. We have several measures of funding outcomes: the amount of money collected per day, the number of days it took to raise the capital (campaigns generally stay active until they collect all funds), and the number of lenders that loaned money to the borrower. We primarily focus on money collected per day as an outcome metric because we are interested in analyzing how lenders allocate their capital between borrowers, to estimate their preferences for specific features of borrowing campaigns. We characterize the competitive landscape by exploiting the fact that our data contains all borrowers active in the covered period. For each borrower, we compute the number of borrowers from the same country and sector listed concurrently. We also include the share of borrowers of the same *race* and *gender*.¹¹ Finally, to flexibly capture time trends, we introduce interactions between the month in which a campaign was posted and the sector, and interaction between the country and the month.

Images used in borrowing campaigns are also publicly available. We use the feature detection algorithm, Convolutional Neural Network (CNN), described in Appendix A, to detect features in profile images and enrich the funding outcomes dataset. The algorithm that we use takes the image as an input and returns a vector of probabilities associated with a pre-defined list of features. From the CNN, we obtain around 140 features in images: various objects in the image, technical aspects of the photo (*blurry, flash*), individuals' facial expressions, and other personal characteristics like *race* and *age*. However, not all of these features are useful in our context. First, many of the features do not or very rarely appear in our dataset. To reduce data size, we remove features that appear in only 0.01% of images or less. By doing so, we mostly drop features describing specific objects in the image (e.g., *cup*). Second, we drop several features that are highly correlated (e.g., *frowning* and *smiling*) since such features mostly duplicate information.¹² In the end, we focus on 55 features in images. The full list is

¹⁰See here: <https://www.kiva.org/build/data-snapshots> for the publicly available dataset.

¹¹*Race* and *gender* are predictions based on campaigns' profile images.

¹²When several features have a Pearson correlation coefficient above 0.75, we select one of them.

available in Appendix A.

In Appendix B, we document the results of an audit study of the model performance. We compare image labels annotated by recruited human raters with the algorithm prediction. We conclude that the algorithm’s predictions are highly correlated with labels given by humans. Additionally, we document that the model is more likely to make false negative errors than false positives (i.e., the model predicts that a feature is not present in the dataset when human raters label the image as having this feature). Consequently, the estimates based on the observational data are likely to understate the impact of features on outcomes. We evaluate the extent of this potential bias in Appendix B.

The feature detection algorithm estimates the probability that a feature is present in the image; depending on the application, we either use the continuous value or a binary indicator taking the value of one when the probability exceeds 50% and zero otherwise. We use *italics* when referring to demographic features predicted using CNN.

Out of the 55 features we obtain, some cannot be changed when borrowers create their profiles (e.g., demographics).¹³ We categorize such features as *type*. We treat *types* as the collection of features from profile images and the description of the borrowing campaign (e.g., country). We label features that borrowers can change when creating their profiles as *style*. We categorize the following features as *style*: *No Eyewear*, *Sunglasses*, *Smile*, *Blurry*, *Eyes Open*, *Mouth Wide Open*, *Blurry*, *Harsh Lighting*, *Flash*, *Soft Lighting*, *Outdoor*, *Partially Visible Forehead*, *Color Photo*, *Posed Photo*, *Flushed Face*, *Top* (person’s face in the top part of the image), *Right* (person’s face in the right part of the image), *Bottle* (there is a bottle in the image), *Chair* (there is a chair in the image), *Person* (there is another person in the image), and *Body-shot* (the body of the borrower occupies a substantial part of the image).

Finally, the data on defaults spans from 2006 until 2016 and contains approximately 420 thousand borrowing campaigns. The unit of observation in this dataset is a loan from an individual lender (note that generally borrowing campaigns have multiple lenders). We have information on whether each loan has been repaid or not.¹⁴ We aggregate the dataset to the borrowing campaign level and consider the loan to be repaid if the borrower paid back the money to all lenders. The main variable is whether borrowers defaulted or paid all their loans.

¹³Of course, it may be possible for borrowers to shift perceptions to create ambiguity about their type features. We abstract from such manipulation, but caution that if platform participants learn that they could benefit from it, they may modify how they present their *type* features.

¹⁴We also have information on who defaulted on the loan: (i) defaults by the borrower are 75% of cases, (ii) defaults by the micro-finance partner are 23% of cases, and (iii) defaults by both are 2% of cases. Each of these categories has the same impact on the lender: the loan is not repaid. Thus, in the main analysis, we do not distinguish between the reasons for the default. See Appendix F for further discussion.

We merge the three datasets, and our complete dataset captures the period from 2006 to 2016.

Table 1 presents summary statistics of the main variables. The full list is available in Appendix 7.

Table 1: Summary statistics of the main variables

Statistic	N	Mean	St. Dev.	Min	Pctl(25)	Pctl(75)	Max
cash per day	420,765	104.587	136.378	1	25	116.7	621
days to raise	420,765	13.175	10.947	1	5	20	38
default	420,765	0.050	0.218	0	0	0	1
loan amount	420,765	800.107	993.370	25	275	950	50,000
no. competitors	420,765	0.091	0.173	0.003	0.006	0.075	1.000
share same race and gender	420,765	0.665	0.294	0	0.4	1	1
<i>male</i>	420,765	0.198	0.398	0	0	0	1
<i>smile</i>	420,765	0.498	0.177	0	0	1	1
<i>body-shot</i>	420,765	0.406	0.491	0	0	1	1

Note: Summary statistics of selected variables. Cash per day and days to raise are Winsorized at the top 97th percentile. Cash per day and loan amount in USD dollars; male and smile take the value of 1 when CNN predicted probability is above 0.5 and zero otherwise. No. competitors is the number of borrowing campaigns from the same sector and country posted concurrently; the value is standardized by the maximum.

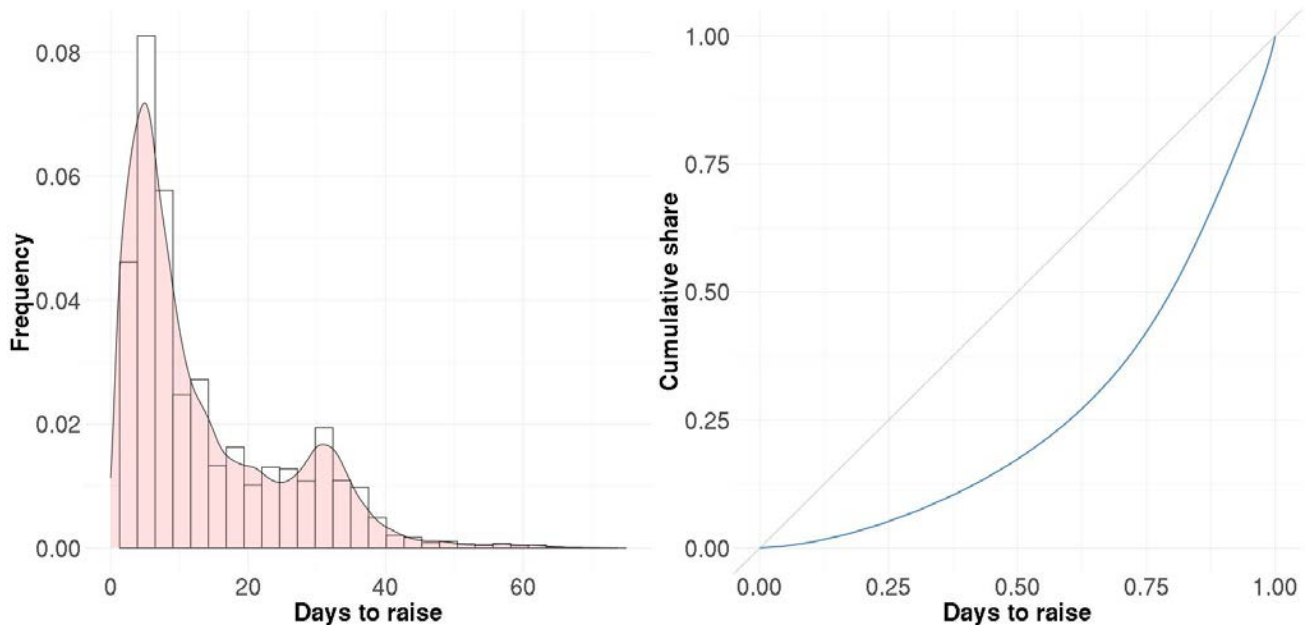
4.2 Inequities in funding outcomes

While loans on Kiva stay active for a long time so that the vast majority of them eventually get funded, there is substantial variation in how long it takes to reach campaigns' funding goals or how much money is collected per day. In Figure 2, we show a histogram of the number of days it takes to collect the entire amount (*days to raise*) and a Lorenz curve documenting inequity in this outcome. If every borrowing campaign took the same number of days to get funded, the blue (actual distribution) and gray (perfect equality) curves would overlap.

From the left panel of Figure 2 we can observe that there is substantial dispersion in how long it takes to receive commitments for the entire amount of the loan. The mean outcome is 14.5 days, but many campaigns are fully funded almost immediately while others take over a month to reach their funding goals.

An important driver of how long it takes to raise all of the funds is the size of the loan. Thus, a useful measure of how quickly borrowers raise funds is the amount of money raised per campaign per day. In Figure 3, we show a histogram of funds in dollars collected per day (*cash per day*), y-axis shows the share of campaigns attracting the sum of money shown on the x-axis, and an associated Lorenz curve. There is a substantial variation in *cash per day*. The mean is 118 USD, but many campaigns raise just a few dollars per day. Focusing on the Lorenz curve (right panel) we observe even higher

Figure 2: *Days to raise*: histogram and Lorenz curve



Note: Left panel - histogram of days to raise capped at 75 USD. The fitted density curve is shown in pink. Right panel - Lorenz curve of days to raise.

inequities than in Figure 2.

The evidence presented in Figures 2 and 3 is based on data collected over ten years. Much of the variation can be due to time trends such as differences in the number of available lenders or borrowers. In Figure 4, we group campaigns into weekly intervals such that campaigns that were listed online during the same week are in the same group. Thus, a group of borrowing campaigns approximates a choice set available to lenders that were active in that week.¹⁵ We use two measures of inequity, the Gini coefficient and the sum of market shares of the 33% of borrowers with the lowest amount of money collected per day.¹⁶ The Gini coefficient of 0 expresses perfect equality, where all values are the same. The market share of the bottom third amounting to 33% would indicate that the outcomes are equally distributed across tertiles. We can observe that both metrics reveal that outcomes are far from being equally distributed.

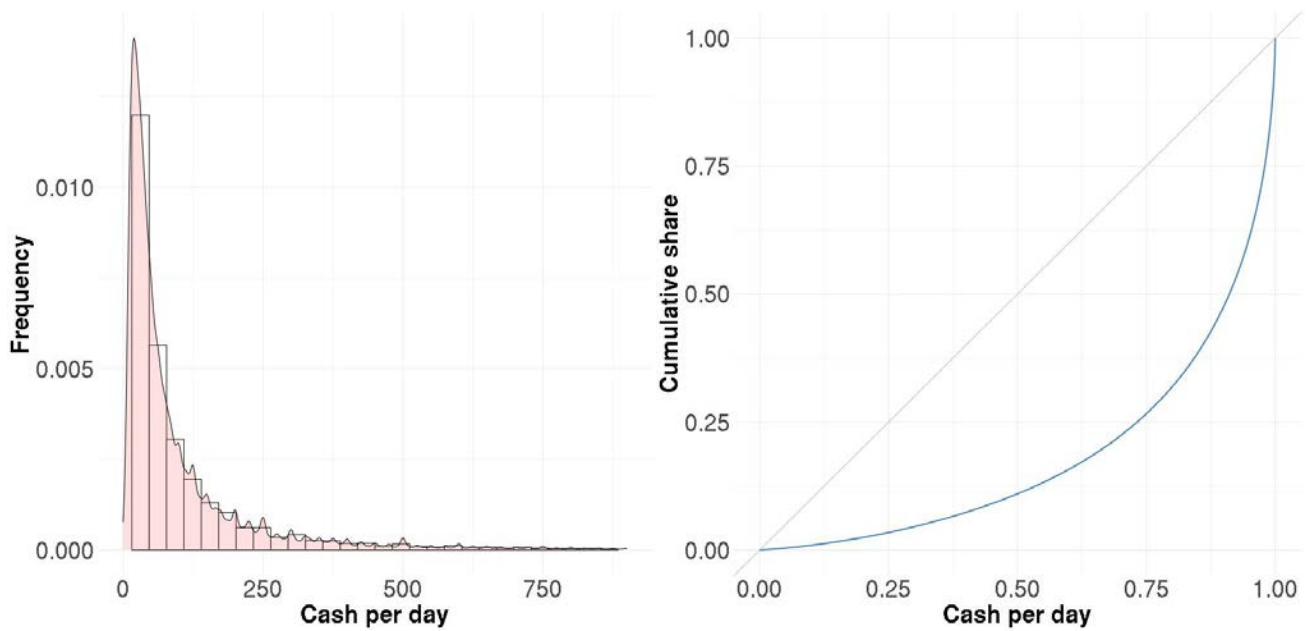
¹⁵We say that this only approximates the choice set of lenders because, in some cases, a borrowing campaign could have been posted at the beginning of the week, quickly collect all the funds and disappear from the platform. Thus, a lender active only at the end of the week would not see such a campaign.

¹⁶The Gini coefficient is defined as

$$Gini = \frac{\sum_{j=1}^n \sum_{j'=1}^n |x_j - x_{j'}|}{2n\bar{x}}$$

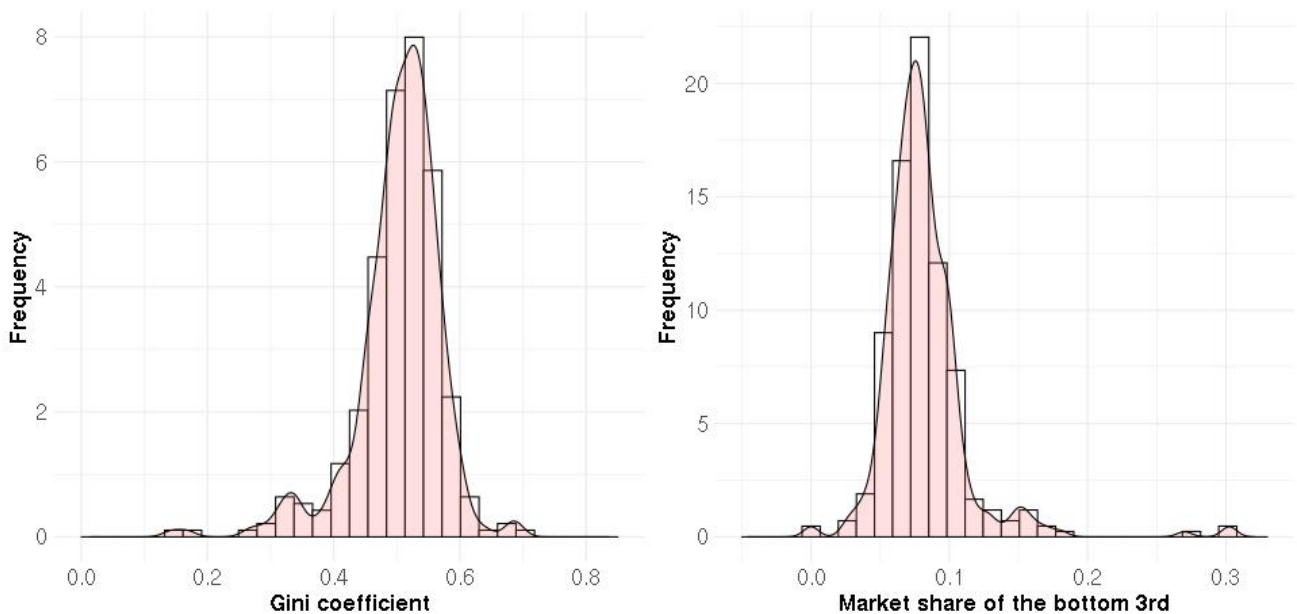
where x_j is the outcome for borrower j and $x_{j'}$ for borrower j' , n is the number of borrowers available in that week and \bar{x} is the average cash collected per day.

Figure 3: Cash per day: histogram and Lorenz curve



Note: Left panel - histogram of cash per day capped at 1250 USD. The fitted density curve is shown in pink. Right panel - Lorenz curve of cash per day. On average there are 450 borrowing campaigns active in a given week, which together raise on average over USD 400,000 in loans.

Figure 4: Cash per day distribution within weeks: Gini coefficient and share of the bottom tertile.



Note: Statistics in both panels are computed on a weekly basis. Left panel - Gini coefficients of weekly distributions of cash collected per day. Right panel - weekly sums of cash collected per day by the 33% lowest performing borrowers.

4.3 Profile images and outcomes

Funding outcomes. Evidence presented in Figures 2, 3, and 4 indicates that there are substantial inequities in funding outcomes across borrowing campaigns. The inequities can be due to differences in borrowers' *types* and their *style*. *Style* features are central to this analysis because a platform can design interventions to modify them. In Table 2, we show that part of the variation in outcome can be explained by *style* features in images.

