

NBER WORKING PAPER SERIES

MACHINE LEARNING AND MOBILE PHONE DATA CAN IMPROVE THE TARGETING
OF HUMANITARIAN ASSISTANCE

Emily Aiken
Suzanne Bellue
Dean Karlan
Christopher R. Udry
Joshua Blumenstock

Working Paper 29070
<http://www.nber.org/papers/w29070>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
July 2021, Revised June 2022

Isabel Onate Falomir, Shikhar Mehra, Suraj Nair, Adrian Dar Serapio, Nathaniel Ver Steeg, and Rachel Warren provided invaluable research assistance on this project. This project would not have been possible without the dedication of our project partners in Togo, especially Minister Cina Lawson, Shegun Bakari, Stanislas Telou, Leslie Mills, Kafui Ekouhoho, Morlé Koudeka, and Attia Byll. The team at GiveDirectly was instrumental in implementing the Novissi expansion studied in this paper (especially Han Sheng Chia, Michael Cooke, Kristen Lee, Alex Nawar, and Daniel Quinn). We thank Esther Duflo, Luis Encinas, Tina George, Rema Hanna, Ethan Ligon, and Ben Olken for helpful feedback. We are grateful for financial support from Google.org, data.org, the Center for Effective Global Action, the Jameel Poverty Action Lab, the Global Poverty Research Lab at Northwestern University, and the World Bank, which financed the phone surveys and data collection under the WURI program. Blumenstock is supported by NSF award IIS – 1942702. Authors retained full intellectual freedom in conducting this research, and as such all opinions and errors are our own. The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

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 They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2021 by Emily Aiken, Suzanne Bellue, Dean Karlan, Christopher R. Udry, and Joshua Blumenstock. 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.

Machine Learning and Mobile Phone Data Can Improve the Targeting of Humanitarian Assistance
Emily Aiken, Suzanne Bellue, Dean Karlan, Christopher R. Udry, and Joshua Blumenstock NBER
Working Paper No. 29070
July 2021, Revised June 2022
JEL No. C55,I32,I38,O12,O38

ABSTRACT

The COVID-19 pandemic has devastated many low- and middle-income countries, causing widespread food insecurity and a sharp decline in living standards. In response to this crisis, governments and humanitarian organizations worldwide have distributed social assistance to over 1.5 billion people. Targeting is a central challenge in administering these programs: given available data, how does one rapidly identify those with the greatest need? Here we show that data from mobile phone networks can improve the targeting of humanitarian assistance. Our approach uses traditional survey data to train machine-learning algorithms to recognize patterns of poverty in mobile phone data; the trained algorithms can then prioritize aid to the poorest mobile subscribers. We evaluate this approach by studying Togo's flagship emergency cash transfer program, which used these algorithms to disburse millions of dollars in COVID-19 relief aid. Our analysis compares outcomes – including exclusion errors, total social welfare, and measures of fairness – under different targeting regimes. Relative to the geographic targeting options considered by the Government of Togo, the machine learning approach reduces errors of exclusion by 4-21%. Relative to methods requiring a comprehensive social registry (a hypothetical exercise; no such registry exists in Togo), the machine learning approach increases exclusion errors by 9-35%. These results highlight the potential for new data sources to complement traditional methods for targeting humanitarian assistance, particularly in crisis settings when traditional data are missing or out of date.

Emily Aiken
102 South Hall
University of California, Berkeley
Berkeley, CA 94720
USA
emilyaiken@berkeley.edu

Suzanne Bellue
University of Mannheim
sbellue@mail.uni-mannheim.de

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

Christopher R. Udry
Northwestern University
Department of Economics
Weinberg College of Arts and Sciences
2211 Campus Drive #3247
Evanston, IL 60208
and NBER
christopher.udry@northwestern.edu

Joshua Blumenstock
School of Information
102 South Hall
University of California, Berkeley
jblumenstock@berkeley.edu

The COVID-19 pandemic has led to a sharp decline in living standards across the world, as policies designed to stop the spread of the disease have disrupted ordinary economic activity. Economically vulnerable households in low- and middle-income countries are among the hardest hit, with over 100 million individuals estimated to have transitioned into extreme poverty since the onset of the pandemic⁴.

To offset the most severe consequences of this sudden income decline, governments and humanitarian organizations around the world have mobilized relief efforts. Gentilini et al. (2021) estimates that over 3,300 new social assistance programs have been launched since early 2020², providing over \$800 billion dollars in cash transfer payments to over 1.5 billion people (roughly one fifth of the world's population).

The overwhelming majority of COVID-19 response efforts – and the majority of cash transfer programs globally – provide *targeted* social assistance. In other words, specific criteria are used to determine potential eligibility, typically some proxy for socioeconomic status. In most wealthy nations, governments rely on recent household income data to determine program eligibility⁵. However, in low and lower-middle income countries (LMICs), where economic activity is often informal and based on home-produced agriculture, governments typically do not observe income for the vast majority of the population³. Other potential sources of targeting data are often incomplete or out of date^{6,7}; for example, only half of the poorest countries having completed a census in the past 10 years⁸. In such contexts, data gaps preclude governments from implementing well-targeted social assistance programs^{9,10}.

Here we develop, implement, and evaluate a new approach to targeting social assistance based on machine learning algorithms and non-traditional “big data” from satellites and mobile phone networks. This approach leverages recent advances in machine learning that show that such data can help accurately estimate the wealth of small geographic regions^{11–13} and individual mobile subscribers^{14,15}. It also builds on a rich economics literature on the design of appropriate mechanisms for targeting social assistance^{3,16–19}. See Supplementary Discussion section 1 for a summary of prior work.

Humanitarian Response to COVID-19 in Togo

Our results are based on the design and evaluation of Novissi, Togo's flagship emergency social assistance program. The Government of Togo launched Novissi in April 2020, shortly after the first COVID-19 cases appeared in the country. As economic lockdown orders forced many Togolese to stop working and led to widespread food insecurity (Supplementary Fig. 1), Novissi aimed to provide subsistence cash relief to those most impacted (see <https://novissi.gouv.tg/>). Eligible beneficiaries received bi-weekly payments of roughly USD \$10. In an effort to minimize in-person contact, Novissi enrollment and payment was digital: beneficiaries registered using their mobile phones and transfers were made via mobile money. Full details on the Novissi program are provided in Methods, ‘The COVID-19 Pandemic in Togo.’

When the government first launched Novissi, it did not have a traditional social registry that could be used to assess program eligibility, and had neither the time nor the resources to build such a registry in the middle of the pandemic. The most recent census, which was completed in 2011, did not contain information on household wealth or poverty; more recent national surveys

on living standards only contacted a small fraction of all households (see Methods, ‘The COVID-19 Pandemic in Togo’). Instead, Novissi eligibility was determined based on data contained in a national voter registry that had been updated in late 2019. Specifically, benefits were initially disbursed to individuals who met three criteria: (1) “self-targeted”¹⁶ by dialing in to the Novissi platform and entering basic information from their mobile phone and; (2) registered to vote in specific regions (the program initially focused on the Greater Lomé region around the capital city); and (3) self-declared to work in an informal occupation in their voter registration. The decision to target informal occupations helped prioritize benefits to people who were forced to stop working at the onset of the crisis. However, this approach does not necessarily target benefits to the poorest households in the country (Supplementary Fig. 2).