We train three models to predict *cash per day*. The first is a full model which includes all variables in *Kiva data*; the second is a restricted model that contains only *style* features, and the third is a benchmark mean model. We use a gradient boosted machine (Friedman (2001)) for the predictive model.¹⁷ We split the dataset 70:30 into train and test. Table 2 reports predictive performance in the test set measured using mean squared error. We find that including *style* features improves the predictive performance of this model compared to a mean model. Additionally, a full model based on all covariates in *Kiva data* performs better than the model with only *style* features.

Table 2: Image features as predictors of *cash per day*.

specification	MSE	SE
Mean	22367	252
<i>Style</i> features	19373	224
Full model	10996	138

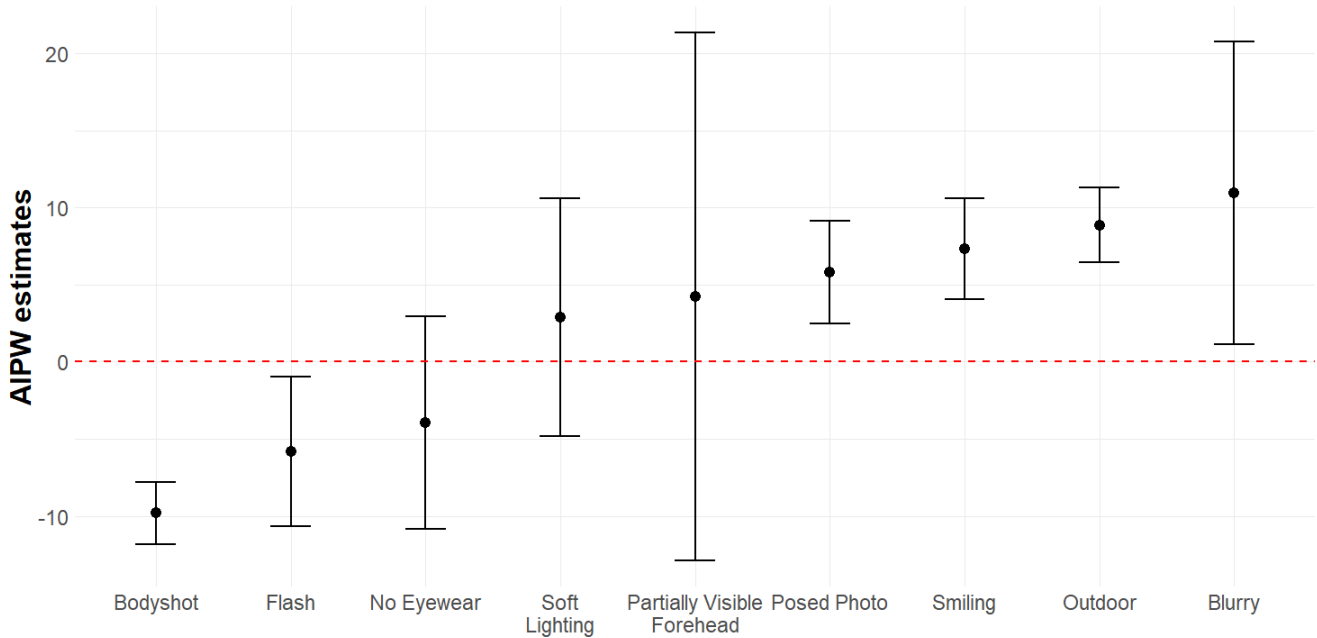
Note: Test set performance of a gradient boosted machine (GBM) trained using all available covariates (full model) models with only image style features and a mean model. Models trained on 70% of data and tested on 30%. Mean squared errors are in the second column. Standard errors of MSE are in the third column.

Specific style features. Results presented in Table 2 show that both *style* and other image features jointly are predictive of funding outcomes. However, if Kiva is to construct platform policies around *style* features, it is important to select individual impactful features. In other words, we want to know: "what would happen if a profile was presented with a change in one characteristic and remained unchanged otherwise." We want to know the average treatment effect (ATE) of a specific feature in an image.

To estimate ATEs we use the Augmented Inverse Propensity Weighing (AIPW) estimator (Robins et al., 1994; Glynn and Quinn, 2010; Wager and Athey, 2018). AIPW is a doubly-robust method: it

¹⁷Unless stated otherwise, we use the gradient boosted machine for all predictive tasks. We selected the model based on a comparison of test set performance with other popular predictive models. See Appendix E for a detailed discussion.

Figure 5: Estimates of the average treatment effect of selected style features



Note: Estimates of the average treatment effect of selected features on cash collected per day. x-axis selected features, y-axis ATE estimates. 99.9% confidence interval. Propensity and outcome model estimates using Regression Forest. We transform the treatment variable to a binary variable that takes the value of one when the predicted probability of the feature is above 0.5 and zero otherwise. See Appendix G for results using GBM and diagnostics.

adjusts for covariates in the outcome model and the propensity score. We use the *grf* implementation of the AIPW estimator (Athey et al. (2019)).

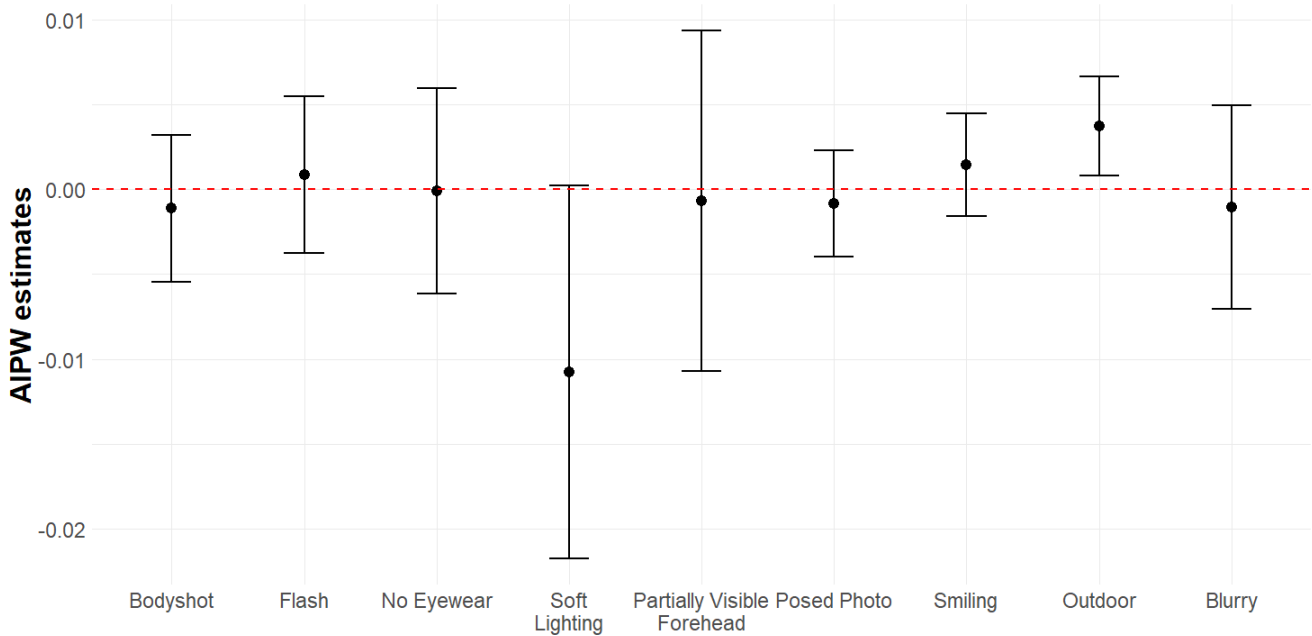
The feature detection algorithm detects a very rich set of characteristics, but we may still be missing some variables (both features in images and other variables) that might influence lenders' decisions. Thus, the ATE estimates should be interpreted as comparing profiles that are similar in all observed dimensions other than that which is studied. We return to this issue in Section 5, where we present experimental evidence corroborating the impact of selected features on outcomes.

Figure 5 shows estimates of average treatment effects on *cash per day* for selected features. We find that several features have negative ATE e.g., *flash*, *body-shot*, while others like *posed photo* or *smile* have positive and statistically significant effects.¹⁸

Loan repayment. Lenders might use *style* features as signals of the probability of repayment: we show that *style* features do not predict whether the loan will be repaid. In Table 3, we compare the predictive performance of default models with and without image features. We see that the inclusion

¹⁸See Appendix G for diagnostics.

Figure 6: Estimates of the average treatment effect of selected style features



Note: Estimates of the average treatment effect of selected features on probability to default. *x*-axis selected features, *y*-axis ATE estimates. 99.9% confidence interval. We transform the treatment variable to a binary variable that takes the value of one when the predicted probability of the feature is above 0.5 and zero otherwise.

of style features does not improve the predictive performance of the model compared to a mean model.

Table 3: Image features as predictors of default probability.

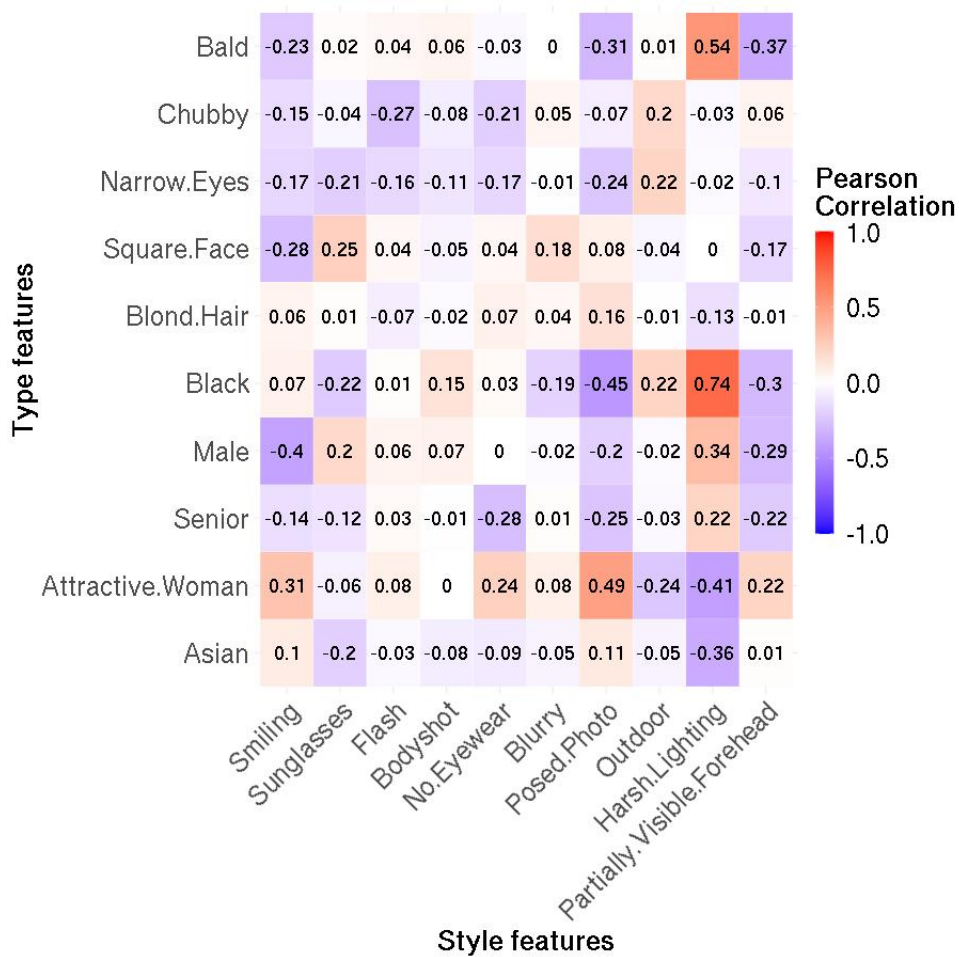
specification	MSE	SE
Mean model	0.065	0.001
Style features	0.064	0.001
Full model	0.059	0.001

Note: Test set performance of a gradient boosted machine (GBM) trained using all available covariates (full model) and simplified model using image style features (Style features) and a model with only an intercept (Mean model). Models trained on 70% of data and tested on 30%. Mean squared errors are in the second column. Standard errors of MSE are in the third column.

We also estimate the average treatment effects of individual features on the probability of default. Figure 6 shows the results. The only feature with a non-zero statistically significant impact on repayment probability is *outdoor*. We take these results as evidence that *style* features are not useful signals of repayment probability.¹⁹

¹⁹Our results contrast with findings in Netzer et al. (2019), which shows that free text in listings on a crowd-funding platform Prosper is predictive of the default probability, adjusting for financial and socio-demographic characteristics.

Figure 7: Correlation between selected *type* and *style* features.



Note: Pearson correlation coefficient between selected style features in columns and type features in rows.

4.4 Correlation of *types* and *style* features

Style features can aggravate or mitigate inequities in outcomes due to differences in *types*. When borrowers with *type* features associated with high outcomes build profiles that attract lenders, the disparities due to *types* will increase further; in contrast, if borrowers with less desirable *type* features choose attractive profile *styles*, outcomes will be more equitable. In this section, we document the correlation between *types* and *styles*.

Figure 7 shows a correlation between selected *type* and *style* features. Some of the features are highly correlated; for example, *smiling* is less common amongst *male* and *bald* and more common for *attractive woman*.

To argue that the choices of *style* features exacerbate inequities due to *types*, we need to show that the distribution of *types* across borrowing campaigns results in inequity of outcomes and that the

desirable *style* features are more prevalent amongst borrowing campaigns whose *types* lead to better funding outcomes.

To do that we carry out a *Gelbach Decomposition* (Gelbach, 2016) of selected *type* variables. This is a method for measuring the extent to which adjusting for a group of variables changes the coefficient of a selected variable. It tells us what the coefficient would be if the means of the adjusted variables were the same across the levels of the evaluated variables.

Table 4 presents the results for selected *type* variables. The first column shows the name of the variable. The second column shows the coefficient associated with the selected variable from a linear regression of the variable on *cash per day*. In the third column, we see the coefficient adjusted for all variables in *Kiva data*. In the final column, we have the adjustment due to *style* features. For example, if we partial out differences in the distribution of *style* features, profiles with images of *bald type* receive USD 10.21 less than those of *not-bald type*. Thus, differentially distributed *style* features aggravate the disparity between *bald* and *not bald*. Finally, as evidenced in Table 3, the additional inequity is not explainable by differences in repayment probability, and in this sense is unfair to the borrower.

Table 4: Impact of style features on coefficients associated with *types*

feature	Coefficient base	Std. error base	Coefficient full	Std. error full	Delta style
<i>Bald</i>	-71.57	4.86	-22.14	7.31	-10.21
<i>Chubby</i>	16.92	2.02	-7.65	4.17	2.08
<i>Narrow Eyes</i>	-9.14	1.87	3.77	3.66	2.01
<i>Square Face</i>	-87.70	8.06	-27.64	10.46	-52.59
<i>Black</i>	-12.06	1.37	-1.61	3.51	4.61
<i>Senior</i>	-57.18	4.70	-6.02	5.96	-6.02
<i>Attractive Woman</i>	75.32	2.43	10.72	4.13	9.48
<i>Asian</i>	12.39	1.52	0.72	2.74	2.30

Note: Gelbach decomposition of selected type features (Gelbach, 2016). Coefficient base refers to coefficient of a univariate model with the selected type feature; coefficient full coefficient from a model adjusting for all covariates in Kiva data; delta style impact of style features on the disparity between types. R implementation by Stigler (2018)

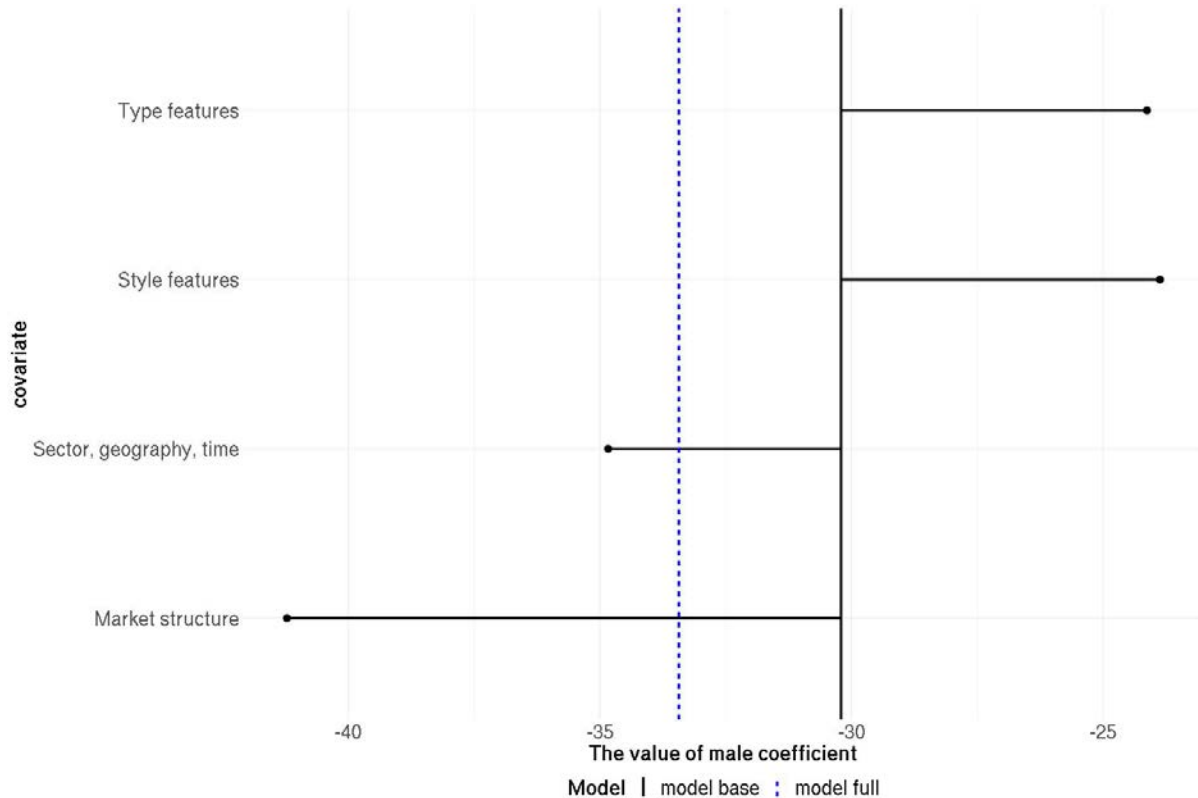
Results presented in Table 4 indicate that *style* choices aggravate disparities due to *bald*, *square face*, *senior*, *attractive woman*, and *Asian types* and mitigate for *chubby*, *narrow eyes*, and *black types*.

4.5 Gender gap

A specific *type* feature that matters in our context is *gender*.²⁰ We find that campaigns classified as *male* raise on average USD 36 less per day and take 5.8 more days to be funded fully (differences in

²⁰We use an algorithmic prediction of *male*. Thus, the variable *male* indicates that the feature detection algorithm assigns a probability of at least 0.5 that the person in the image is a male.

Figure 8: Gelbach decomposition of *male* coefficient



Note: The solid line is an estimate of the coefficient associated with male from a univariate linear regression; the dashed line is the coefficient adjusted for all variables in Kiva data; OLS estimator. Horizontal lines represent contributions of the variables group to the coefficient associated with male; type features include all other type features from the image; sector, geography, time includes sector, country, week fixed effects, loan amount, and repayment details, and market structure includes interactions of month and country, month and sector, number of lenders in the week, number of competing campaigns, and share of campaigns of the same race and gender. The R implementation of Gelbach (2016) by Stigler (2018).

means).²¹

As evidenced by the results presented in Table 4, the differences between *types* can be due to the non-equal distribution of other characteristics, e.g., *style*. To decompose the unadjusted difference we perform another Gelbach Decomposition (Gelbach, 2016). The proposed method compares a baseline model with only a *male* indicator variable with a full model that includes all the variables in the *Kiva data* and decomposes the contribution of all the added variables to the changes in the coefficient of interest. Figure 8 presents the results.

²¹In the context of microfinance, the gender gap might be driven by users that aim to correct for discrimination against women in traditional finance. There is a rich literature documenting discrimination against women in traditional entrepreneurial lending. Alesina et al. (2013) shows that women entrepreneurs pay higher rates for access to credit and Brock and De Haas (2021) use a randomized experiment to show that loan officers grant loans to women under less favorable conditions than to men. The phenomenon of over-correcting for discrimination is well documented in experimental psychology (Mendes and Koslov (2013), Nosek et al. (2007)). It is also plausible that Kiva lenders follow broader policy discussions, where the emphasis on developmental policies and aid targeting women is common (Kristof and WuDunn, 2010).

Figure 8 depicts the differences in the coefficient associated with *male* between a univariate linear regression (solid line) and a full model which includes all variables from the *Kiva data* (dashed line). The length of the horizontal arms going from the base model indicates the contribution of each variable group to the *male* coefficient in the full model. Thus the length of the horizontal line is the partial effect of the unequal distribution of features within the group. We can observe that changing the distributions of *style* features would decrease the gender gap; additionally, we find that *male* campaigns also have a non-desirable distribution of other *type* features, but the distribution of *sector*, *geography*, *time* and *market structure* decreases the *gender* gap.

Additionally, we look at the prevalence of the desirable and non-desirable features. For example, the frequencies of *body-shot* and *smile* differ substantially between genders: 77% of *female* borrowers *smile* in the image, compared to 33% of *male* borrowers. 26% of *male* borrowers use a *body-shot*, compared to 22% of *female* borrowers (both are statistically significant).

5 Recruited experiment

Our efficiency-fairness framework has two components: the differences in distributions of *style* features across borrowers who succeed in quickly funding their projects and those who do not and the effects of these features on securing funding. The causal interpretation of the estimates of treatment effects from Section 4 rests on the assumption of unconfoundedness, which is difficult to verify. This section provides experimental evidence of the impact of two selected style features, *smile* and *body-shot*, on outcomes. We selected these features because of high and statistically significant estimates of their impact on *cash per day*, high correlation with borrowers' types, and differences in their prevalence amongst *male* and *female* borrowers.

5.1 Experiment design

In the experiment, recruited subjects are presented with a series of pairs of images of borrowers and asked to select one out of each pair. The images are fabricated to have exogenous variation in *smile*, *body-shot*, and *male*. The objective of the experiment is to estimate the treatment effects of the three features, *smile*, *body-shot*, and *male*. The design of the experiment builds on the literature on conjoint analysis (Hainmueller et al. (2014)) with the novelty that variation in features of interest is encoded in images.

Figure 9: Variation in *smile*



Note: Two versions of an image with variation in smile. Both images were generated using GANs.

Figure 10: Variation in *male*



Note: Two versions of an image with variation in male. Both images were generated using GANs.


Images. To generate images that differ in the selected features we use Generative Adversarial Networks (GANs). GANs, designed by Goodfellow et al. (2014), is an approach to generative modeling that uses deep-learning methods. The key objective of GANs is to generate fabricated data that are similar to real data. GANs are frequently used to modify images and generate so-called "deep fakes". GANs have been used in social sciences to generate realistic images (Ludwig and Mullainathan, 2022) and synthetic datasets (Athey et al., 2021). For our task, we apply Style-GAN developed by Karras et al. (2019). Specifically, we use GANs to vectorize a selected feature in images. Once the feature has been vectorized, we can adjust the vector in the desired direction. The modified attribute is then embedded into the original image while the rest of the image stays unchanged. Finally, we ensure that images look realistic by deblurring, inpainting, and auto-blending. See Appendix C for further discussion, and Figures 9 and 10 for examples of GAN-generated images.

Experiment implementation. In the experiment, subjects are introduced to the concept of micro-loans and then presented with six pairs of borrowers and asked to select one in each pair. Figure


Figure 11: An example of a choice instance

Question 2

This page shows two loan campaigns. Please choose the one you would prefer to lend money to.



Sai, India. Total loan: \$900. Invest in his grocery stall.



Bayu, Indonesia. Total loan: \$1000. Invest in his farm.

Note: An example of choice instance from the recruited experiment. Both images show borrowing campaign profiles featuring males. The left profile is not a body-shot and not smiling. The right panel is not a body-shot and the borrower is smiling.

11 shows an example of a choice instance. Participation in the experiment took approximately five minutes. The survey included several features to encourage thoughtful responses. We provide more details in Appendix H.

To generate experimental protocols, we first create a pool of images. Starting with 20 original images, we generated artificial versions with variations in *male*, *smile*, and *body-shot*. This gives us eight versions of each image.²² All images used in the experiment were artificial, GAN-generated versions of the original images. This means that subjects were choosing between two fabricated images. In the rest of this section, we use the term ‘profile’ when referring to all eight variants derived from the same image.

Second, to allocate images to protocols we draw the first image, without replacement, and pair it with a second image that was not generated from the same original image and that differs in at least one feature. We repeat this until we have six pairs. In total, we created 15 protocols.

²²For privacy and ethical reasons we do not modify images of Kiva borrowers, whom we cannot contact to seek consent to modify their images. Instead, we purchase images from *Shutterstock.com*, a website that sells images. We select images that are similar to images used by Kiva borrowers and use these purchased images to train GANs and alter images in the desired way.

Sample recruitment. The experiment was hosted on Prolific.co. We recruited 400 subjects who declared that they had contributed to a charitable cause in the previous year. We considered subjects from developed countries with high socioeconomic status (self-reported). We impose these criteria to ensure that our subjects are similar to Kiva lenders.

The mean age of subjects in the experiment is 33 years and 51% are women. Subjects were asked to declare the amount they donated to charity in the previous year. We retained as eligible those who donated at least USD 1; 60% of respondents donated less than USD 75. The United Kingdom is our subjects' most common country of residence, with 40% of subjects, followed by Spain with 30%, and France with 9%. In Appendix I, we present summary statistics on employment and self-reported socioeconomic status.

Note that the subjects in our experiment are not actual lenders on Kiva. Thus, the preference estimates obtained in our experiment might not reflect the preferences of Kiva users. The estimates should be viewed as additional evidence corroborating the importance of selected features in choosing between microlending campaigns, but not as definitive evidence of their impact on outcomes on the Kiva platform.

5.2 Experiment results

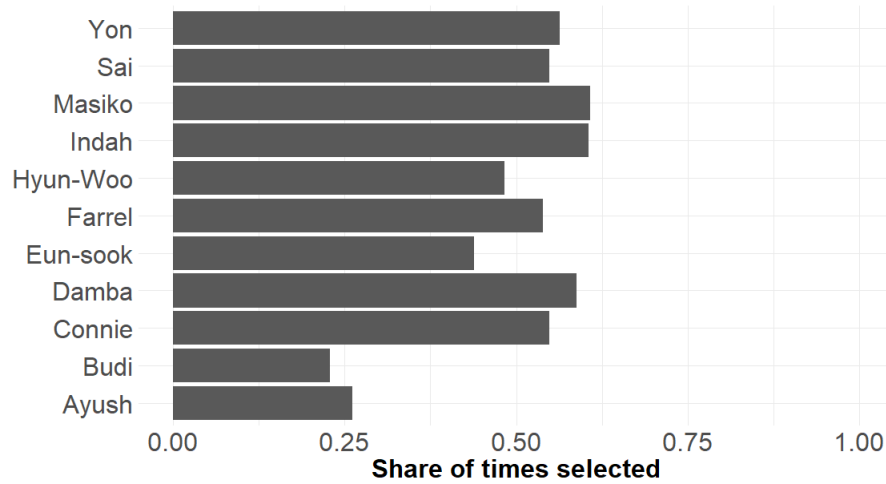
Subjects in the experiment choose between two profiles in each choice instance. The outcome is whether the profile is chosen or not, and its mean is 0.5. We start by evaluating whether some profiles are, on average, selected more frequently than others. Figure 12 shows the average outcome per profile. We observe clear differences in the popularity of profiles, but none of them was selected almost always or almost never. We keep all of them for the analysis (for robustness, we also carry out the analysis without the two least popular profiles).

Average treatment effects. Subjects were given a pair of profile pictures and asked which one they preferred. Suppose the preference has a systematic component, which depends on *male*, *smile*, and *body-shot*, and there is also an idiosyncratic random component ϵ_{ij} . The utility of subject i choosing the option j is written:

$$u_{ij} = \alpha * male_j + \beta * smile_j + \gamma * body - shot_j + \mu_j + \epsilon_{ij}. \quad (1)$$

Assuming that ϵ is distributed following a type I extreme value distribution, the probability of a

Figure 12: Mean outcomes per profile



Note: Vertical axis is the borrower profile name (we group male and female variants together but show one name). Horizontal axis is the number of times the profile has been selected divided by the number of times the profile was shown.

subject i choosing option j , is written:

$$u_{ij} = \frac{\exp(\alpha * male_j + \beta * smile_j + \gamma * body - shot_j + \mu_j)}{\sum_{k=j,j'} \exp(\alpha * male_k + \beta * smile_k + \gamma * body - shot_k + \mu_j)}. \quad (2)$$

We are interested in the estimates of parameters α , β , and γ . We obtain them by estimating a logistic regression model by maximizing the conditional likelihood.

Table 5 presents the results; the baseline specification is in column (1), column (2) additionally adjusts for subject-specific covariates, column (3) excludes the least liked profiles (Ayush and Budi), column (4) divides profiles into those with high and low fixed effects and interacts profile features and fixed effects, and column (5) adds subjects' characteristics to column 4.

We find that all the features of interest are statistically significant and have high magnitudes: *male* and *body-shot* lead to lower outcomes, while *smile* leads to higher outcomes. In column (3) where we exclude low fixed effects profiles *body-shot* is no longer statistically significant. In columns (4) and (5), we divide the profiles into highly attractive and least attractive, based on mean outcomes. Two profiles fall into the latter category: Budi and Ayush. We find that the point estimates of the average treatment effects go in the same direction for both profile types. The average marginal effect of *male* is a 31% reduction in the probability of being selected; there is a 17% drop for the *body-shot*, and a 34% increase for *smile*.²³

²³The estimates of the impact of *male* are very similar to those based on observational data. While in the experiment, we estimate a considerably higher impact of *smile* on outcomes compared to the observational data. In Appendix B, we document that the feature prediction model has a substantial rate of false negatives in the task of detecting *smile*, and that

Table 5: Average treatment effects estimates

	<i>Dependent variable:</i>				
	chosen		chosen	chosen	
	(1)	(2)	(3)	(4)	(5)
male	-0.385*** (0.079)	-0.373*** (0.080)	-0.747*** (0.092)		
smile	0.298*** (0.074)	0.331*** (0.078)	0.554*** (0.088)		
body shot	-0.191** (0.079)	-0.160* (0.084)	-0.121 (0.096)		
male × high FE				-0.372*** (0.091)	-0.366*** (0.093)
smile × high FE				0.233*** (0.081)	0.252*** (0.087)
body shot × high FE				-0.176** (0.084)	-0.142 (0.089)
male × low FE				-0.337* (0.195)	-0.318 (0.197)
smile × low FE				0.822*** (0.276)	0.944*** (0.307)
body shot × low FE				-0.381* (0.226)	-0.414* (0.241)
Image FE	x	x	x	x	x
Subject’s characteristics		x	x		x
Restricted sample			x		
Observations	4,920	4,644	4,142	4,920	4,644

Note: Estimates of the logistic regression. Columns (2) and (5) include features of subjects. Columns (4) and (5) divide profiles into high and low fixed effects (low FE are the two least popular profiles - Budi and Ayush). Column (3) restricts the sample to high fixed effects profiles only. * $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$

To summarize, evidence from the recruited experiment corroborates findings from the observational data. While neither analysis is conclusive, together they support a conclusion that the selected *style* features do have an impact on outcomes on Kiva.

6 Efficiency-fairness tradeoff: counterfactual simulations

There are many different platform policies that could exploit our finding that certain *style* features impact outcomes. In this section, we propose several such policies, simulate counterfactual outcomes, and evaluate their impact on fairness and efficiency. To do that, we consider a simplified model of interactions on Kiva characterized by the parameters from the recruited experiment.

Although we use a stylized approach, our findings can be used to assess what types of policies are likely to succeed. In practice, our method can be useful in suggesting which policies to prioritize for a randomized experiment.

6.1 A model of a micro-lending platform

Pool of borrowers. The pool of available borrowers is a set of borrowing campaigns from which Kiva selects a subset to display to lenders. The pool of borrowing campaigns can be summarized by a vector of profiles \mathbf{x} , where each element is a profile $x_i = (male_i, smile_i, body - shot_i, \eta_i)$, which

might result in the negative bias in the estimate of the impact of *smile* on outcomes. These false negatives might be one of the reasons for the higher experimental estimates.

describes features of the borrower i ; the first element corresponds to the borrower's *type*: *male* or *female*. The two latter features, *smile* and *body-shot*, pertain to the borrower's *style*, and η_i is a fixed effect which summarizes all other characteristics of the borrower.

The pool of borrowers is exogenously determined and the joint distribution of borrowers' characteristics is denoted as G .

Policy and markets. A market is a set of borrowers shown to a lender. Platform policy \mathbb{H} transforms the joint distribution of borrowers' characteristics from G to H . Specifically, the policy defines $\mathbf{E}_H[f|male, \eta]$ for $f \in \{smile, body - shot\}$ the conditional probability of *style* features in the pool of borrowers. Additionally, a policy applies the probability of being shown to lenders $h : (\eta, male, smile, body - shot) \rightarrow [0, 1]$ to the pool of borrowers. Thus, a policy can be summarized as $\mathbb{H} = \{\mathbf{E}_H[f|male, \eta], \mathbf{h}\}$.