Our research efforts focused on helping the government expand the Novissi program from informal workers in Greater Lomé to poorer individuals in rural regions of the country, and were designed to meet the government’s two stated policy objectives: first, to direct benefits to the poorest geographic regions of the country; and second, to prioritize benefits to the poorest mobile subscribers in those regions. (Individuals without access to a mobile phone could not receive Novissi payments, which were digitally delivered using mobile money – see Methods, ‘Program Exclusions’ for details). The approach we developed, which uses machine learning to analyze non-traditional data from satellites and mobile phone networks, has two distinct steps (Extended Data Fig. 1).

Targeting with Mobile Phone Data

In the first step, we obtained public micro-estimates of the relative wealth of every 2.4km by 2.4km region in Togo, which were constructed by applying machine learning algorithms to high-resolution satellite imagery¹³. These estimates provide an indication of the relative wealth of all the households in each small grid cell; we take the population-weighted average of these grid cells to estimate the average wealth of every *canton*, Togo’s smallest administrative unit (see Methods, ‘Poverty Maps’).

In the second step, we estimated the average daily consumption of each mobile phone subscriber by applying machine learning algorithms to mobile phone metadata provided by Togo’s two mobile phone operators (see Methods, ‘Data Privacy Concerns’). Specifically, we conducted surveys with a large and representative sample of mobile phone subscribers, used the surveys to measure the wealth and/or consumption of each subscriber, and then matched the survey-based estimates to detailed metadata on each subscriber’s history of phone use. This sample was used to train supervised machine learning algorithms that predict wealth and consumption from phone use (Pearson ρ ranges from 0.41-0.46 – see Methods, ‘Predicting Poverty from Phone Data’)^{14,15,20}. This second step is similar in spirit to a traditional proxy means test (PMT), with two main differences: we used a high-dimensional vector of mobile phone features instead of a low-dimensional vector of assets to estimate wealth; and we used machine learning algorithms designed to maximize out-of-sample predictive power instead of the traditional linear regression that maximizes in-sample goodness-of-fit²¹.

Results

Our main analysis evaluates the performance of this new targeting approach that combines machine learning and mobile phone data – which we refer to more succinctly as the *phone-based approach* – by comparing targeting errors using the phone-based approach to targeting errors under three counterfactual approaches: a geographic targeting approach that the government piloted in summer 2020 (where all individuals are eligible within the poorest *prefectures*, Togo’s admin-2 level, or poorest *cantons*, Togo’s admin-3 level); occupation-based targeting (including Novissi’s original approach to targeting informal workers as well as an “optimal” approach to targeting the poorest occupation categories in the country); and a parsimonious method based on phone data without machine learning (that uses total expenditures on calling and texting as a proxy for wealth).

We present results that compare the effectiveness of these different targeting mechanisms in two different scenarios. First, we evaluate the actual policy scenario faced by the government of Togo in September of 2020, which involved distributing cash to 60,000 beneficiaries within Togo’s 100 poorest cantons. This first scenario is evaluated using data collected in a large phone survey we designed for this purpose and conducted in September 2020. The “ground truth” measure of poverty in this first scenario is a PMT, as consumption data could not be feasibly collected in the phone survey. (The PMT is based on a stepwise regression procedure, described in Supplementary Methods section 3, which captures roughly 48% of the variation in consumption). Thus, for the first scenario focused on the rural Novissi program, all targeting methods are evaluated with respect to this PMT. The phone-based machine-learning model is likewise trained using the PMT as ground truth. Second, we simulate and evaluate a more general and hypothetical policy scenario in which the government is interested in targeting the poorest individuals nationwide; this scenario is evaluated using national household survey data collected in the field by the government in 2018-2019. The second simulation uses consumption as the “ground truth” measure of poverty. These data are described in Methods, ‘Data’ and details on the evaluation are in Methods, ‘Targeting Evaluations’.

In the first scenario focused on reaching the poorest people in the 100 poorest cantons, we find that the phone-based approach to targeting significantly reduces errors of exclusion (true poor who are mistakenly deemed ineligible) and errors of inclusion (non-poor who are mistakenly deemed eligible), relative to the other feasible approaches to targeting available to the Government of Togo (Figure 1a and first four columns of Table 1). We focus on the ability of each targeting method to reach the poorest 29% in each of the two survey datasets, since the rural Novissi expansion only had sufficient funding to provide benefits to 29% of individuals in eligible geographies (Extended Data Table 1 and Extended Data Table 2 evaluate performance using alternative poverty thresholds). Using a PMT as a measure of “true” poverty status, phone-based targeting (AUC= 0.70) outperforms the other feasible methods of targeting rural Novissi aid (e.g., AUC=0.59-0.64 for geographic blanket targeting). As a result, errors of exclusion (defined as $1 - \text{Recall}$) are lower for the phone-based approach (53%) than for feasible alternatives (59%-78%).

Phone-based targeting likewise outperforms most feasible methods when we simulate the targeting of a hypothetical national anti-poverty program (Figure 1b and last four columns of Table 1). Here, the phone-based approach is more effective at prioritizing the poor (AUC = 0.73) than geography-based alternatives (AUC = 0.66 - 0.68), and likewise leads to lower exclusion

errors (50%) than most feasible alternatives (52%-76%). One exception in this hypothetical program is occupation-based targeting: while the Novissi program's original criteria of targeting informal workers would not scale well to a national program (76% exclusion errors), an alternative "optimal" occupation-based approach that we develop (Methods, 'Experimental Design'), which assigns all transfers to the poorest occupational category (agricultural workers), slightly outperforms phone-based targeting (48% exclusion errors).

Taken together, the results in Table 1 indicate that the phone-based targeting approach was more effective in the actual rural Novissi program than it would be in a hypothetical nationwide program. Our analysis suggests that the benefits of phone-based targeting are greatest when the population under consideration is more homogeneous, and when there is less variation in other factors (such as place of residence) that are used in more traditional approaches to targeting (Methods, 'Targeting Methods and Counterfactuals'). For instance, when we restrict the simulation of the hypothetical national program to households in rural areas, the gains from phone-based targeting increase (Supplementary Table 1).

We likewise find that the performance benefits of phone-based targeting increase as programs seek to target the most extreme poor. This increase can be seen by comparing Table 1, where targeting performance is measured by how many of the poorest 29% receive benefits, to Extended Data Table 1, which measures whether households below the extreme poverty line (\$1.43 per capita daily consumption) receive benefits, and Extended Data Table 2, which measures whether households below the poverty line (\$1.90 per capita daily consumption) receive benefits. While all targeting methods perform better at targeting the extreme poor, the differential between the phone-based approach and other methods is greater when the consumption threshold is lower. (In this analysis, the wealth distribution of the underlying population is important: since over half of the Togolese population is below the poverty line, the targeting methods are attempting to differentiate between different gradations of poverty. Just as precision increases as the target population grows – i.e., from Table 1 to Extended Data Table 1 to Extended Data Table 2 -- results may differ in contexts where the target population is much smaller).

The phone-based approach we develop relies heavily on machine learning to construct a poverty score for each mobile subscriber, where eligibility is a complex function of how the subscriber uses their phone (Extended Data Table 3). We also consider an alternative approach that does not use machine learning, but instead simply targets mobile phone subscribers with the lowest mobile phone expenditures over the preceding months (Methods, 'Parsimonious Phone Expenditure Method'). We find that this "phone expenditure" method (AUC= 0.57 for rural Novissi and 0.63 in for the hypothetical national program – see Table 1) performs substantially worse than the ML-based model (AUC= 0.70 for rural Novissi and 0.73 in for the hypothetical national program). While the phone expenditure model requires much less data and may be easier to implement, this parsimony increases targeting errors – and may also introduce scope for strategic "gaming" if used repeatedly over time.

An important factor in the success of the machine learning model is the fact that it was trained on representative survey data collected immediately prior to the program's expansion. Since an individual's poverty status can change over time, and since the best phone-based predictors of wealth may also change, a model trained in one year or season may not perform well if applied in a different year or season. In Togo, we find that when the machine learning model or the

mobile phone data are roughly 18 months out of date, predictive accuracy decreases by 4-6% and precision drops by 10-14% (Extended Data Table 4 and Methods, ‘Temporal Stability of Results’). These losses are nearly as large as the gains that phone-based targeting provides over geographic targeting – a finding that underscores the importance of training the model with current and representative data.

We also compare the phone-based approach to alternative targeting approaches that require a recent and comprehensive social registry. While the Government of Togo did not have such a registry, this comparison helps situate this method relative to other methods commonly used by development researchers and policymakers. These results, shown in Panel B of Table 1, can only be simulated using the national in-person survey, since the phone survey did not collect consumption data. The results are more ambiguous: the phone-based approach (AUC = 0.70-0.73) is approximately as accurate as targeting using an *asset-based wealth index* (AUC = 0.55-0.75), but less accurate than using a *poverty probability index* (AUC = 0.81) or perfectly-calibrated *proxy-means test* (AUC = 0.85) – see Methods, ‘Survey Data’ for the differences between these indices. We note, however, that the performance of the “perfectly calibrated” PMT may substantially over-estimate the performance of a real-world PMT, which declines steadily over time since calibration (Methods, ‘Targeting Methods and Counterfactuals’)^{22,23}.

Social welfare and fairness

Improvements in targeting performance translate to an increase in social welfare. Using the constant relative risk-aversion (CRRA) utility function, we calculate aggregate welfare under the phone-based approach and each of the counterfactual targeting approaches. Under the CRRA assumptions, individual utility is a concave function of consumption. By assuming a fixed budget – which we fix at a size analogous to that of the Novissi rural aid program, which had a budget of USD 4 million to distribute among 154,238 program registrants – and equal transfer sizes to all beneficiaries, we simulate the distribution of benefits among eligible individuals at counterfactual targeting thresholds to trace out social welfare curves for each targeting method. This social welfare analysis also allows us to identify the optimal beneficiary share and corresponding transfer size. **Figure 2** shows the utility curves for each of the targeting methods simulated, separately for the two datasets. Note that phone-based targeting, geographic blanketing, and an asset-based wealth index all achieve approximately the same maximum utility in the hypothetical national program, but phone-based targeting dominates in the rural Novissi program. Also note that all targeting methods outperform a universal basic income scheme if the beneficiary share and transfer size is well-calibrated.

These utilitarian welfare gains suggest that society as a whole will benefit from improved targeting, but do not imply that all subgroups of the population will benefit equally. Indeed, there is growing concern that algorithmic decision-making can unfairly discriminate against vulnerable groups²⁴⁻²⁶. To address these concerns in the context of the Novissi program, we audit the fairness of each targeting method across a set of potentially sensitive characteristics, while noting that notions of fairness and parity are contested and often in tension.²⁷ Figure 3a shows, as an example, that the phone-based approach does not cause women to be systematically more likely to be incorrectly excluded by the targeting mechanism from receiving benefits than men (see also Methods, ‘Fairness’). Likewise, the phone-based approach does not create significant exclusion errors for specific ethnic groups (Figure 3b), religions, age groups, or type of household, though

there are small differences in targeting accuracy between groups (Extended Data Figure 2). We also compare the fairness of the phone-based approach to several other targeting approaches by evaluating each method’s demographic parity, i.e., the extent to which each method under- or over-targets specific demographic subgroups, relative to that group’s true poverty rate (Figure 3c-d and Extended Data Figure 3). Overall, we find that none of the targeting methods analyzed naively achieves perfect parity across subgroups; a phenomenon referred to as “no fairness through unawareness.”²⁸ The largest parity differences occur with geographic targeting methods.

Exclusions and limitations

This novel approach to targeting requires careful consideration of the ways in which individuals can be incorrectly excluded from receiving program benefits (Methods, ‘Program Exclusions’). Our analysis highlights six main sources of exclusion errors for the expansion of Novissi (Table 2): (i) beneficiaries must have a SIM card and access to a mobile phone (2018-2019 field survey data indicate that 65% of adults and 85% of households have a phone; see also Supplementary Fig. 3); (ii) they must have used their SIM card recently, in order to generate a poverty score (between 72% and 97% of program registrants) (iii) they must be a registered voter (roughly 87% of adults); (iv) they must self-target and attempt to register (roughly 40% of eligible individuals attempted); (v) they must succeed in registering, which requires basic reading and digital literacy (72% succeed); and (vi) they must be successfully identified as eligible by the ML algorithm (47% recall, per Table 1). Many of these sources of possible exclusion overlap; Extended Data Table 5 thus estimates, based on the 2020 phone survey, the extent to which each successive step in registration creates additional exclusions. These results highlight the fact that algorithmic targeting errors are an important source of program exclusion, but that real-world programs also face structural and environmental constraints to inclusion.

More broadly, our analysis shows how non-traditional “big” data and machine learning can improve the targeting of humanitarian assistance. Beyond the gains in targeting performance, a key advantage of this approach is that it can be deployed quickly and responsively. In Togo, the government’s objective was to deliver benefits to the poorest people in the country, so our efforts focused on training a machine learning model to target the poor. In other settings, such as following natural disasters, the people most impacted by adverse events may not be the poorest²⁹. With high-frequency phone data available in near real-time, related techniques might be used to more dynamically prioritize the people with the greatest need. For example, it may be possible to train a machine learning algorithm to identify people whose consumption fell by the greatest amount, based on changes in patterns of phone use following a crisis. Another possibility would be to simply use location information from mobile phone data to prioritize people likely to live in impacted regions (Methods, ‘Location-Based Targeting’).

It is important to emphasize that our phone-based approach is far from perfect, and may lead to important errors of both exclusion and inclusion. There are also practical limitations to this approach, for instance regarding data access and privacy; several such considerations are addressed in Supplementary Discussion section 2. Moreover, our results do not imply that mobile phone-based targeting should replace traditional approaches reliant on proxy means tests or community-based targeting. Rather, these new methods provide a rapid and cost-effective supplement that may be most useful in crisis settings or in contexts where traditional data sources are incomplete or out of date. We believe future work should explore how real-time data

sources, such as the phone data used by Novissi, can be best combined with more traditional field-based measurements, so that these complementary data sources can be best integrated in the design of inclusive systems for social protection²⁰.