The policies that we consider have two elements: they can impact the distribution of *style* features in the pool of borrowers. Examples of this would be advice on profile creation, a protocol that requires borrowers to upload several images and selects the most compliant one, or behavioral interventions that nudge borrowers to create compliant profiles. Second, a policy can differentiate the probabilities with which borrowers in the pool appear in the market. We allow the platform to condition these probabilities on borrowers' characteristics.

Lenders. Lenders, indexed by j , arrive to the platform and observe available borrowers. They choose the option that maximizes their utility: a borrower or the outside option. The utility associated with choosing one of the borrowers is written:

$$u_{ij} = \alpha_j * male_i + \beta_j * smile_i + \gamma_j * body - shot_i + \eta_i + \epsilon_{ij}, \quad (3)$$

where $(\alpha_j, \beta_j, \gamma_j)$ are random preference parameters, and ϵ_{ij} is a random utility parameter, which is independent across lenders and borrowers, GEV distributed. The utility from choosing the outside option is $u_{oj} = \omega + \epsilon_{oj}$. Lenders choose an option that maximizes their utility.

6.2 Implementation

Markets. We consider a pool of 22 borrowing campaigns and assume that a market can have a maximum of ten borrowers. Campaigns' fixed effects take values of fixed effects estimated in the experi-

ment and the conditional distribution of features G follows the distribution in *Kiva data*.²⁴ Fixed effects in *Kiva data* are estimated as predicted *cash per day* net of the impact of *male*, *smile* and *body-shot*. To construct borrowers' profiles, we treat a fixed effect as a random variable $\hat{\eta}$ drawn from \mathcal{N} a set of fixed effects estimated in the experiment; $\hat{\eta}$ is its realization. Second,

$$\mathbf{E}_G [\textit{male} | D(\hat{\eta})] = \mathbf{E}_K [\textit{male} | D(\eta_k)],$$

where K stands for distribution in *Kiva data*, $D(\cdot)$ is the decile of the fixed effect and η_k is the fixed effect from *Kiva data*. $\mathbf{E}_K [\textit{male} | D(\eta_k)]$ is the share of *male* profiles in *Kiva data* per decile; thus, the share of *male* borrowers with the fixed effect in the first decile of fixed effects estimated from the recruited experiment equals the share of *male* borrowers in the lowest decile of *Kiva data* fixed effects. Finally, *style* features are also distributed following the conditional distribution in *Kiva data*:

$$\mathbf{E}_G [\textit{smile} | \textit{male}, D(\hat{\eta})] = \mathbf{E}_K [\textit{smile} | \textit{male}, D(\eta_k)].$$

Thus, we allow the smiling rates to differ across *male* and fixed effects. Analogously *body-shots* are distributed such that:

$$\mathbf{E}_G [\textit{body-shot} | \textit{male}, D(\hat{\eta})] = \mathbf{E}_K [\textit{body-shot} | \textit{male}, D(\eta_k)].$$

Lenders' preferences. We assume that lenders' preferences $(\alpha_j, \beta_j, \gamma_j)$ are parameters drawn from distributions estimated using experimental data, such that $\alpha_j \sim N(\alpha, sd_\alpha)$, where α is the estimate of the average treatment effect and sd_α is its standard error, and ϵ_{ij} is a random utility parameter, which is iid across lenders and borrowers, GEV distributed. We set the utility from choosing the outside option to one (the highest FE estimated in the experiment is 0.64).

Outcome metrics. We propose two metrics of fairness: first, to capture the overall distribution of outcomes, we use the Gini coefficient defined as

$$\textit{Gini} = \frac{\sum_{j=1}^n \sum_{j'=1}^n |x_j - x_{j'}|}{2n\bar{x}},$$

²⁴Kiva's existing policy is based on when the borrower posted the campaign. Thus, assuming that arrival time is independent of characteristics, a lender sees each borrower in the pool with equal probability. In reality, this is an approximation, because campaigns that reach their funding outcomes are removed from the platform. Thus, the less attractive campaigns stay longer on the platform, so lenders have a higher chance of observing them.

where x_j is the outcome of borrower j and $x_{j'}$ of borrower j' , n is the number of borrowers, and \bar{x} the average outcome. Second, to analyze how the outcomes of the worst performing borrowers depend on the platform policy, we use the sum of market shares of borrowers in the bottom tercile. We will compare the outcomes under various policies to a *fair* benchmark, where the distribution of *style* features does not impact the distribution of outcomes.

We measure efficiency as the share of lenders that chose a borrower instead of an outside option. To compute all metrics we consider all borrowers in the pool. Thus, we capture both borrowers that were included in the market and those that stayed out.

Market outcomes. To determine market outcomes, we simulate markets and choices by lenders. Based on the distribution of outcomes, we compute fairness and efficiency metrics.

Each simulation proceeds in three steps: first, we simulate the pool of borrowers. To do that we draw 22 fixed effects from \mathcal{N} , the pool of fixed effects estimated using data from the experiment, and assign *male* to profiles with the frequency from *Kiva data*. After that, we assign *smile* and *body-shot* following their conditional frequencies.

Second, we construct markets from the pool of borrowers. A policy determines $h(\eta_i, \text{male}_i, \text{smile}_i, \text{body-shot}_i)$, the probability that a borrower in a pool appears in the market. A market is constructed per lender. This means that in one simulation there is one pool of borrowers, from which borrowers are sampled for each lender.

Finally, we simulate lenders' preferences and their choices as described in Equation 3. We perform 50 simulations of 500 lenders' choices for each policy. We use the outcomes to compute our metrics of fairness and efficiency. We consider all borrowers in the pool, irrespective of whether they were shown to lenders or not. Appendix J presents the algorithm that we used.

6.3 Counterfactual policies

Baseline. Baseline policy represents the existing policy on Kiva. In the baseline policy, each borrower in the pool is assigned an equal probability of being included in the market and the joint distribution of features is G , that is, Kiva does influence how *style* features are distributed.

Benchmark. In the *Benchmark* policy, every borrower in a pool has a profile image with a *style* featuring *smile* and without *body-shot*. All borrowers have the same probability of appearing in the market. In the *Benchmark*, we keep the probability of choosing an inside option fixed at the level in *Baseline*. By

doing so we can isolate the role played by the distribution of *style* features in shaping the inequity of outcomes and showcase a fairer outcomes distribution.

Naive. In this policy, we show what happens when a platform realizes that profiles with *smile* and without *body-shot* are more attractive and over-samples them. In practice, when the number of borrowers with *smile* and *body-shot* is more than ten, the platform randomly samples from them. Otherwise, the platform includes all compliant borrowers and fills in the empty slots by randomly drawing additional borrowers. In expectation (i.e., before the pool of borrowers is determined), some non-compliant borrowers are always included.

Partial Compliance. In this policy, the platform recommends to all borrowers to ensure that their profile images have *smile* and do not have *body-shot*. In practice, we assume that previously non-compliant borrowers become compliant with a probability of 75%.²⁵ After a pool of borrowers is determined, the platform assigns an equal probability of being included in the market to all borrowers.

Low-type support. This policy promotes borrowers who are predicted to have low funding outcomes based on their *types*, by ensuring that they are always included in the market. We focus on *gender* in this application. Practically, the approach is analogous to *Naive*: when the number of *male* campaigns is above ten, the platform samples randomly from them. Otherwise, the platform includes all *male* profiles and fills in other slots by randomly selecting from available profiles. In expectation, there are some *female* profiles included in the market.

Restrict Competition. In this policy, the platform promotes fairness by reducing the competition between borrowers. To implement this the platform randomly selects five borrowers from the pool to form the market (instead of ten).

Hybrid. The hybrid policy combines *Partial Compliance* and *Low-type support*.

All the policies that we propose in expectation give non-zero probabilities of being included in the market to any borrower.

Note that *Hybrid* and *Partial Compliance* policies require that borrowers comply with the policy recommendation. In the analysis, we assume that 25% of the borrowers do not adhere to the rec-

²⁵Such a profile feature recommendation can be implemented in various ways, for example, through behavioral nudges or a script requiring that several images need to be uploaded from which platform selects the ones to be shown to lenders.

ommendation. A particular type of non-compliance in the case of *smile* might be that borrowers attempt to create an image with *smile*, but they do not succeed; for example, the *smile* does not appear genuine. In Appendix K, we develop an additional algorithm that distinguishes between fake and genuine smiles and apply it to the Kiva observational data. We show that only genuine smiles lead to higher outcomes. Consequently, the policy will be less effective if some of the newly added *smile*'s are perceived as non-genuine. This analysis highlights the importance of clear instructions and a well-designed system that supports borrowers when they create profiles.²⁶

6.4 Results

Figure 13 presents the results from simulations of the proposed policies. On the horizontal axis, we show the mean of Gini coefficients across all simulations of each policy. On the vertical axis, we show the mean of lenders' shares choosing a borrower rather than the outside options.

We find that for the parameters that we used the proposed policies impact both metrics, the Gini coefficient and efficiency, considerably. First, in the *Baseline* policy, the Gini coefficient is around 0.67 and efficiency is 0.54. Second, the *Benchmark* policy showcases the impact of the unequal distribution of *style* features on fairness; when all borrowers have profiles with the desired features the Gini coefficient reduces to 0.58. Next, *Naive* policy has a strong negative impact on fairness. The Gini coefficient increases to almost 0.8. Recall that under this policy the platform includes more profiles with *smile* and *body-shot* in the market. Unfortunately, this policy has the unintended consequence of reducing the prominence of profiles with *types* associated with lower outcomes, further increasing inequities in outcomes. This is due to the correlation between *types* and *style*. The upside of this policy is that it boosts efficiency, as lenders prefer profiles with the selected features. The other policy focused on increasing prominence, *Low-type support*, has exactly the opposite effect. We observe a decrease in efficiency because the platform now includes in the market more borrowers with *types* leading to lower outcomes. As a consequence, this policy leads to more equitable outcomes: because of random preference components of utility, lenders will choose these borrowers more frequently than in the *Baseline* policy. The alternative pro-fairness policy, *Restrict competition* leads to a small reduction in the Gini coefficient, however, there is a substantial cost to efficiency. With fewer campaigns to choose from, lenders are more likely to choose the outside option.

Two policies, marked in blue, stand out. *Partial Compliance* delivers gains on both dimensions.

²⁶A computer vision algorithm showcased in Appendix K could be a component of such a system, where borrowers could be prompted if their smile is at risk of being perceived as not genuine.

Figure 13: Fairness - Efficiency tradeoff: Gini coefficient



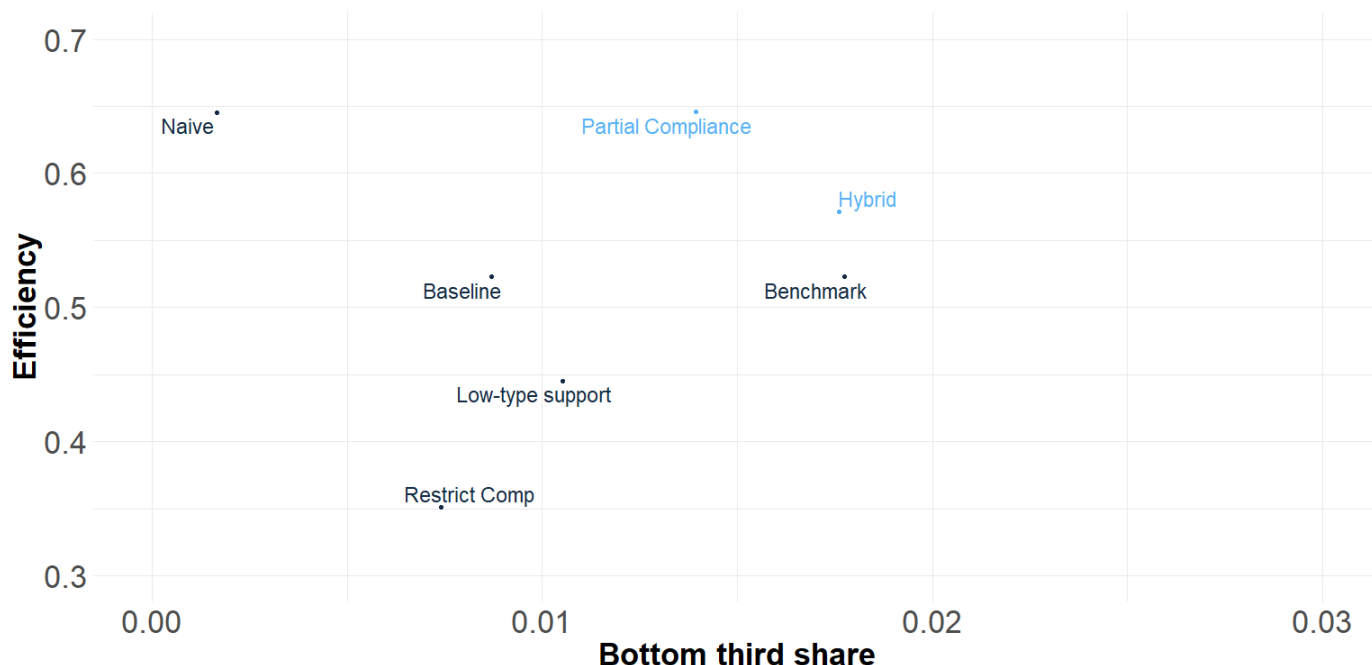
Note: Gini coefficients and efficiency. Each point represents the mean of 100 simulations with 500 lenders each. The horizontal axis presents Gini coefficients while the vertical axis reports the share of lenders choosing an outside option.

On the one side, higher frequencies of *smile* and *body-shot* lead to higher average desirability of the borrowers. On the other hand, due to the lower initial prominence of these features amongst borrowers with *type* features associated with lower outcomes, we observe a decrease in the Gini coefficient. Finally, the *Hybrid* policy combines the effects of *Partial Compliance* and *Low-type support*. There is a substantial reduction in the Gini coefficient of the distribution outcomes and a moderate gain in efficiency.

In Figure 14 we compare the proposed policies using the alternative fairness metric: the sum of market shares of the 33% of the least popular borrowers in the pool. The vertical axis - the share of borrowers choosing an inside option - is unchanged. First, in the *Baseline* policy, only 1% of lenders choose a borrower in the bottom third. Second, all the proposed policies increase the share of the borrowers with the lowest outcomes, except for *Naive* and *Restrict Competition*, where we see a reduction in the market share of the bottom third. The *Partial Compliance* policy improves on both dimensions. Finally, the *Hybrid* policy leads to a pronounced increase in the share of the bottom third, up to 2%.

To better visualize the impact of proposed policies, we return to the histograms presented in Section 4.2. In Figure 15, we present histograms of observed and simulated *cash per day*. In red (baseline) we show the distribution in *Kiva data*. In blue we present simulated outcomes.

Figure 14: Fairness - Efficiency tradeoff: bottom third market share



Note: Bottom third market share and efficiency. Each point represents the mean of 100 simulations with 500 lenders each. The horizontal axis presents the sum of the market shares of borrowers in the bottom third by outcomes. The vertical axis reports the share of lenders choosing an outside option.