References

1. Egger, D. *et al.* Falling living standards during the COVID-19 crisis: Quantitative evidence from nine developing countries. *Science Advances* **7**, eabe0997 (2021).
2. Gentilini, U., Almenfi, M., Orton, I. & Dale, P. Social Protection and Jobs Responses to COVID-19: A Real-Time Review of Country Measures. *World Bank Policy Brief* (2020).
3. Hanna, R. & Olken, B. A. Universal Basic Incomes versus Targeted Transfers: Anti-Poverty Programs in Developing Countries. *Journal of Economic Perspectives* **32**, 201–226 (2018).
4. Lakner, C., Yonzan, N., Mahler, D., Aguilar, R. A. & Wu, H. *Updated estimates of the impact of COVID-19 on global poverty: Looking back at 2020 and the outlook for 2021*. <https://blogs.worldbank.org/opendata/updated-estimates-impact-covid-19-global-poverty-looking-back-2020-and-outlook-2021>.
5. Mirrlees, J. A. An Exploration in the Theory of Optimum Income Taxation. *The Review of Economic Studies* **38**, 175–208 (1971).
6. Jerven, M. *Poor numbers: how we are misled by African development statistics and what to do about it*. (Cornell University Press, 2013).
7. Serajuddin, U., Wieser, C., Uematsu, H., Dabalén, A. L. & Yoshida, N. *Data deprivation : another deprivation to end*. 1–24 <http://documents.worldbank.org/curated/en/700611468172787967/Data-deprivation-another-deprivation-to-end> (2015).
8. Yeh, C. *et al.* Using publicly available satellite imagery and deep learning to understand economic well-being in Africa. *Nature Communications* **11**, 2583 (2020).
9. Bank, W. *World Development Report 2021: Data for Better Lives*. (The World Bank, 2021).
10. Coady, D., Grosh, M. & Hoddinott, J. Targeting outcomes redux. *The World Bank Research Observer* **19**, 61–85 (2004).
11. Jean, N. *et al.* Combining satellite imagery and machine learning to predict poverty. *Science* **353**, 790–794 (2016).
12. Engstrom, R., Hersh, J. S. & Newhouse, D. L. *Poverty from space : using high-resolution satellite imagery for estimating economic well-being*. 1–36 <http://documents.worldbank.org/curated/en/610771513691888412/Poverty-from-space-using-high-resolution-satellite-imagery-for-estimating-economic-well-being> (2017).
13. Chi, G., Fang, H., Chatterjee, S. & Blumenstock, J. E. Micro-Estimates of Wealth for all Low- and Middle-Income Countries. *In Submission* (2021).
14. Blumenstock, J. E., Cadamuro, G. & On, R. Predicting poverty and wealth from mobile phone metadata. *Science* **350**, 1073–1076 (2015).
15. Blumenstock, J. E. Estimating Economic Characteristics with Phone Data. *American Economic Review: Papers and Proceedings* **108**, 72–76 (2018).
16. Nichols, A. L. & Zeckhauser, R. J. Targeting Transfers through Restrictions on Recipients. *The American Economic Review* **72**, 372–377 (1982).
17. Grosh, M. E. & Baker, J. L. *Proxy means tests for targeting social programs*. (The World Bank, 1995). doi:10.1596/0-8213-3313-5.

18. Alatas, V., Banerjee, A., Hanna, R., Olken, B. A. & Tobias, J. Targeting the Poor: Evidence from a Field Experiment in Indonesia. *American Economic Review* **102**, 1206–40 (2012).
19. Blumenstock, J. Machine learning can help get COVID-19 aid to those who need it most. *Nature* **581**, (2020).
20. Aiken, E., Bedoya, G., Coville, A. & Blumenstock, J. E. Program Targeting with Machine Learning and Mobile Phone Data: Evidence from an Anti-Poverty Intervention in Afghanistan. *Working Paper* (2020).

Methods

THE COVID-19 PANDEMIC IN TOGO

Togo is a small country of roughly 8 million in West Africa. Over 50% of the population lives below the international poverty line. Shortly after the first COVID-19 cases were confirmed in Togo in early March 2020, the government imposed economic lockdown orders to prevent the spread of the disease. These lockdowns forced many Togolese to stop working, raising concerns about the potential for rising food insecurity (Supplementary Fig. 1).

On April 8, 2020, the government launched the Novissi program, where “Novissi” means “solidarity” in the Ewé language. According to Minister Cina Lawson, Novissi “was built and designed in order to help those people who are the most vulnerable population and the most impacted by the anti-COVID measures.”³⁰ Novissi was initially designed to provide benefits to informal workers in Greater Lomé, the large metropolitan area surrounding the capital city where the lockdown orders were initially focused. The rationale for targeting informal workers was that they were more likely to be vulnerable, and more likely to be impacted by the lockdown orders.

To determine eligibility for Novissi, the government relied upon a national voter registry that was updated in late 2019, in which individuals indicated their home location and occupation. At the time, the voter registry contained 3,633,898 entries, which the electoral commission reports is equivalent to 87% of the total adult population (see Table 2 for details).

Receiving Novissi benefits required that individuals register by dialing in to the Novissi USSD platform from a mobile phone. Thus, registration initially required (i) a valid and unique voter ID linked to an eligible occupation from an eligible location; (ii) a valid SIM card, and (iii) access to a mobile phone. A smartphone was not required for registration; the USSD platform was accessible from a basic phone. Since phone sharing is common in Togo, multiple SIM cards could be registered through a single phone (so long as each SIM was then linked to a valid voter ID). See Methods, ‘Program Exclusions’ for a discussion of the extent to which voter and phone requirement may have led to program exclusions.

Eligible female beneficiaries were then paid 12,250 FCFA (USD \$22.50) per month; men received 10,500 FCFA (USD \$20) per month. The payments were disbursed in two bi-weekly installments, for three months, using existing mobile money infrastructure managed by the country’s two mobile network operators. The system was designed to be 100% digital, so that registration, eligibility determination, and payment could all be accomplished without face-to-face contact. Novissi was promoted actively through radio advertisements and community leaders, and 4.4 million registration attempts were reported on the day the program launched. In this first phase of Novissi, which focused on Greater Lomé, roughly 510,000 beneficiaries received payments.

During the summer of 2020, in response to localized outbreaks of COVID-19, the government piloted an expansion of Novissi based on geographic targeting. In this geographically targeted expansion, all individuals registered to vote in the Soudou canton were made eligible for Novissi benefits. The geographic targeting was determined primarily by public health considerations, and not by poverty rates. In total, roughly 5,800 beneficiaries were paid through this geographically targeted program.

Our analysis focuses on a second phase of Novissi, which was initiated after the Novissi program in Greater Lomé had terminated. Specifically, in partnership with the NGO GiveDirectly, the government wished to expand Novissi eligibility to the rural poor. The policy mandate from the government was to (i) prioritize benefits to people living in Togo's 100 poorest cantons (of the 397 cantons nationally), where the number 100 was selected by the government in order to balance the desire to focus on the poorest villages, without focusing excessively on specific regions; and (ii) prioritize the poorest individuals in those 100 cantons.

During the second phase of Novissi, registration and enrollment used several of the same steps described above: individuals were required to have a voter ID registered in one of the 100 poorest cantons, and they had to self-register using a mobile phone with a unique SIM card. However, the individual's occupation was not used to determine eligibility; instead, the estimated wealth of the individual, based on the ML methods described in this paper, were used to limit eligibility to the estimated poorest subscribers in those 100 cantons.

DATA SOURCES

Survey Data

Our core analysis relies heavily on two surveys conducted by Togo's Institut National de la Statistique et des Études Economiques et Démographiques (INSEED). The first survey, which is nationally representative, was conducted in the field in 2018 and 2019 ($N = 6,171$). The second survey was conducted over the phone in September 2020, and is representative of mobile network subscribers inferred to be living in rural cantons eligible for Novissi aid ($N = 8,915$). We use these two different survey datasets because neither dataset is sufficient by itself for the analysis we require: The 2020 survey did not collect consumption data, which is important for evaluating certain counterfactuals; the 2018-19 survey is representative only at the prefecture level, and only surveyed a small number of households in the 100 poorest cantons that were eligible for Novissi. (We had planned to conduct a large in-person survey in early 2021 that would provide the single point of focus for this paper, but were forced to postpone the survey indefinitely due to a resurgence in COVID-19).