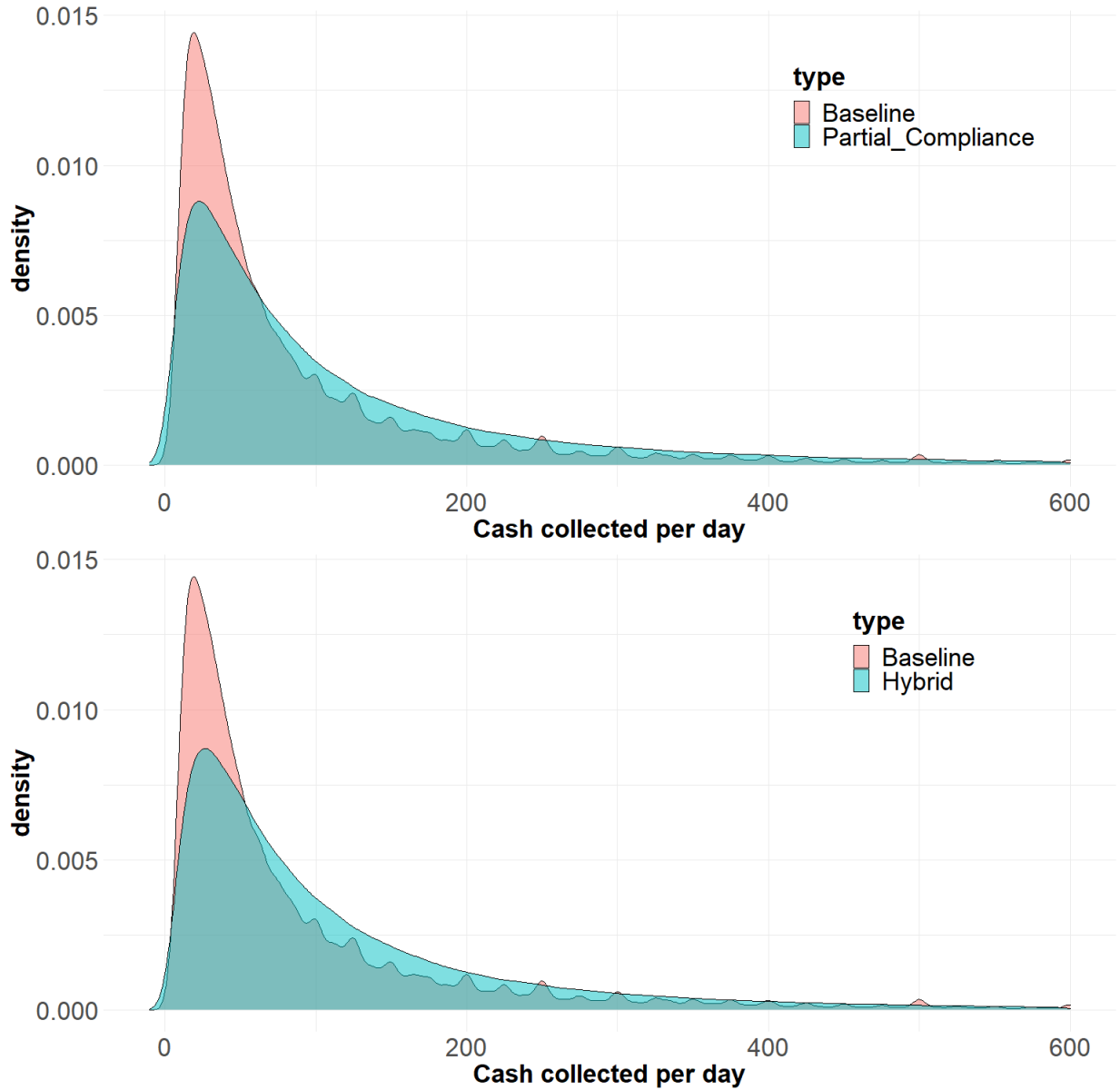
To obtain simulated outcomes, we randomly draw from a log-normal distribution such that the mean value equals the average *cash per day* in *Kiva data* and the variation of the distribution is selected such that the difference in Gini coefficients between baseline and counterfactual is the same, in percentage terms, as in simulations presented in Figure 13. Next, to adjust for the change in efficiency, we increase all values by the percentage difference in efficiency from Figure 13. Both policies move some borrowers from the low to moderate outcomes sections of the distribution. However, some high-performing borrowers experience lower outcomes.

6.5 The impact of proposed policies on gender inequity

In Section 4.5, we documented a substantial disparity in outcomes across *genders*.²⁷ In Figure 16, we present the impact of proposed policies on *gender equity* and efficiency. In the *Baseline* policy only about 25% of lenders choose a borrower with a *male* profile (we adjust the share of selected *male* borrowers by the share of *male* borrowers in the pool). Part of the gender gap can be attributed to different frequencies of the selected *style* features. From *Benchmark* we can observe that when all borrowers have profiles with *smile* and without *body-shot*, the share of lenders choosing a borrower

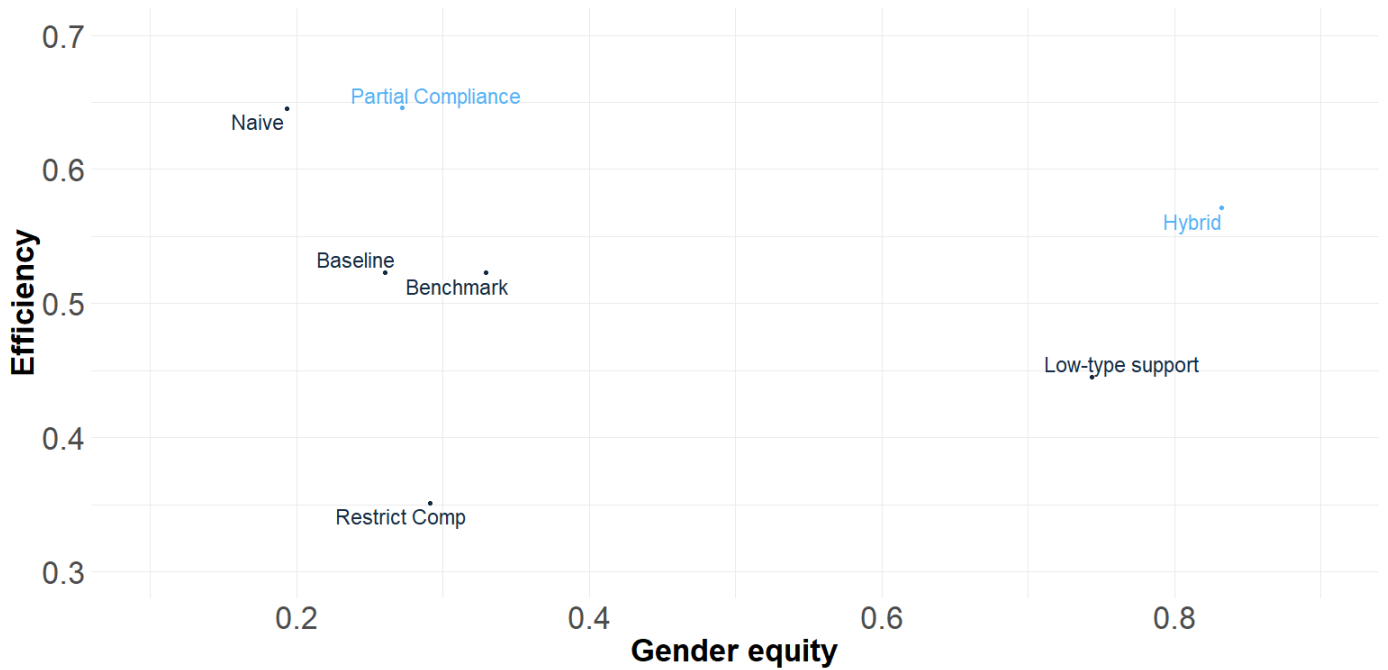
²⁷Recall that we use the algorithmic prediction of gender; thus, the inequity that we consider relates to the disparity between profiles that the algorithm classifies as *male* and *female*.

Figure 15: Histograms of distribution of cash per day under baseline and selected counterfactual policies



Note: Histograms of cash collected per day. The baseline in red is the distribution observed in Kiva data; the blue histogram is based on data simulated from counterfactual policies. Partial Compliance on the top panel and Hybrid on the bottom.) Values for the counterfactual policies simulated from log-normal distribution to match the values in Figure 13.

Figure 16: Fairness - Efficiency tradeoff: gender disparity



Note: Gender disparity and efficiency. Each point represents the mean of 50 simulations with 1000 lenders each. The horizontal axis presents the share of lenders choosing a borrower with a male profile, adjusted for the share of male profiles in the borrower pool. The vertical axis reports the share of lenders choosing an outside option.

with a *male* profile increases to approximately 30%.

All of our counterfactual policies, with the exception of *Restrict Competition*, increase the share of lenders choosing a *male* borrower. However, we observe that *Low-type support* and *Hybrid* policies result in the disparity going in the opposite direction; under those policies, campaigns of *male* type substantially outperform *female* campaigns. The *Hybrid* policy boosts the share of lenders choosing a borrower with a *male* profile to above 84% and leads to a moderate gain in efficiency. Finally, the *Partial compliance* policy increases both fairness and efficiency.

The results of the counterfactual simulations confirm the logic described earlier: when desirable *style* features are positively correlated with *type* characteristics leading to better funding outcomes, platform policies that promote the selected *style* characteristics by increasing the prominence of profiles that have them aggravate inequities in outcomes. In contrast, policies that change the distribution of attractive *style* features in a way that increases their prevalence amongst borrowers with high *type* characteristics lead to a more equitable distribution of outcomes. Furthermore, policies that increase the overall share of profiles with desirable *style* also boost efficiency.

The proposed model is a simplification of interactions between lenders, borrowers, and the micro-

lending platform. Amongst other assumptions, we abstract from supply-side responses or the possibility of lenders choosing several campaigns. Therefore, the specific magnitudes of the impact that the policies we discuss have on fairness and efficiency should be measured in a randomized experiment. Nevertheless, the proposed approach provides strong support for testing policies based on *style* recommendations.

7 Conclusion

We study how to improve fairness and efficiency in an online marketplace in which users have preferences for certain features in profile images. We introduce a distinction between *type* and *style* features: *type* is fixed when users create their online profiles, and *style* features are determined during the profile creation.

Using observational data from a large microfinance platform, we first show high inequities in funding outcomes. Second, we demonstrate that *style* features predict funding outcomes but not defaults, and third, we showcase selected *style* features that have an impact on funding outcomes and are correlated with borrowers' *types*. We highlight that these *style* features exacerbate inequity in outcomes in a way that - we argue - is unfair to borrowers. To corroborate these findings, we carry out a recruited experiment. In the experiment, we use Generative Adversarial Networks to generate fabricated images with a variation in selected features. We document that subjects prefer profile images with *smile* that are not *body-shots*. Finally, we evaluate counterfactual platform policies exploiting the estimates of the impact of selected *style* features on outcomes and their correlation with *types*. We show that a policy that encourages profiles with desirable features increases fairness and leads to more transactions.

The mechanism underlying our findings is that unchangeable *types* and amenable *styles* can both lead to inequities, and their correlation determines which platform policies will promote fairness without sacrificing efficiency. In the case of a positive correlation between desirable *style* features and high *types*, platform policies which promote profiles with selected features are likely to lead to less fair outcomes. In contrast, policies that shift distributions of these features by increasing their adoption by low-performing users can promote fairness.

Our approach can be used to determine which types of policies have the potential to increase fairness and efficiency. However, the extent to which the proposed policies are effective depends on specific parameters of lenders' demand (magnitudes and their stability) and borrowers' respon-

siveness to recommendations. Therefore, we believe that the pipeline we propose can be useful to prioritize policies for a randomized experiment. Carrying out such an experiment is a natural next step in this research agenda.

Throughout the paper, we argued that inequities in outcomes which are created by *style* choices and are not justified by different repayment probabilities are unfair. However, it's possible that *style* features reveal information about the developmental impact of the loan. Unfortunately, we do not have metrics that would allow us to adjust for this type of consideration; doing so would be a valuable extension of this research.

References

- Abbey, J. D. and Meloy, M. G. (2017). Attention by design: Using attention checks to detect inattentive respondents and improve data quality. *Journal of Operations Management*, 53:63–70.
- Aggarwal, R., Goodell, J. W., and Selleck, L. J. (2015). Lending to women in microfinance: Role of social trust. *International Business Review*, 24(1):55–65.
- Ai, W., Chen, R., Chen, Y., Mei, Q., and Phillips, W. (2016). Recommending teams promotes prosocial lending in online microfinance. *Proceedings of the National Academy of Sciences*, 113(52):14944–14948.
- Alesina, A. F., Lotti, F., and Mistrulli, P. E. (2013). Do women pay more for credit? Evidence from Italy. *Journal of the European Economic Association*, 11:45–66.
- Ash, E., Durante, R., Grebenshchikova, M., and Schwarz, C. (2022). Visual representation and stereotypes in news media.
- Athey, S., Imbens, G. W., Metzger, J., and Munro, E. (2021). Using Wasserstein Generative Adversarial Networks for the Design of Monte Carlo Simulations. *Journal of Econometrics*.
- Athey, S., Tibshirani, J., and Wager, S. (2019). Generalized random forests. *The Annals of Statistics*, 47(2):1148–1178.
- Banerjee, A., Karlan, D., and Zinman, J. (2015). Six randomized evaluations of microcredit: Introduction and further steps. *American Economic Journal: Applied Economics*, 7(1):1–21.
- Brock, J. M. and De Haas, R. (2021). Discriminatory lending: Evidence from bankers in the lab. *CentER Discussion Paper*.
- Duarte, J., Siegel, S., and Young, L. (2012). Trust and credit: The role of appearance in peer-to-peer lending. *The Review of Financial Studies*, 25(8):2455–2484.
- D’Espallier, B., Guérin, I., and Mersland, R. (2011). Women and repayment in microfinance: A global analysis. *World Development*, 39(5):758–772.
- Edelman, B., Luca, M., and Svirsky, D. (2017). Racial discrimination in the sharing economy: Evidence from a field experiment. *American Economic Journal: Applied Economics*, 9(2):1–22.

- Edelman, B. G. and Luca, M. (2014). Digital discrimination: The case of Airbnb. com. *Harvard Business School NOM Unit Working Paper*, (14-054).
- Ert, E., Fleischer, A., and Magen, N. (2016). Trust and reputation in the sharing economy: The role of personal photos in Airbnb. *Tourism management*, 55:62–73.
- Flores-Macías, G. and Zarkin, J. (2022). Militarization and perceptions of law enforcement in the developing world: Evidence from a conjoint experiment in Mexico. *British Journal of Political Science*, 52(3):1377–1397.
- Fong, C. M. and Luttmer, E. F. (2009). What determines giving to Hurricane Katrina victims? Experimental evidence on racial group loyalty. *American Economic Journal: Applied Economics*, 1(2):64–87.
- Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of Statistics*, pages 1189–1232.
- Galak, J., Small, D., and Stephen, A. T. (2011). Microfinance decision making: A field study of prosocial lending. *Journal of Marketing Research*, 48(SPL):S130–S137.
- Ge, Y., Knittel, C. R., MacKenzie, D., and Zoepf, S. (2016). Racial and gender discrimination in transportation network companies. Technical report, National Bureau of Economic Research.
- Gelbach, J. B. (2016). When do covariates matter? And which ones, and how much? *Journal of Labor Economics*, 34(2):509–543.
- Glynn, A. N. and Quinn, K. M. (2010). An introduction to the augmented inverse propensity weighted estimator. *Political Analysis*, 18(1):36–56.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative Adversarial Networks. *arXiv preprint arXiv:1406.2661*.
- Hainmueller, J., Hopkins, D. J., and Yamamoto, T. (2014). Causal inference in conjoint analysis: Understanding multidimensional choices via stated preference experiments. *Political Analysis*, 22(1):1–30.
- Jenq, C., Pan, J., and Theseira, W. (2015). Beauty, weight, and skin color in charitable giving. *Journal of Economic Behavior & Organization*, 119:234–253.
- Johannemann, J., Hadad, V., Athey, S., and Wager, S. (2019). Sufficient representations for categorical variables. *arXiv preprint arXiv:1908.09874*.

- Jordan, S. R., Rudeen, S., Hu, D., Diotalevi, J. L., Brown, F. I., Miskovic, P., Yang, H., Colonna, M., and Draper, D. (2019). The difference a smile makes: Effective use of imagery by children's nonprofit organizations. *Journal of Nonprofit & Public Sector Marketing*, 31(3):227–248.
- Karlan, D. and Morduch, J. (2009). Access to Finance. In Rodrick, D. and Rosenzweig, M. R., editors, *Handbook of Development Economics*, volume 5. Elsevier.
- Karras, T., Laine, S., and Aila, T. (2019). A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410.
- Kasy, M. and Abebe, R. (2021). Fairness, equality, and power in algorithmic decision-making. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 576–586.
- Kline, P. M., Rose, E. K., and Walters, C. R. (2021). Systemic discrimination among large US employers. Technical report, National Bureau of Economic Research.
- Kristof, N. D. and WuDunn, S. (2010). *Half the sky: Turning oppression into opportunity for women worldwide*. Vintage.
- Krumhuber, E., Manstead, A. S., Cosker, D., Marshall, D., Rosin, P. L., and Kappas, A. (2007). Facial dynamics as indicators of trustworthiness and cooperative behavior. *Emotion*, 7(4):730.
- Kung, F. Y., Kwok, N., and Brown, D. J. (2018). Are attention check questions a threat to scale validity? *Applied Psychology*, 67(2):264–283.
- Landry, C. E., Lange, A., List, J. A., Price, M. K., and Rupp, N. G. (2006). Toward an Understanding of the Economics of Charity: Evidence from a Field Experiment. *The Quarterly Journal of Economics*, 121(2):747–782.
- Lepri, B., Oliver, N., Letouzé, E., Pentland, A., and Vinck, P. (2018). Fair, transparent, and accountable algorithmic decision-making processes. *Philosophy & Technology*, 31(4):611–627.
- Ludwig, J. and Mullainathan, S. (2022). Algorithmic Behavioral Science: Machine Learning as a Tool for Scientific Discovery. *Chicago Booth Research Paper*, (22-15).
- Mendes, W. B. and Koslov, K. (2013). Brittle smiles: positive biases toward stigmatized and outgroup targets. *Journal of Experimental Psychology: General*, 142(3):923.