2018-2019 Field Survey: Our first survey dataset was obtained from a nationally representative household survey. Specifically, 540 enumeration areas (EAs) were drawn at random from Togo's approximately 6,000 EAs, with weight proportional to the size of the EA in the last national census (conducted in 2011). 12 households were then drawn at random from each of the selected EAs to be interviewed, for a total of 6,172 households. Surveys, which lasted about three hours, were conducted in two waves, with the first wave between October and December 2018 and the second wave between April and June 2019. We remove one observation that is missing consumption expenditure and asset data, leaving 6,171 observations. Interviews took place with the head of household when possible, and alternatively with the most knowledgeable adult present. Answers were recorded by enumerators on tablets using SurveyCTO software.

As part of the survey's recontact protocol, phone numbers were requested from a representative of each household; 4,618 households (75%) of households are matched to a phone number. The data do not include an identifier for which member of the household the phone number belongs to. 4,171 households have phone numbers that contain at least one transaction in our mobile phone transaction logs in the three months prior to their survey date (90% of households with

phone numbers), leading to a matched survey-mobile phone dataset of $N = 4,171$. Note that this matched dataset is not nationally representative nor necessarily representative of mobile phone subscribers, as there is selection in which households and household members provide phone numbers.

2020 Phone Survey: Our second survey dataset is obtained from a phone survey conducted over two weeks in September 2020. The survey lasted approximately 40 minutes, and covered demographics, asset ownership, and well-being. Answers were recorded by enumerators on tablets using SurveyCTO software. Phone numbers for the 2020 phone survey were drawn from mobile phone transaction logs and the sample is representative of subscribers inferred based on their mobile phone data to be living in rural cantons eligible for Novissi aid (See Supplementary Methods section 4). Note that since the sample is drawn based on inferred location, not all interviewees necessarily reside in an aid-eligible canton. The survey includes a question on canton of residence, and 68% of observations report living in a Novissi-eligible canton.

Of the phone numbers drawn, 35% respond, consent to the survey, and complete the entire survey. In total, after removing low-quality surveys and those missing poverty outcomes, the dataset contains 8,915 observations corresponding to individual subscribers. We reweight the survey for nonresponse using the same mobile phone features and machine learning methods described in Methods, ‘Predicting Poverty from Phone Data.’ Our sample weights consist of the inverse of the draw probability and the inverse of the predicted probability of response. More details on the content of the 2020 phone survey, the sampling procedure, and the reweighting procedure are available in Supplementary Methods section 5.

Construction of Poverty Outcomes: We construct four poverty outcomes from the survey data: consumption expenditure (captured in the 2018-2019 field survey only), an asset-based wealth index, a poverty probability index (PPI), and a proxy-means test (PMT).

- *Consumption expenditure:* The consumption expenditure outcome is only available in the dataset from the 2018-2019 field survey. Disaggregated expenditures for more than 200 food and nonfood items are elicited in each household interview. The consumption aggregate is then adjusted for a price index calculated at the prefecture level. The final outcome measure is per-capita adult equivalent household consumption expenditure, which we transform to USD per day.
- *Asset index:* We calculate a PCA asset index for households in the 2018-2019 field survey and for the households associated with individuals interviewed in the 2020 phone survey. Asset indices are constructed with Principal Component Analysis (PCA). The asset index is constructed from 24 underlying binary asset variables in the 2018-2019 field survey and 10 underlying binary asset variables in the 2020 phone survey. The asset indices for the two surveys are constructed independently, from different sets of assets, and therefore do not share a basis vector. The basis vector for each index is shown in Supplementary Table 2. The asset index explains 31.50% of the variance in asset ownership in the 2018-2019 field survey, and 53.45% of the variance in asset ownership in the 2020 phone survey. However, the variance explained in the two indices should not be directly compared since there are far fewer assets recorded in the 2020 phone survey than in the 2018-2019 field survey. We also note that the asset index for the 2020 phone survey dataset is dominated by variation in ownership of three assets (toilet, radio, and

motorcycle; see Supplementary Table 2) and is therefore considerably less smooth than the asset index in the 2018-2019 phone survey dataset.

- *Poverty probability index (PPI)*: We use the scorecard for the current poverty probability index used by Innovations for Poverty Action (<https://www.povertyindex.org/country/togo>). The index is calibrated based on a nationally representative survey conducted by INSEED in 2015 ($N = 2,335$). “Poverty probability” is scored based on ten household questions, including region of residence, education of adults and children, asset ownership, and consumption of sugar. We calculate the PPI only for households in the 2018-2019 field survey, as the data necessary for all components were not collected in the 2020 phone survey.
- *Proxy-means test (PMT)*: Using the data from the 2018-2019 field survey, we follow a stepwise forward selection process to select the 12 asset and demographic variables that are jointly most predictive of per-capita household consumption (see Supplementary Fig. 4 and Supplementary Methods section 3 for details). We use these variables to construct a consistent proxy-means test (PMT) for the 2018-2019 field survey and the 2020 phone survey. Following recent literature, we use a regularized linear model (Ridge regression) rather than a simple linear regression to maximize out-of-sample accuracy^{21,26}. For the 2018-2019 field survey, PMT “consumption” estimates are produced out-of-sample over 10-fold cross validation. For the 2020 phone survey, we train the Ridge regression on the entire 2018-2019 field survey sample and use the fitted model to produce PMT “consumption” estimates for each phone survey observation. Over 10-fold cross validation, the PMT explains 48.35% of the variance in log-transformed consumption expenditure in the 2018-2019 field survey. This explanatory power is similar to that of other national-scale PMTs reported in Indonesia, Peru and Jamaica (41%-66%)^{3,17,31}. The weights for the PMT are included in Supplementary Table 3. Since they are trained to predict consumption, PMT “consumption” estimates can be interpreted as estimated USD/day.
- *Rural-specific PMT*: We follow another stepwise forward selection process using the 2018-2019 field survey restricted to households in rural areas ($N = 3,895$) to create a PMT specific to rural areas with 12 components. The weights for the rural-specific PMT are shown in Supplementary Table 4. Over 10-fold cross-validation the rural-specific PMT explains 17% of the variation in log-transformed consumption expenditure in the 2018-2019 field survey restricted to rural areas. We note that this explanatory power is substantially lower than that of other rural-specific PMTs evaluated in past work in Jamaica and Burkina Faso (36%-45%)^{32,33}. We produce out-of-sample values for the rural-specific PMT over cross validation for the 2018-2019 field survey, and use the fitted model to produce values for the 2020 phone survey. We mean-impute the rural-specific PMT for observations that do not have all necessary components in the 2020 phone survey dataset ($N = 18$). The correlation between the rural-specific PMT and general PMT is 0.75 in the 2018-2019 survey dataset restricted to rural areas, and 0.76 in the 2020 phone survey dataset.

Construction of Occupation Categories: We use self-reported occupation (of the household head for the 2018-2019 field survey, and of the respondent for the 2020 phone survey) to categorize occupations and later simulate occupation-based targeting. We first classify each of the self-

reported occupations according to the occupation categories in the Novissi registry. We identify which of these categories are informal (in the Novissi registry, more than 2,000 unique occupations are considered informal – some of the most common ones are vendors, hairdressers, taxi drivers, tailors, construction workers, and the unemployed). We further classify occupations in 10 broad categories according to the Afrostat system (<https://www.afrostat.org/nomenclatures/>). Supplementary Table 5 records these categories, along with the proportion in each category in each of the two surveys and associated average consumption.

Summary Statistics: Supplementary Table 6 presents summary statistics on each of the two surveys; for the 2018-2019 household survey, results are presented separately for households who provide phone numbers (further broken down into those with phone numbers that match to the mobile phone metadata and those whose phone numbers do not match), and those without phone numbers. Note that since phone numbers for the 2018-2019 household survey were collected for a recontact protocol, a household without a phone number could represent a household without a phone or one that refused to be contacted for further surveys. We find that households providing phone numbers (average consumption = \$2.56/day) are less poor than households not providing them (average consumption = \$1.75/day); among those associated with a phone number, households that do not match to mobile phone metadata (average consumption = \$2.21/day) are poorer than those that do (average consumption = \$2.59/day). These patterns are consistent with related work in Afghanistan in which phone numbers were collected for the purpose of matching to mobile phone metadata. That study found that households with phones were wealthier than those without, and households associated with a matched phone number were wealthier than those that did not match²⁰.

Comparing summary statistics from the 2020 phone survey and 2018-2019 household survey, respondents to the 2020 survey tend to be poorer (average PMT = 1.62 vs. 2.10), younger (average age = 33 vs. 44), and more predominantly male (23% women vs. 28% women). These differences are not surprising given that the 2020 survey was conducted in rural areas whereas the 2018-2019 household survey was designed to be nationally representative.

Poverty Maps

To simulate geographic targeting, we rely on poverty maps of Togo's prefectures (admin-2 level, 40 prefectures) and cantons (admin-3 level, 397 cantons). In the 2018-2019 field survey, the latitude and longitude of each household were recorded by enumerators as part of the interview, so we map each observation to a prefecture and canton using the geographic coordinates. For the 2020 phone survey, we ask each respondent to report their prefecture and canton of residence.

Prefecture poverty map: INSEED completed a survey-based poverty mapping exercise in 2017. Specifically, a proxy-means test was calibrated on a small consumption sample survey conducted in 2015 ($N = 2,335$). 26,902 households were then surveyed in the field over three weeks in 530 EAs, sampled to be representative at the prefecture level. The interview included questions on demographics, education, asset ownership, and household characteristics that made up the PMT. The calibrated PMT was then used to infer the "consumption" of each household, and observations were aggregated to estimate the percentage of the population living under the Togo-specific poverty line of USD 1.79/day in each prefecture. Supplementary Fig. 5 shows the resulting poverty map. For validation, we evaluate the correlation between prefecture-level

poverty rates from the poverty mapping exercise and average consumption in the 2018-2019 field survey. The Pearson correlation coefficient is -0.78, and the Spearman correlation coefficient is -0.70.

Canton poverty map: When COVID-19 first appeared in Togo in early 2020, it had been at least ten years since a household survey had been conducted in Togo that was representative at the canton level. Togo's last census was conducted in 2011, but did not include information on income, consumption, or asset ownership. We therefore rely on recently-produced publicly available satellite-based estimates of poverty which use deep learning models trained on DHS data from neighboring countries to estimate the average relative wealth of each 2.4km tile in Togo.¹³ We overlay the resulting tile-level wealth estimates with high-resolution estimates of population density inferred from satellite imagery³⁴ to obtain population-weighted average wealth estimates for each canton, shown in Supplementary Fig. 5. As noted in Chi et al. (2021)¹³, the relative wealth measures are estimated with uncertainty. Thus, for validation, we evaluate the canton-level correlation between average wealth from the satellite-based poverty map and average consumption in the 2018-2019 field survey (though note that the latter survey is not representative at the canton level). The Pearson correlation coefficient is 0.57, and the Spearman correlation coefficient is 0.52.

Mobile Phone Metadata

We obtain mobile phone metadata (call detail records, or CDR) from Togo's two mobile network operators for certain time periods in 2018-2021. We focus on three slices of mobile network data: October - December 2018, April - June 2019, and March - September 2020. The three-month periods in 2018 and 2019 are matched to households interviewed in the first and second wave of the field survey, respectively. The seven-month period in 2020 is matched to outcomes for individuals interviewed in the phone survey in September 2020. Summary statistics on network activity in these periods are shown in Supplementary Fig. 6.

Our CDR data contain the following information:

- *Calls:* Caller phone number, recipient phone number, date and time of call, duration of call, ID of the cell tower through which the call is placed
- *SMS messages:* Sender phone number, recipient phone number, date and time of the message, ID of the antenna through which the message is sent
- *Mobile data usage:* Phone number, date and time of transaction, amount of data consumed (upload and download combined)
- *Mobile money transactions:* Sender phone number, recipient phone number (if peer-to-peer), date and time of the transaction, amount of transaction, and broad category of transaction type (cash in, cash out, peer-to-peer, or bill pay)

October-December 2018 and April-June 2019 CDR: Between October 1 and December 30, 2018, there were a total of 4.84 million unique mobile network subscribers between the two mobile phone networks (where a subscriber is any phone number that places at least one call or SMS on a network). Between April 1 and June 31, 2019, there were a total of 4.89 million mobile network subscribers. We identify spammers on the network as any phone number that

placed an average of over 100 calls or 100 SMS messages per day, and remove any transactions associated with these numbers from our dataset. We remove 232 spammers in the 2018 time period and 162 spammers in the 2019 time period. In the 2018-2019 CDR, we observe only calls, SMS messages, and mobile money transactions (we do not observe mobile data usage).

March-September 2020 CDR: For data between March 1 and September 30, 2020, we observe a total of 5.83 million mobile network subscribers (note that this subscriber population does not necessarily reflect a 19% increase in subscribers from 2018-2019, since the slice is seven months rather than three months and there is significant month-to-month churn in subscribers; during the 3-month period from July-September 2020 we observe 5.20 million unique subscribers, a 6% increase from the 2019 period). We identify spammers as described above, resulting in the removal of transactions associated with 107 spammers from the 2020 CDR dataset. In the 2020 CDR, we observe calls, SMS messages, mobile data usage, and mobile money transactions.

Featurization: For each subscriber observed on the network in each of the three time periods, we calculate a set of 857-1,042 “CDR features” that describe aspects of the subscriber’s mobile phone behavior. These include:

- *Call and SMS features:* We use open-source library bandicoot³⁵ to produce around 700 features relating to the calls and SMS messages each subscriber places and receives. These range from general statistics (e.g. number of calls/SMS messages, balance of incoming vs. outgoing transactions), to social network characteristics (e.g. number and diversity of contacts), to measures of mobility based on cell tower locations (e.g. number of unique towers, radius of gyration).
- *Location features:* Based on the locations of each of the cell towers in Togo, we calculate information about where each subscriber places their transactions. Specifically, we calculate the number and percentage of calls placed in each of Togo’s 40 prefectures, and the number of unique antennas, cantons, prefectures, and regions that each subscriber visits.
- *International transaction features:* Using country codes associated with phone numbers, we calculate the number of outgoing international transactions, separately for calls and SMS messages. We also calculate the total time spent on outgoing international calls.
- *Mobile money features:* For each of four variables relating to transaction size --- transaction amount, percent of balance, balance before transaction, and balance after transaction --- we calculate the mean, median, minimum, and maximum, separately for incoming and outgoing mobile money transactions. We also calculate the total transaction count for each subscriber (separately for incoming and outgoing) and the total number of unique mobile money contacts (separately for incoming and outgoing). We perform these calculations for all transactions together, as well as separately by transaction type (cash in, cash out, peer-to-peer, bill payments, and other transactions).
- *Mobile data features:* We calculate the total, mean, median, minimum, and maximum mobile data transaction for each subscriber, as well as the standard deviation in transaction size. We also calculate the total number of mobile data transactions and the number of unique days on which data is consumed. Note that mobile data features are

only calculated for the 2020 CDR period, as our 2018-2019 CDR does not include mobile data records.

- *Operator*: In our feature dataset we include a dummy variable for which of the two mobile network operators each subscriber is associated with.