- Mullainathan, S., Noeth, M., and Schoar, A. (2012). The market for financial advice: An audit study. Technical report, National Bureau of Economic Research.
- Netzer, O., Lemaire, A., and Herzenstein, M. (2019). When words sweat: Identifying signals for loan default in the text of loan applications. *Journal of Marketing Research*, 56(6):960–980.
- Nosek, B. A., Smyth, F. L., Hansen, J. J., Devos, T., Lindner, N. M., Ranganath, K. A., Smith, C. T., Olson, K. R., Chugh, D., Greenwald, A. G., et al. (2007). Pervasiveness and correlates of implicit attitudes and stereotypes. *European Review of Social Psychology*, 18(1):36–88.
- Park, J., Kim, K., and Hong, Y.-Y. (2019). Beauty, gender, and online charitable giving. *Available at SSRN 3405823*.
- Pham, C. and Septianto, F. (2019). A smile—the key to everybody’s heart? the interactive effects of image and message in increasing charitable behavior. *European Journal of Marketing*.
- Pope, D. G. and Sydnor, J. R. (2011). What’s in a picture? Evidence of Discrimination from Prosper.com. *Journal of Human resources*, 46(1):53–92.
- Ravina, E. (2019). Love & loans: The effect of beauty and personal characteristics in credit markets. *Available at SSRN 1107307*.
- Rhue, L. and Clark, J. (2020). Automatically Signaling Quality? A Study of the Fairness-Economic Tradeoffs in Reducing Bias through AI/ML on Digital Platforms. *Working Paper, NYU Stern School of Business*.
- Robins, J. M., Rotnitzky, A., and Zhao, L. P. (1994). Estimation of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association*, 89(427):846–866.
- Salminen, J., Şengün, S., Santos, J. M., Jung, S.-G., and Jansen, B. (2022). Can Unhappy Pictures Enhance the Effect of Personas? A User Experiment. *ACM Transactions on Computer-Human Interaction*, 29(2):1–59.
- Septianto, F. and Paramita, W. (2021). Sad but smiling? how the combination of happy victim images and sad message appeals increase prosocial behavior. *Marketing Letters*, 32(1):91–110.
- Stigler, M. (2018). `dec_covar`: R implementation of Gelbach covariate decomposition. https://github.com/MatthieuStigler/Misconometrics/tree/master/Gelbach_decompo.

- Sun, L., Kraut, R. E., and Yang, D. (2019). Multi-level modeling of social roles in online micro-lending platforms. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–25.
- Theseira, W. (2009). *Competition to default: Racial discrimination in the market for online peer-to-peer lending*. PhD thesis, Dissertation, Wharton.
- Troncoso, I. and Luo, L. (2022). Look the part? the role of profile pictures in online labor markets. *Marketing Science*.
- Wager, S. and Athey, S. (2018). Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242.
- Williams, B. A., Brooks, C. F., and Shmargad, Y. (2018). How algorithms discriminate based on data they lack: Challenges, solutions, and policy implications. *Journal of Information Policy*, 8(1):78–115.
- Younkin, P. and Kuppuswamy, V. (2018). The colorblind crowd? Founder race and performance in crowdfunding. *Management Science*, 64(7):3269–3287.
- Zhang, S., Mehta, N., Singh, P. V., and Srinivasan, K. (2021). Can an AI algorithm mitigate racial economic inequality? An analysis in the context of Airbnb. *Working Paper*.
- Zhang, S. and Yang, Y. (2021). The unintended consequences of raising awareness: Knowing about the existence of algorithmic racial bias widens racial inequality. *Available at SSRN*.

Appendix

A Feature Detection Algorithms

Mask-RCNN. To structurally obtain features of images we use Mask-RCNN. The Mask-RCNN algorithm, developed by Facebook, detects objects from images. As shown in Figure 17, It takes in an image in the input layer and returns an estimated "package" for each object, including the class name, the bounding box, and the mask of each object detected, and those predictions are jointly optimized through the loss function.

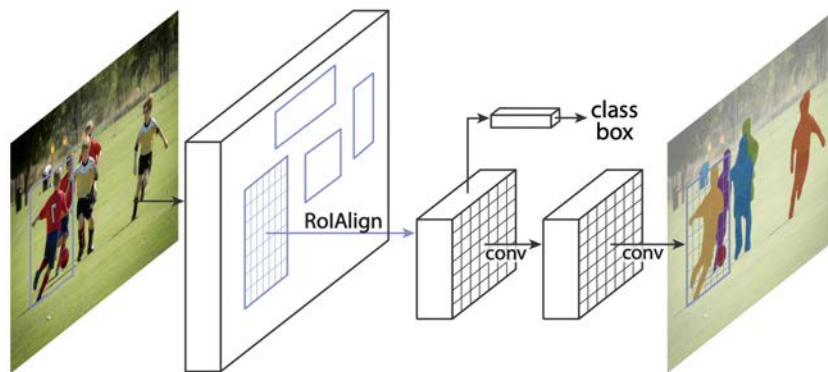


Figure 17: The Mask R-CNN framework: <https://arxiv.org/pdf/1703.06870.pdf>

Object detection. We apply this pre-trained model and estimate a score for each object which varies from 0 to 1. The score represents the algorithm's confidence in the existence of a specific feature, such as a tree, person, animal, digital items, etc. Figure 18 visualizes the output. We also apply this algorithm to detect human body-shot.²⁸

Facial feature classification. We detect facial features using the *face-classification* algorithm that takes in one face image and outputs a face embedding vector, evaluated by a pre-trained neural network.²⁹ Then, the embedding vector, as well as the feature labels, enter another neural network model (Multi-layer Perceptron). This model takes in one facial embedding vector and assigns a score for each unique facial feature such as *race*, *gender*, *smile*, etc. It is a supervised learning process, and the training label is pre-annotated.

The features that we obtain from images can be informally classified into three categories: (i)

²⁸<https://github.com/facebookresearch/detectron2>

²⁹<https://github.com/wondonghyeon/face-classification>



Figure 18: An example outcome of image detection using Mask-RCNN. Each detected object was given a label, put on a mask, and given the corresponding probability score: <https://github.com/facebookresearch/detectron2>.

technical aspects of the image (e.g., *blurry, flash, harsh light*), (ii) personal characteristics (e.g., *straight hair, eyes open, pale skin*), (iii) objects in the image (e.g., *chair, clock*).

Image and personal characteristics (e.g., race, age, hair color, facial shape, eyes/nose characteristics) are detected by FaceNet model which was pre-trained and tested on the large dataset CelebA with over 200,000 facial images. The algorithm detects the person’s face and then identifies its features.

B Auditing the Feature Detection Algorithm

To test whether the features detected by our algorithm correspond to the human perception of the image, we carried out an audit study. The audit study proceeded in two steps; first, we recruited human raters and asked them to label a sample of Kiva images. Second, we compared these labels with the prediction of the feature detection algorithm. We focused on two features: smiling and gender.

We carry out three analyses; first, we analyze the overall correlation between how human raters annotate the image and the model’s prediction of that annotation. Second, we consider the correlation of the model’s prediction of whether the person in the image smiles or not with human labels, separately for images labeled by humans as men and women. Finally, we divide the model’s errors into false positives and false negatives. The false positives are instances when human raters indicated that the feature is not present, while the model prediction was that it is; false negatives are when the

model predicts that the feature is not present, even though, according to raters, it is. The two types of errors create different types of bias, false positives lead to overestimation of the impact of the feature on outcomes, while false negatives lead to underestimation.

To create human-made labels, we randomly drew 100 images from the Kiva dataset and organized them into ten protocols, so each protocol had ten images. After that, we recruited 30 subjects per protocol on Prolific. Subjects were asked about the gender of the person in the image and whether the person smiles or not. We carried out attention checks and asked about the level of confidence the person had in the response. To compute a label, we average subjects’ responses per image.

Table 6 presents the comparison of the mean label per feature in the sample with the mean predicted probability from the feature-detection algorithm.

Table 6: Comparison of mean scores from the audit and CNN output.

Feature	Mean prob. CNN	Mean score audit	Difference	CI low	CI high
Man	0.48	0.56	-0.08	-0.20	0.05
Smiling	0.45	0.45	0.00	-0.09	0.09
Smiling amongst man	0.30	0.31	-0.01	-0.12	0.11
Smiling amongst woman	0.58	0.58	0.00	-0.13	0.13

Note: Average CNN is the average of predicted probabilities per image using the feature-detection algorithm, and mean score audit is the average label by human raters. Rows one and two show the values in the entire sample. Rows three and four present the values when considering the subsamples based on gender determined by human subjects. For example, the third row shows the values amongst the images labeled as a man by at least 6 out of 10 human subjects.

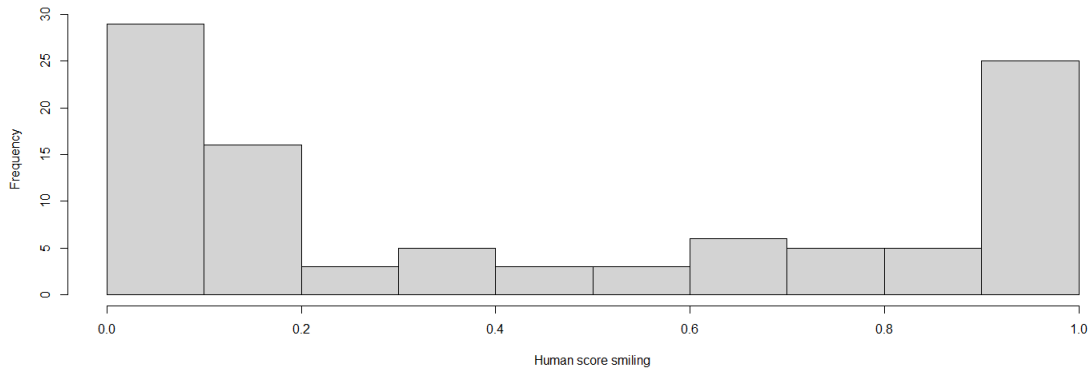
Table 6 shows that human raters and the algorithm detect similar frequencies of the selected features. Additionally, as shown in rows three and four, the frequency of smiling across images labeled by humans as man or woman is consistent across the two methods. Thus, the two methods lead to the same conclusion that in the Kiva images, men are less likely to smile than women.³⁰ Additionally, the Pearson correlation coefficient between the human label and the algorithm’s prediction is 0.92 for gender and 0.71 for smiling.

Finally, we transform the labels and the algorithm’s predictions into binary indicators, taking the value of 1 when the continuous variable is above 0.5 and 0 otherwise. We analyze how often the two indicators coincide. We find that human scores and the algorithm prediction lead to the same gender classification in 93% of cases and for smiling in 81% of images.

In the case of many images, human raters disagree about whether the person in the image smiles

³⁰We reach the same conclusions when considering only subjects that have passed all attention checks and when excluding images that less than 7 out of 10 subjects rated similarly.

Figure 19: Histogram of human scores for smiling



Note: Score is the average per image label by human raters.

or not. Figure 19 documents this. Many of the prediction errors are concentrated in these intermediate cases. When considering the subsample of images that seven or more human raters labeled similarly, the error rate declines to 14%.

We also consider the error rates of smile detection separately across genders as defined by human raters. The error rate amongst images labeled as man is 81.2% and woman 80.8%. We conclude that the algorithm performance is similar across genders. Finally, we group errors into false positives (the algorithm detects a feature that does not exist according to human raters) and false negatives (the algorithm does not detect an existing feature). We find that the rate of false positives is 7% and of false negatives is 34%.

To conclude, we find that the algorithm's predictions highly indicate the feature's presence in an image as perceived by human raters. The main concern is the high rate of false negatives for the smile detection task. We find that 34% of images labeled as smiling by humans are predicted as not smiling by the algorithm. Thus, we might be underestimating the impact of smiling on outcomes in the observational data. To assess the extent of the bias caused by these false negatives, we consider the following example: (i) there is a population in which half of the individuals have images with a *smile* feature, (ii) the outcome is a random variable normally distributed with a mean of 1 and a standard deviation of 1 for individuals that do not have an image with *smile*, and with a mean of 1.3 and a standard deviation of 1, for those that do *smile* in an image, (iii) we assume that the algorithm which detects features has the false negative error rate of 34%, (iv) we estimate a linear probability model using the OLS estimator. We simulate this example 1000 times and estimate the degree of bias caused

by the false negatives. We find that the false negative error rate of 34% results in the negative bias of the smiling coefficient of 22% (s.e. 0.16%); the average coefficient, across simulations, is 0.23 instead of 0.3. Consequently, the estimates used in the observational data understate the impact of *smile* on outcomes.

C Generative Adversarial Networks

The algorithm described above detects our features of interest. To modify images with respect to these features we use another tool known as Generative Adversarial Networks (GANs). GANs designed by Goodfellow et al. (2014) are an approach to generative modeling using deep learning methods. The key objective of GANs is to generate fabricated data that are similar to particular data, such as realistic images (Ludwig and Mullainathan, 2022) and synthetic datasets (Athey et al., 2021). GANs, although do not directly produce estimates of the density or distribution function at a particular point, can be thought of as implicitly estimating the distribution of latent features, and they can be used to generate or output new examples that plausibly could have been drawn from the original dataset.

The core idea of GAN is to have two models: a generator G and a discriminator D . As illustrated by Goodfellow et al. (2014), to learn the generator’s (image) distribution p_g over data x , we define a prior on input noise variables $p_z(z)$, then represent a mapping to data space as $G(z; \theta_g)$. Discriminator $D(x; \theta_d)$ outputs a single scalar, representing the probability that x came from the data rather than p_g . D and G play the two-player minimax with the value function:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [1 - \log D(G(z))]$$

GANs are frequently used to modify images and generate so-called "deep fakes" - fabricated images based on input images that have been altered in a specific way. In our work, we apply Style-GAN developed by Karras et al. (2019) to generate fake images that differ in a specific feature. Conditioning on the attributes (areas) we make the change, the algorithm detects the key image area to leave its counterpart unchanged.

The key image is fed into a pre-trained GAN generator and embedded, as a latent vector V , into a latent space. We compute the direction of the gradient ∇V of our feature of interest W (e.g. smiling), determined from our ATE analysis, by computing the “difference in means” of the latent vector encoded into the latent space from images with and without such feature³¹.

³¹We used around 200 images labeled by CNN and verified by human audit

$$\nabla V = \mathbb{E}[V_i[W_i = 1, X_i = x]] - \mathbb{E}[V_i[W_i = 0, X_i = x]]$$

We have hyper-parameters to decide the extent to which we want to alter the images in the desired direction. We fine-tune the hyper-parameters image by image to offset the correlation in image features bleeding into the pre-trained GAN model. The modified attribute is embedded into its unchanged counterpart, and we ensure that images look realistic by deblurring, inpainting, and auto-blending.

D Summary statistics of *Kiva data*

Table 7: Summary statistics of *Kiva data*

Statistic	N	Mean	St. Dev.	Min	Pctl(25)	Pctl(75)	Max
Loan amount	420,765	800.107	993.370	25	275	950	50,000
Cash per day	420,765	123.522	270.186	1	25	116.7	8,750
Days to raise	420,765	13.427	11.667	1	5	20	83
Total number of lenders	420,765	0.012	0.015	0.001	0.005	0.015	0.967
default	420,765	0.050	0.218	0	0	0	1
Male	420,765	0.198	0.398	0	0	0	1
Number of borrowers	420,765	1.958	3.171	1	1	1	50
No. competitors	420,765	0.091	0.173	0.003	0.006	0.075	1.000
Same race gender share	420,765	0.665	0.294	0	0.4	1	1
Asian	420,765	0.191	0.261	0.0001	0.016	0.266	0.995
White	420,765	0.218	0.265	0.001	0.031	0.323	0.999
Black	420,765	0.167	0.281	0.0001	0.006	0.148	0.990
Baby	420,765	0.004	0.003	0.0001	0.002	0.006	0.067
Child	420,765	0.073	0.056	0.001	0.034	0.095	0.609
Youth	420,765	0.264	0.211	0.0002	0.092	0.391	0.982
Middle.Aged	420,765	0.084	0.093	0.0004	0.026	0.104	0.898
Senior	420,765	0.041	0.079	0.0001	0.004	0.039	0.950
Black.Hair	420,765	0.388	0.242	0.0005	0.171	0.589	0.970
Blond.Hair	420,765	0.007	0.029	0.00000	0.001	0.004	0.943
Brown.Hair	420,765	0.405	0.156	0.012	0.288	0.517	0.919
Bald	420,765	0.037	0.073	0.0001	0.004	0.030	0.835
No.Eyewear	420,765	0.865	0.148	0.007	0.830	0.959	1.000
Sunglasses	420,765	0.017	0.014	0.001	0.008	0.020	0.327
Mustache	420,765	0.072	0.160	0.00003	0.004	0.046	0.998
Smiling	420,765	0.549	0.177	0.013	0.424	0.685	0.966
Chubby	420,765	0.339	0.190	0.012	0.185	0.466	0.972
Blurry	420,765	0.162	0.095	0.006	0.090	0.214	0.758
Harsh.Lighting	420,765	0.339	0.165	0.031	0.217	0.430	0.930
Flash	420,765	0.245	0.126	0.010	0.148	0.322	0.855
Soft.Lighting	420,765	0.677	0.090	0.222	0.623	0.742	0.943
Outdoor	420,765	0.447	0.140	0.045	0.343	0.545	0.914
Curly.Hair	420,765	0.394	0.155	0.031	0.275	0.499	0.932
Wavy.Hair	420,765	0.226	0.170	0.004	0.095	0.312	0.991
Straight.Hair	420,765	0.606	0.178	0.034	0.489	0.741	0.982
Receding.Hairline	420,765	0.205	0.235	0.0004	0.039	0.282	0.995
Bangs	420,765	0.171	0.171	0.001	0.052	0.229	0.993
Sideburns	420,765	0.145	0.195	0.001	0.025	0.168	0.977
Partially.Visible.Forehead	420,765	0.094	0.090	0.001	0.032	0.125	0.834
Arched.Eyebrows	420,765	0.451	0.213	0.004	0.282	0.618	0.978
Narrow.Eyes	420,765	0.588	0.204	0.031	0.432	0.755	0.992
Eyes.Open	420,765	0.871	0.073	0.338	0.834	0.925	0.991
Big.Nose	420,765	0.730	0.190	0.042	0.606	0.886	0.998
Big.Lips	420,765	0.586	0.215	0.014	0.425	0.766	0.986
Mouth.Closed	420,765	0.303	0.146	0.018	0.193	0.390	0.944
Mouth.Wide.Open	420,765	0.057	0.040	0.002	0.030	0.072	0.516
Square.Face	420,765	0.019	0.041	0.00005	0.002	0.015	0.759
Round.Face	420,765	0.201	0.155	0.002	0.078	0.287	0.908
Color.Photo	420,765	0.948	0.026	0.632	0.935	0.966	0.997
Posed.Photo	420,765	0.486	0.132	0.069	0.391	0.581	0.925
Attractive.Woman	420,765	0.125	0.151	0.001	0.028	0.158	0.989
Indian	420,765	0.061	0.098	0.00002	0.009	0.066	0.962
Bags.Under.Eyes	420,765	0.586	0.170	0.016	0.468	0.717	0.967
Rosy.Cheeks	420,765	0.122	0.069	0.011	0.072	0.155	0.729
Shiny.Skin	420,765	0.215	0.121	0.004	0.121	0.288	0.808
Pale.Skin	420,765	0.334	0.171	0.014	0.192	0.460	0.908
Strong.Nose.Mouth.Lines	420,765	0.611	0.172	0.026	0.496	0.746	0.966
Flushed.Face	420,765	0.102	0.050	0.009	0.067	0.126	0.573
Top	420,765	157.544	106.715	0	80	204	1,598
Right	420,765	410.062	174.165	29	271	534	960
Bottle	420,765	0.503	2.259	0	0	0	99
Chair	420,765	0.125	0.498	0	0	0	24
Person	420,765	2.119	3.002	1	1	2	39
Bodyshot	420,765	0.406	0.491	0	0	1	1