Matching survey and CDR datasets: Using phone numbers collected in surveys, we match survey observations to CDR features. As noted in Methods, ‘Survey Data,’ there are 4,618 households in the 2018-2019 field survey that provide a phone number, of which 4,171 match to CDR (90% of households with phone numbers, and 68% of households overall). We match households surveyed in the first survey wave to features generated in the October-December 2018 CDR period, and households surveyed in the second survey wave to features generated in the April-June 2019 CDR period. To build intuition on the relationships between phone-related features and poverty, Supplementary Fig. 7 compares four CDR features for those above and below the poverty line in the 2018-2019 household survey. Since the 2020 survey was sampled based on the CDR dataset, all 8,915 observations in the 2020 survey dataset are matched to CDR.

Data Privacy Concerns

The CDR data we obtained for each subscriber contain personally identifying information (PII) in the form of the subscriber’s phone number (it does not contain the individual’s name, address, or other PII), as well as other potentially sensitive information such as data about the subscriber’s network and cell tower locations. To protect the confidentiality of these data, we pseudonymized the CDR prior to analysis by hash-encoding each phone number into a unique ID. The data are stored on secure university servers to which access is limited based on a data management plan approved by U.C. Berkeley’s Committee for the Protection of Human Subjects.

We obtained informed consent from all research subjects in the phone survey prior to matching CDR records to survey responses. However, there are still open concerns around the use of CDR by bad actors, particularly as even pseudonymized datasets can frequently be de-anonymized for a subset of observations.^{36,37} Active research on applying the guarantees of *differential privacy* to CDR datasets and associated machine learning models holds promise for balancing the utility of CDR data with privacy concerns.^{38,39} For additional discussion of these considerations, see Supplementary Discussion section 2.

PREDICTING POVERTY FROM PHONE DATA

Machine Learning Methods

We follow the machine learning methods described in prior work^{14,15,20} to train models that predict poverty from CDR features. Specifically, we train a gradient boosting regressor with Microsoft’s LightGBM for the two matched survey-CDR datasets separately. We tune hyperparameters for the model over 3-fold cross validation, with parameters chosen from the following grid:

- *Winsorization of features*: {No winsorization, 1% limit}
- *Minimum data in leaf*: {10, 20, 50}

- *Number of leaves*: {5, 10, 20}
- *Number of estimators*: {20, 50, 100}
- *Learning rate*: {0.05, 0.075, 0.1}

We train and evaluate the model over 5-fold cross validation, with hyperparameters tuned independently on each fold, to obtain out-of-sample estimates of accuracy and out-of-sample predictions of poverty for each observation in our matched survey datasets. We then re-train the model on all survey data (for each of the two datasets separately), record feature importances (the total number of times a feature is split on over the entire forest), and use the final model to generate wealth predictions for every subscriber on the mobile phone network during the relevant time period.

We experiment with training models in this way for each of the relevant poverty outcomes: consumption expenditure, PMT, and asset index for the 2018-2019 field survey dataset and PMT and asset index for the 2020 phone survey dataset. Evaluations of model accuracy are found in Extended Data Table 6. The correlation between the phone-based poverty predictions and a traditional PMT is 0.41, as trained and evaluated on the 2020 phone survey dataset (Extended Data Table 6, Panel C). When trained and evaluated using the national 2018-2019 household survey with consumption data, the correlation between the phone-based poverty predictions and consumption is 0.46 (Extended Data Table 6, Panel A).

Feature Importances: Feature importances for each model are presented in Extended Data Table 3. We note that in examining the feature importances, location-related features (number and percent of calls placed in each prefecture of the country) are very important. The correlation between phone-based poverty predictions using only these location features and a standard PMT is 0.35 when trained and evaluated with the 2020 phone survey (vs. 0.41 using all features). When trained and evaluated with the 2018-2019 field survey, the correlation between location-only phone-based poverty predictions and consumption is 0.42 (vs. 0.46 when using all features). Given the relative importance of location features, we provide more in-depth analysis of the role of geography in phone-based targeting approaches in Methods, “Location-Based Targeting.” Other important features in the full phone-based poverty scores relate to nighttime calling behavior, mobile data usage, and mobile money usage.

Aggregate Validation of CDR-Based Poverty Estimates: Our machine learning models use cross-validation to help limit the potential that the predictions are overfit to the specific surveys on which they are trained (and on which they are later evaluated in the targeting simulations). To provide a more independent test of the validity of the CDR-based estimates, we compare regional aggregates of wealth based on the CDR model to regional estimates of wealth based on household survey data. In this exercise, we predict the consumption of roughly 5 million subscribers in Togo using the ML model trained to predict consumption using the 2018-2019 national household survey, then calculate the average consumption of each prefecture and canton (where each subscribers’ home location is inferred from CDR using standard methods described in Supplementary Methods section 4).

Results, shown in Supplementary Fig. 8, indicate that the CDR-based estimates of regional poverty correlate with survey-based estimates of regional poverty. At the prefecture level, the Pearson and Spearman correlations of CDR-based consumption with survey-based consumption

are 0.92 and 0.83, respectively; the correlations with the proportion of each prefecture living in poverty are -0.76 and -0.74. At the canton level, comparing the CDR-based estimates to the satellite-inferred canton poverty map from Supplementary Fig. 5, we find Pearson = 0.84 and Spearman = 0.68; compared to the average canton consumption in the 2018-19 field survey, Pearson = 0.57 and Spearman = 0.59. These correlations are toward the lower end of the range of correlations observed in prior efforts to estimate regional poverty with CDR^{14,40,41}.

Parsimonious Phone Expenditure Method

In addition to the machine learning method for wealth prediction described above, we are interested in the performance of an intuitive, parsimonious method for approximating poverty with CDR. We focus on a measure of “phone expenditure” on the basis of costs of all calls placed and SMS messages sent by each subscriber. We apply standard rates for calls and SMS messages in Togo: 30 CFA (USD 0.06) to send an SMS message and 50 CFA (USD 0.09) per minute of call time. (These prices represent a typical Togolese phone plan, though there is considerable diversity in special promotions and friends-and-family plans available from Togo's two mobile phone operators, Moov and Togocom). We use these prices to infer the (approximate) amount spent by each subscriber from their outgoing mobile phone transaction logs. We find that the “phone expenditures” method is significantly less accurate than the ML-based method, with a correlation of 0.13 with both the 2020 phone survey PMT and the 2018-2019 household survey's consumption measure (Extended Data Table 6, Panels A and C).

TARGETING EVALUATIONS

Experimental Design

We simulate phone-based and counterfactual targeting methods for reaching the poorest individuals in Togo, using the two survey datasets described in Methods, ‘Survey Data.’ Specifically, for each dataset, we simulate providing benefits to the poorest 29% of observations in the dataset based on a suite of counterfactual targeting options (with sample weights applied), and compare the population targeted to the population that is “truly poor”, where ground truth poverty is determined using two different measurements. With the 2018-2019 in-person survey dataset, our main ground-truth wealth measure is based on consumption expenditure: we evaluate how well proxy measures of poverty reach those with the lowest consumption. For the 2020 phone survey dataset, our main ground-truth wealth measure is based on the proxy-means test described in Methods, ‘Survey Data’ (this is necessary because consumption information was not collected in the phone survey).

Our main targeting evaluations simulate targeting 29% of individuals because the Novissi program had sufficient funds to target 29% of registrants in eligible cantons. The 29th percentile corresponds to a consumption threshold of USD 1.17/day in the 2018-2019 field survey dataset, and a PMT threshold of USD 1.18/day in the 2020 phone survey dataset. Our analysis shows how accurately each targeting method reaches the 29% truly poorest (Table 1), those below the extreme poverty line, defined as three-quarters of the poverty line, or USD 1.43/day (Extended Data Table 1), and those below the international poverty line of USD 1.90/day (Extended Data Table 2).