E Choice of the predictive model

In this section, we consider several predictive models over three specifications and determine the model to be used in the baseline analysis.

We analyze the performance of models predicting *cash per day*. We consider the following models: Linear Regression, LASSO, Random Forrest (grf), and Boosted Random Forrest (grf and gbm). All models (except for LM) are tuned for the task at hand, we report the performance of the selected best (lowest MSE) model. All models are trained using a 70% sample of *Kiva data* and tested on the 30%.

We consider three specifications differing by the number of covariates: (A) covariates include: details of the loan including amount, repayment scheme, *sector*, *country*, etc. and weekly dummies, (B) details of the photo including both *type* and *style* characteristics, (C) total number of active lenders in this *week*sector*, total number of competitors in this *week*sector*, number of competitors of the same *race* and *gender*, and interaction of *week* and *sector*, and interaction of *week* and *country*. For boosted Forrest we also add a 4th specification where we have a sufficient representation of *week* sector* (D) (Johannemann et al., 2019). Table 8 presents results.

Table 8: Comparison of the test-set predictive performance of selected models

Model	Specification	MSE	SE
Linear regression	A	13840	159
Linear regression	B	13466	155
Linear regression	C	13565	166
LASSO	A	13797	161
LASSO	B	13379	157
LASSO	C	13183	156
Random forest	A	13930	163
Random forest	B	13530	145
Random forest	C	13099	157
Boosted forest (gbm)	A	12235	156
Boosted forest (gbm)	B	11477	141
Boosted forest (gbm)	C	10929	157
Boosted forest (gbm)	D	11406	173
Boosted forest (grf)	A	12665	147
Boosted forest (grf)	B	12003	149
Boosted forest (grf)	C	11777	139
Boosted forest (grf)	D	11962	177

Note: Test set performance of selected predictive models with different sets of covariates.

We conclude that Boosted Forrest has the best test-set predictive performance across all specifica-

tions and we decide to use it as a baseline model for the predictive tasks throughout the paper. *GBM* implementation of the Boosted Forrest has better performance than *GRF*, the difference is moderately small. Sufficient representation does not improve models’ performance and will not be used in the predictive tasks.

F Analysis of defaults across default types

We observe two different reasons for the loan not being repaid: a default by a microfinance organization and a default by the borrower. It’s plausible that image features are predictive of a borrower’s default but not of the microfinance organization. In this section, we separately analyze the predictive performance of a model trained to predict defaults by the borrower with and without image features.

We train a Boosted Forrest (*GBM*) on 70% of data and report the predictive performance on the 30% test set. We consider two specifications a full model, model C in Ewith defaults as the dependent variable, and a model from which we remove image features. Table 9 reports results.

Table 9: Comparison of the test-set predictive performance of models of default with and without image features.

Outcome	Covariates	MSE	Std. error
All defaults	full model	0.059	0.00094
All defaults	no image covariates	0.059	0.00095
Defaults by the borrower	full model	0.046	0.00088
Defaults by the borrower	no image covariates	0.046	0.00088

Note: Test set performance of selected predictive models with different sets of covariates.

These results suggest that image characteristics do not improve the predictive performance of either of the default models.

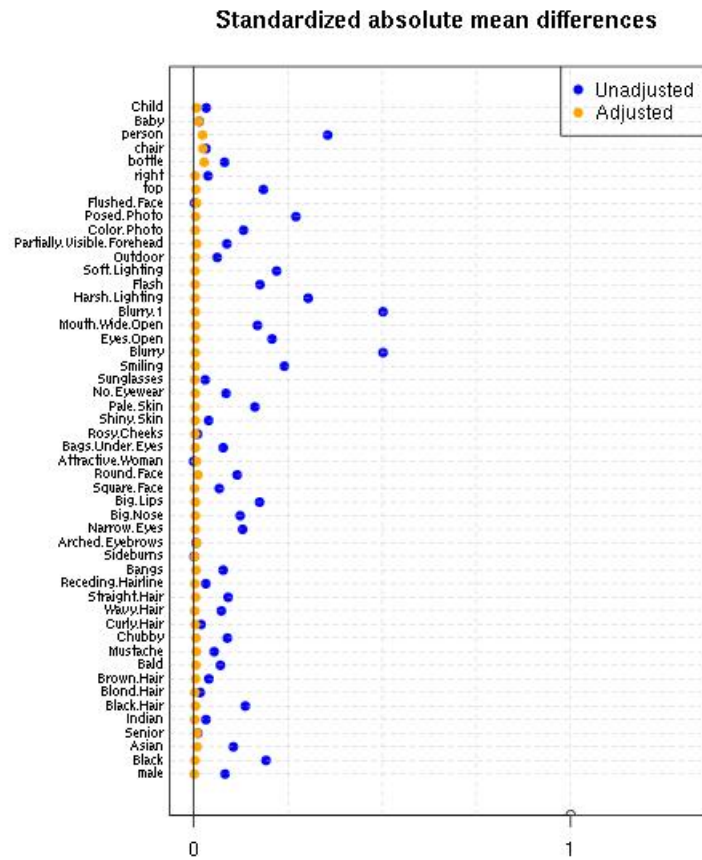
G Supplementary analysis for AIPW estimates of style features

G.1 Diagnostics for selected *style* features

Diagnostics *bodyshot*. Figure 20 shows standardized absolute mean differences of covariates across treatment group (with *Bodyshot*) and control. We see that the adjusted values (yellow) are well balanced.

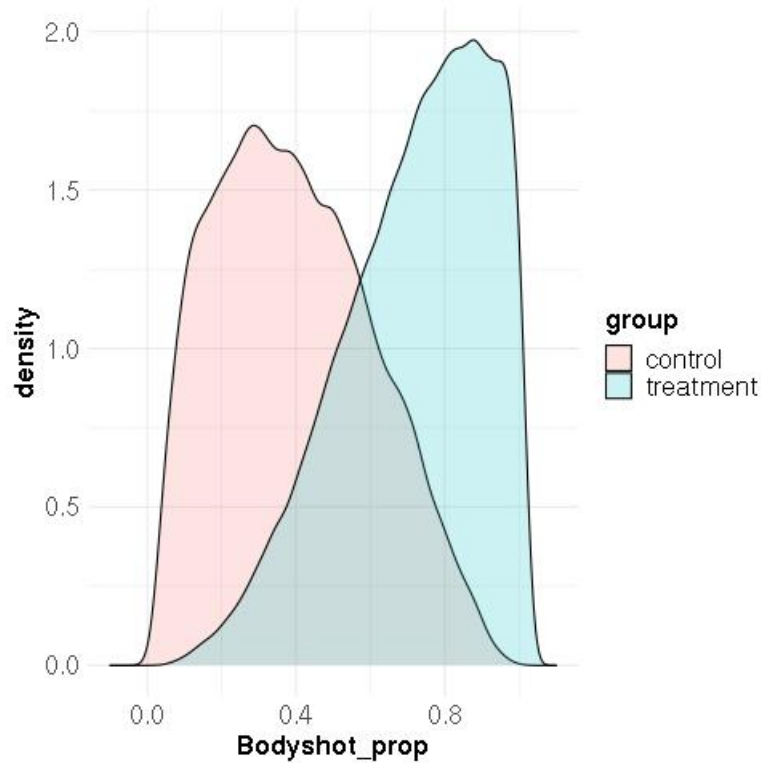
In Figure 21 we show the propensity scores to have a profile with *Bodyshot* across profile images with and without a *Bodyshot*. We find that there is common support between the two groups.

Figure 20: Diagnostics for *bodyshot*



Note: Standardized absolute mean differences of a selected subset of other covariates across profiles with and without bodyshot. Propensity score used for reweighing obtained using GBM model trained on all covariates in Kiva data.

Figure 21: Propensity scores for *Bodyshot*

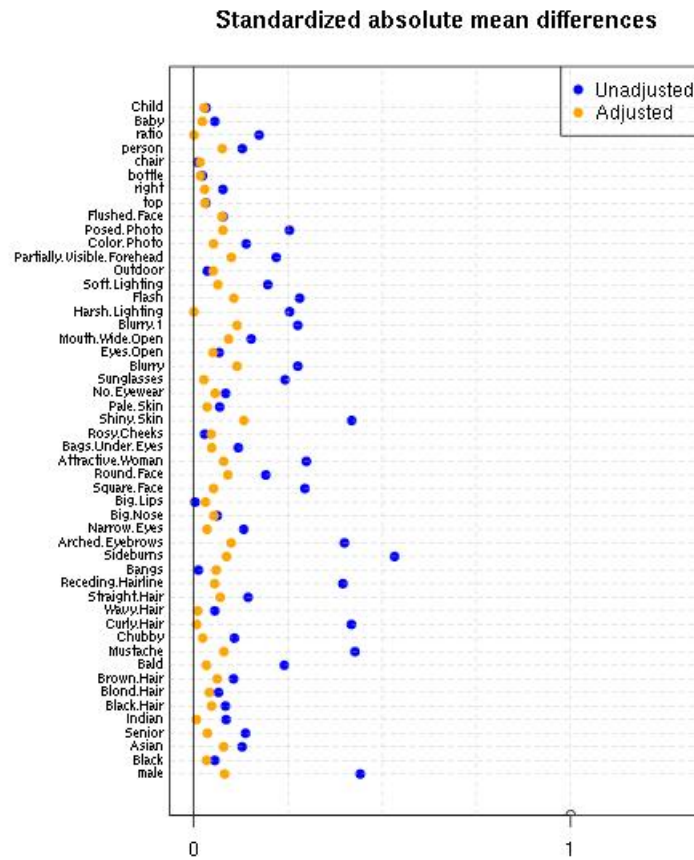


Note: Estimates of the propensity for an image to be a Bodyshot using a GBM-based prediction with full set of controls from Kiva data.

Diagnostics *smile*. Figure 22 shows standardized absolute mean differences of covariates across the treatment group (with *smile*) and control. We see that the adjusted values (yellow) are well balanced.

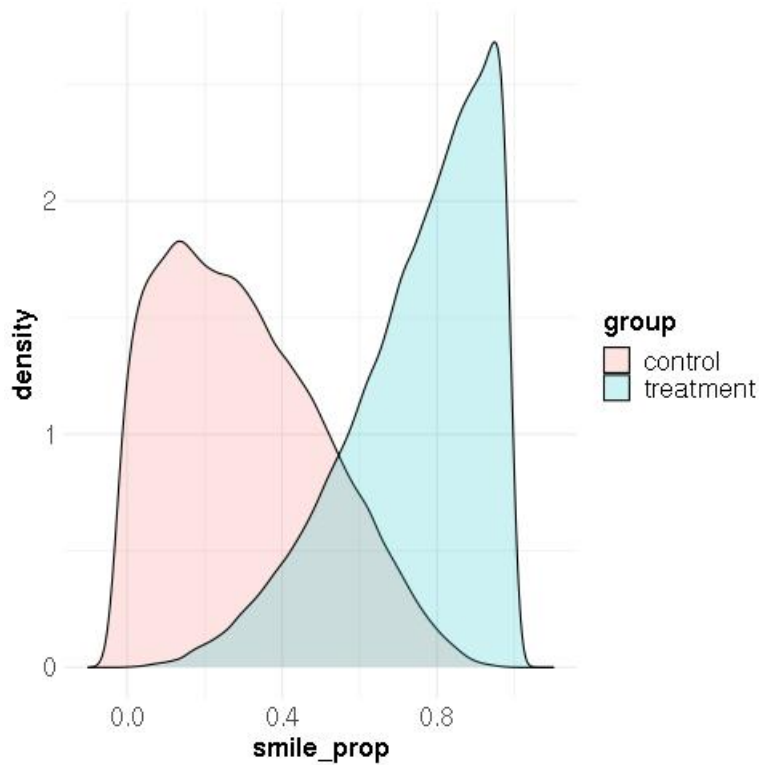
In Figure 23 we show the propensity scores to have a profile with *smile* across profile images with and without a *smile*. We find that there is common support between the two groups.

Figure 22: Diagnostics for *Bodyshot*



Note: Standardized absolute mean differences of a selected subset of other covariates across profiles with and without smile. Propensity score used for reweighing obtained using GBM model trained on all covariates in Kiva data.

Figure 23: Propensity scores for *smile*



Note: Estimates of the propensity for an image to be a smile using a GBM-based prediction with the full set of controls from Kiva data.

H Attention checks in the experiment

To check the quality of experimental data, we included attention checks in the survey. Attention checks are questions designed explicitly to detect inattentive subjects through additional questions (Abbey and Meloy (2017)). There are three purposes of the attention checks in our experimental setting: first, attention checks provide a signal of whether a recruited subject is paying attention to the information on the screen. Second, attention checks encourage the subjects to make thoughtful decisions. In addition, attention checks also give us the flexibility to filter the data in order to have high-quality ones, depending on whether we would like to tighten or loosen our criteria.

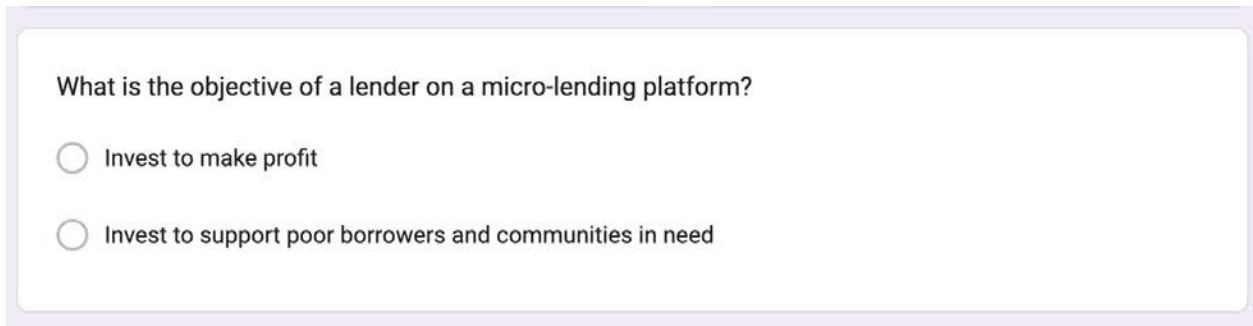
In order to avoid the attention checks themselves inducing a deliberative mindset and becoming a threat to the validity, we try to ask the subjects to recall detail in a previous image after they make the choice and the correct answer to that gives us the reason to believe that people have been paying rational attention to their choices.³²

The attention check in Figure 24 asks *What is the objective of a lender on a micro-lending platform?*

³²Kung et al. (2018) encourage researchers to justify the use of attention checks without compromising scale validity

This question asks about information provided on the first slide of each protocol.

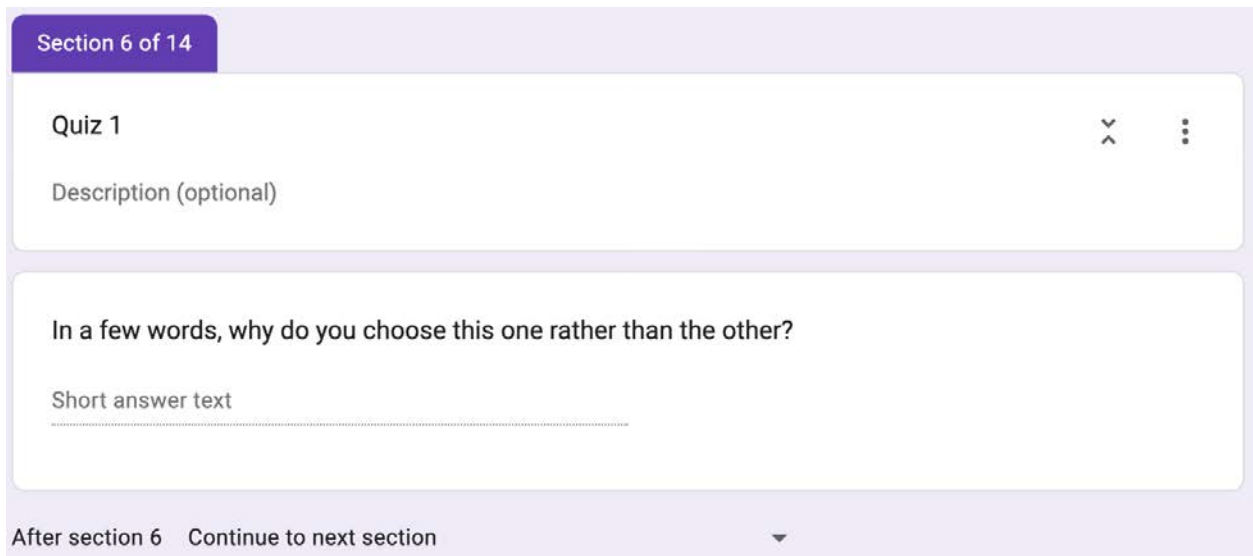
Figure 24: Attention check 1



What is the objective of a lender on a micro-lending platform?

- Invest to make profit
- Invest to support poor borrowers and communities in need

Figure 25: Attention check 2



Section 6 of 14

Quiz 1

Description (optional)

In a few words, why do you choose this one rather than the other?

Short answer text

After section 6 Continue to next section

Attention checks in Figures 25 and 26 are conducted in the format of a quiz. Attention check 2 is an open-ended query asking the subject for the reason for their decisions.³³ The last check is a multiple choice query asking about the occupation of the borrower on the previous slide.

Figures 27 and 28 show shares of subjects that responded correctly to Attention check 1 and 3. In both cases, correct response rates are above 90%. We take this as an indication that subjects were generally paying attention to their choices.

³³Abbey and Meloy (2017) uses this type of attention checks and manipulation validations to detect inattentive respondents in primary empirical data collection

Figure 26: Attention check 3

Section 8 of 14

Quiz 2 ✕ ⋮

Description (optional)

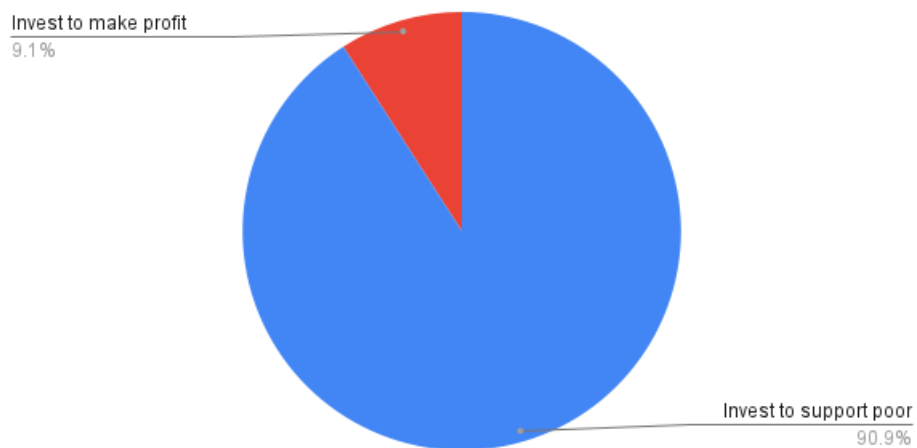
What is the occupation of the borrower on the left side?

Farmer

Teacher

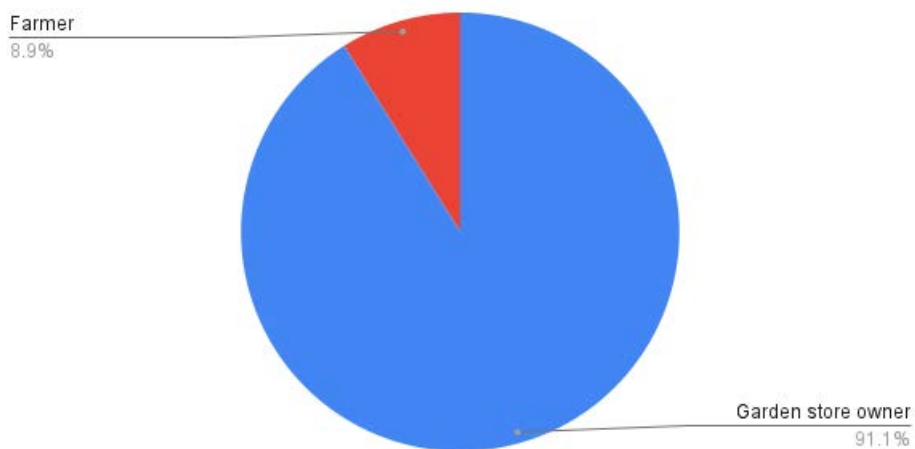
After section 8 Continue to next section ▼

Figure 27: The proportion of correct answers (blue) to the object of a lender



Note: Count of What is the objective of a lender on a micro-lending platform

Figure 28: The proportion of correct answers (blue) to the borrower's occupation is shown on the previous page

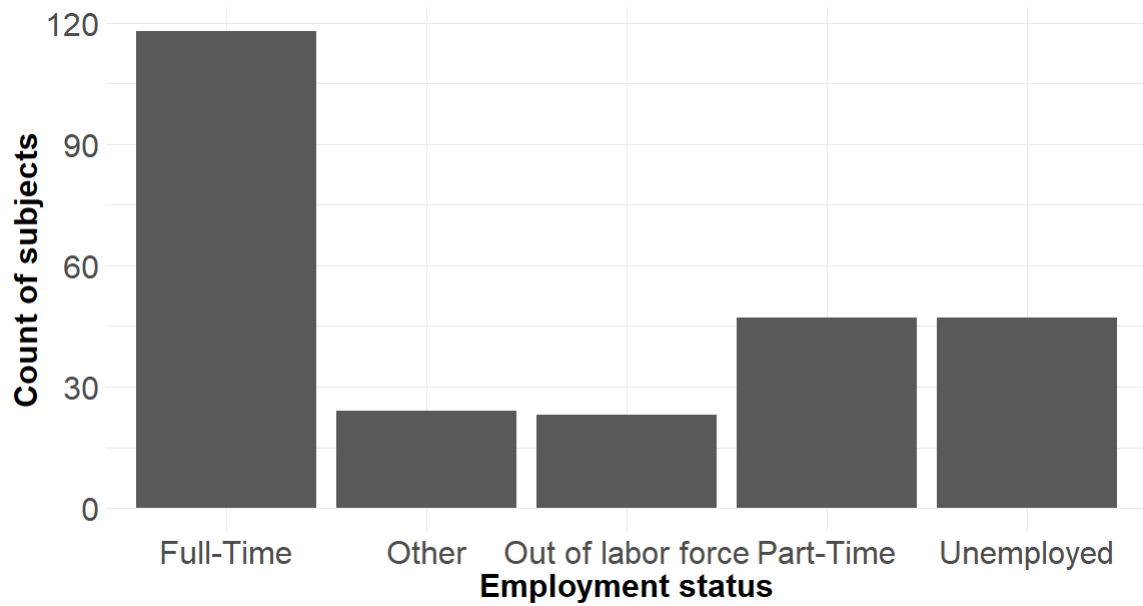


Note: Count of borrower's occupation shown in the previous page

I Summary statistics from the experiment

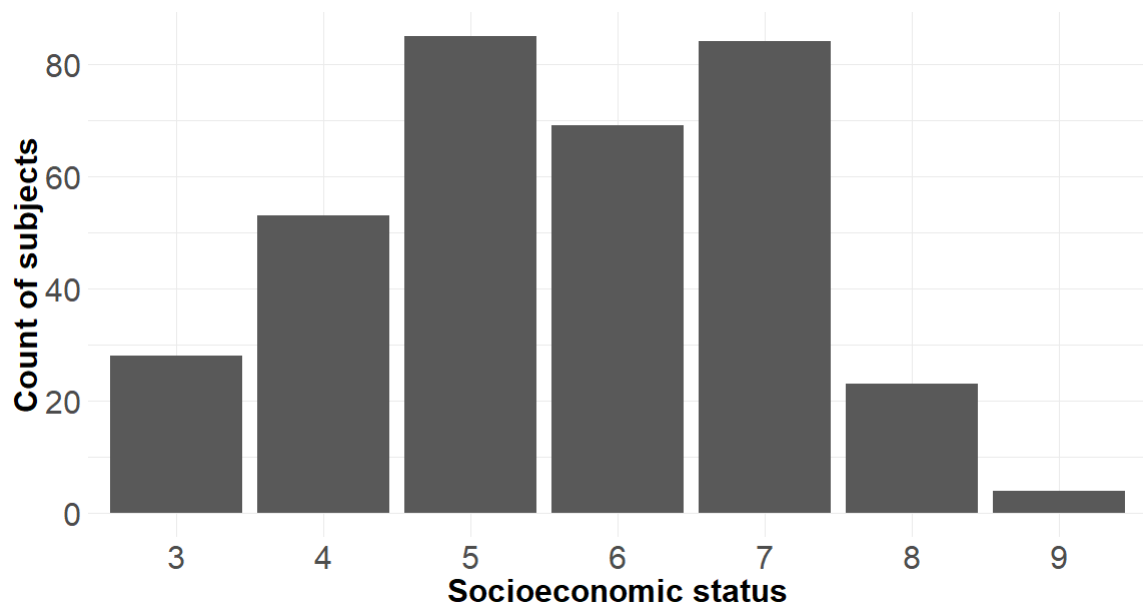
Figure 29 presents the employment status as reported by the subjects, there is a lot of missing data, to a large extent, this is due to the employment information being expired. Figure 30 shows the self-reported socioeconomic status, we screened participants to be of at least status 3.

Figure 29: Current employment status.



Note: Self-reported employment status. We drop observations where the data is unavailable (29%) and we group together full-time and 'starting new job soon' responses.

Figure 30: Self-assessed socio-economic status.



Note: Self-assessed socio-economic status. We required subjects to have at least a score of 3.

J Algorithm for counterfactual simulations

In this section, we describe the algorithm for generating outcomes under counterfactual policies in more detail. We divide the algorithm into two parts: (i) simulation of a market, and (ii) simulation of lenders' choices.

Algorithm 1 Simulation of a market

```

 $\tilde{\eta} \leftarrow U(\mathcal{N}; 22)$  ▷ Draw 22 fixed effects uniformly from the set of estimated fixed effects
 $\tilde{male} \leftarrow \mathbf{E}_G[male|D(\tilde{\eta}); 22]$  ▷ Draw 22 gender realizations
 $\tilde{bodyshot} \leftarrow \mathbf{E}_G[bodyshot|D(\tilde{\eta}), \tilde{male}; 22]$ 
 $\tilde{smile} \leftarrow \mathbf{E}_G[smile|D(\tilde{\eta}), \tilde{male}; 22]$ 
if  $H \in \{Partialcompliance\}$  then
  if  $\tilde{bodyshot} == 1$  then
     $\tilde{bodyshot} = B_{0.25}$  ▷ Bernoulli trial with  $p = 0.25$ 
  end if
  if  $\tilde{smile} == 0$  then
     $\tilde{smile} = B_{0.75}$ 
  end if
end if
 $x \leftarrow (\tilde{\eta}, \tilde{male}, \tilde{bodyshot}, \tilde{smile})$ 

if  $H \in \{RestrictCompetition\}$  then
   $\mathcal{M} \leftarrow h(x; 5)$ 
else
   $\mathcal{M} \leftarrow h(x; 11)$  ▷ Draw borrowers from the pool following the probability function  $h$ 
end if
 $\mathcal{M} \leftarrow (\mathcal{M}, \omega)$  ▷ add outside option
return  $\mathcal{M}$ 

```

Algorithm 1 proceeds in two steps, first, simulates the pool of borrowers and, second, samples from the pool to construct the market. Policies impact the distribution of the features in the pool (*partial compliance*), the size of the market (*Restrict competition*), and the probability of being sampled into the market (through the function h).

Once a market is simulated we determined lenders' choices with Algorithm 2. We first simulate the preferences of a lender, then compute the utility associates from different borrowers, and, finally, determined which borrower is selected.

Algorithm 2 Simulation of a lender choice

```

 $(\tilde{\alpha}, \tilde{\beta}, \tilde{\gamma}) \leftarrow (N(\alpha, sd_\alpha), N(\beta, sd_\beta), N(\gamma, sd_\gamma))$  ▷ draw preference parameters
 $\tilde{\epsilon} \leftarrow GEV$  ▷ draw random utility parameters for each borrowing campaign
 $u \leftarrow U(\mathcal{M}; \tilde{\alpha}, \tilde{\beta}, \tilde{\gamma}, \tilde{\epsilon})$  ▷ compute utilities from choosing any of the borrowers
 $choice \leftarrow \max(u)$ 
return  $choice$ 

```

K Fake and genuine smiles distinction

The effectiveness of policies that encourage borrowers to change facial expressions, specifically to smile, relies on the premise that the previously non-compliant borrowers can create images with desired features and that these newly added features impact lenders' choices; for example, the platform policy might be ineffective if lenders perceive the facial expressions in new images as not genuine. This section argues that this concern is legitimate by showing that non-genuine smiles do not increase funding rates.

To introduce a distinction between genuine and fake (forced) smiles, we train an algorithm that classifies the type of smile. We develop this algorithm using a dataset of 6442 images classified by human annotators as fake or genuine smiles.³⁴

We use the algorithm in a random sample of 45 thousand profiles from the Kiva observational dataset. First, we predict whether the person in the image smiles and whether the smile is genuine or fake. Next, we group the borrowers by the predicted type of smile and compute the average cash collected per day. Finally, we estimate the average impact of each type of smile on cash collected per day; to do that we use the AIPW estimator (we follow the same methodology as in Section 4). Table 10 shows the results.

Table 10: The impact of different types of smile on cash per day.

Estimand	Not smiling	Any smile	Genuine smile	Fake smile
Mean outcome in group	131.8 (0.7)	115.0 (0.7)	136.5 (0.8)	116.4 (0.7)
Average treatment effect	-	7.3 (1.0)	15.3 (1.2)	0.8 (1.2)

Note: The first row shows the mean cash per day across four groups of borrowers: not smiling, having any type of smile, having a genuine smile, and a fake smile. The second row shows the average effect of having a smile on cash collected per day. We estimate the effect using the AIPW estimator, which adjusts for all other observable characteristics. The comparison group includes borrowers that do not have images with a smile. Standard errors are in parentheses.

Results presented in Table 10 indicate that only smiles that our algorithm predicted to be genuine lead to higher outcomes. Specifically, we estimate that a genuine smile increases the cash collected per day by \$15, while a fake smile has no statistically significant impact.

³⁴The dataset and the original model structure are referred here: <https://github.com/vviveks/FakeSmileDetection>; we modified the original algorithm to the task of binary prediction - genuine or fake. Our algorithm predicts the fakeness of smiles using three different detected components of each face: whole face, eyes, and mouth. We train three deep neural networks (ResNet, DenseNet, and AlexNet) jointly and concatenate the learned latent vectors to make a joint prediction of whether the smiling is fake in the last layer. The cross-entropy loss of the prediction in the test set of 0.67 (0.66 in the train set), and the precision (f1-score) of 0.70 in the test set (also 0.71 in the train set).

This analysis showcases an important limitation of the policy based on facial expressions. If lenders perceive some of the smiles created in response to the new policy as not genuine, they might not increase funding rates. In the simulation exercise, we assumed that 75% of the previously non-compliant borrowers become compliant under the new policy. The policy becomes less effective when the share of borrowers that create images with genuine smiles decreases.

To mitigate the risk that a new policy is not effective, a platform might design a system that gives borrowers instant feedback on their images, helping them create profile photos that are impactful. An algorithm similar to the one developed in this section can be a part of such a policy.