Our evaluations are designed to measure how effectively several different targeting methods, described below, are at reaching the poorest individual mobile phone owners in each of the two survey populations. We focus on *individuals* rather than *households* because the Novissi program was designed and paid as an individual benefit. While social assistance programs in other countries typically consider the household to be the unit of analysis that determines program eligibility, there is no notion of a household unit in the Novissi program (in part because the government does not possess data that links individuals to households). See Supplementary Discussion section 2 for additional discussion of the implications of individual versus household-level analysis.

Likewise, our focus on mobile phone owners reflects the fact that the Novissi system in Togo distributed payments via mobile money; as such, anyone without access to a phone could not receive benefits irrespective of the targeting method – see Methods, ‘Program Exclusions’ for a discussion of exclusion errors resulting from this constraint. In practice, this constraint only affects the analysis using the 2018-2019 in-person survey, where 4,171 of 6,171 respondents provided an active phone number. For analysis using the 2020 phone survey, we include all respondents, since every respondent had access to a phone. Future work could compare phone-based targeting to counterfactual targeting methods that could be implemented in-person, and thus account for exclusion errors resulting from phone ownership.

Targeting Methods and Counterfactuals

Our evaluations use the two survey datasets to measure the performance of three targeting methods that were feasible when implementing the Novissi program: geographic blanketing (targeting everyone in certain geographies), occupation-based targeting (targeting everyone in certain occupation categories), and phone-based targeting. The location of subscribers targeted by each of these methods, in both the rural Novissi program and the hypothetical national program, are shown in Supplementary Fig. 9. Note that in the 2020 phone survey the unit of observation is the individual, while in the 2018-2019 field survey the unit of observation is the household: in practice, this means that our simulations with the 2018-2019 field survey dataset reflect a program that would provide benefits only to heads of household, and we do not account for household size in considering exclusion errors or social welfare. Future work could model phone-based targeting on a household basis by collecting phone numbers for all household members and calculating aggregate benefits assigned to each household; given survey data limitations we cannot perform this analysis.

With *geographic targeting*, the primary counterfactual approach considered by the government of Togo in implementing its rural assistance program, we assume that the program would target geographic units in order from poorest to wealthiest, and that all individuals in targeted units would be eligible for benefits. We report results from two different approaches to geographic targeting: (i) a program that targets the poorest prefectures (admin-2 region), defined as those prefectures with the lowest average predicted consumption based on a 2017 INSEED survey PMT; and (ii) a program that targets the poorest cantons (admin-3 region), defined as those cantons with the lowest average wealth based on high-resolution micro-estimates of wealth inferred from satellite imagery. When targeting the n poorest geographic regions would result in more than 29% of individual receiving benefits, then $n-1$ regions are targeted fully, and individuals from the n^{th} poorest region are selected randomly until the 29% threshold is reached. See Supplementary Fig. 5 and Methods, ‘Poverty Maps’ for the poverty maps used for

geographic targeting. (While this purely geographic approach was considered carefully by the Government of Togo, it is less common in non-emergency settings, when other data can inform targeting decisions. For instance, it is common to combine some degree of geographic targeting with community-based targeting and/or proxy means tests).

In *occupation-based targeting*, we first evaluate the effectiveness of targeting *informal workers*, which is the eligibility criteria used by Novissi when it was first launched in April 2020, and which served as the basis for paying roughly 500,000 urban residents. In practice, this process involves categorizing the occupation of every individual respondent in both surveys as either formal or informal (including unemployed), applying the same definition of informality that was used by the Novissi program. In the simulations, informal workers are targeted first (in random order if there are more informal workers than can receive benefits) and formal workers are targeted last (also in random order, if the available benefits exceed the number of informal workers).

We also develop and test a hypothetical occupation-based approach, which we refer to as *optimal occupation-based targeting*, which assumes that the policymaker had high-quality consumption data on the consumption of workers in each occupation and used that information to target the poorest occupations first. While this approach was not considered in Togo's pandemic response, it was feasible with the data sources available in Togo at the time, and represents an upper-bound on the performance of a hypothetical occupation-based targeting system. We simulate this optimal occupation-based approach by calculating the average consumption of each occupation in the 2018-2019 field survey; occupations are then targeted in order of increasing average consumption. The average consumption of each occupation category is shown in Supplementary Table 5. Note that since agricultural workers are the poorest category and make up 29% of the observations in the 2018-2019 field survey dataset and 41% of the observations in the 2020 phone survey dataset, in practice the precision and recall metrics reported in our targeting simulations reflect systems of occupation-based targeting that would prioritize agricultural workers only.

Of primary interest in the targeting evaluation is the performance of the targeting approaches based on mobile phone data. The *phone-based (ML)* approach is the one described in the main text, which uses machine learning to construct a poverty score from rich data on mobile phone use and prioritizes the individuals with the lowest poverty scores (Methods, 'Machine Learning Methods'). For reference, we also calculate the performance of a more parsimonious *phone expenditures* model, which prioritizes the individuals with the smallest total phone expenditures (Methods, 'Parsimonious Phone Expenditure Method').

For completeness, our simulations also include results from targeting methods that were not feasible for the Novissi program, as the data required to implement those methods were not available when Novissi was launched (though Togo plans to create a foundational unique ID system and comprehensive social registry in 2022).⁴² In particular, we simulate targeting using an *asset-based wealth index*, constructed as described in Methods, 'Survey Data.' For the hypothetical national simulations using the 2018-2019 field survey dataset, we also simulate targeting using a *poverty probability index* (PPI) and *proxy-means test* (PMT). Finally, when simulating targeting the hypothetical national program restricted to rural areas (Supplementary Table 1), we also simulate targeting on a *rural-specific PMT* (see Methods, Differences in Rural

and National Evaluations’). We cannot simulate PPI or PMT-based targeting using the 2020 phone survey since the necessary data were not collected.

An important caveat is that the PMT that we use in the 2018-2019 survey is “perfectly calibrated” in the sense that it is both trained and evaluated on the same sample. In real-world settings, the predictive accuracy of a PMT declines as the time increases between the time of calibration and the time of application^{22,23}. As such, the performance of the PMT we report is likely an upper bound of the performance of a real-world PMT.

For the PMT in the 2018-2019 field survey dataset, as well as for CDR-based wealth estimates in both datasets, predictions are produced out-of-sample over cross validation so that they can be fairly evaluated in targeting simulations. Specifically, in each case, the training dataset is divided into 10 cross validation folds; the machine-learning model is trained on 9 of the 10 folds and used to produce predictions for the final fold. The training-and-prediction regime is repeated for all 10 folds.

Measures of Targeting Quality

For each targeting method, we calculate two “threshold-agnostic” metrics of targeting accuracy – metrics that capture relationships between continuous measures of poverty rather than focusing on accuracy for targeting a specific portion of the population. These are:

- *Spearman correlation coefficient*: Spearman’s rank correlation coefficient is the Pearson correlation between the rank values of the true and proxy measures of poverty. We focus on the Spearman correlation rather than standard Pearson correlation as a measure of targeting quality because targeting concerns itself only with the ordering of observations according to poverty. Spearman’s correlation coefficient is calculated as follows:

$$\rho = \frac{6 \sum_{i=1}^N (r_i - \hat{r}_i)^2}{N(N^2 - 1)}$$

where N is the total number of observations, r_i is the rank of observation i according to the ground truth poverty measure, and \hat{r}_i is the rank of observation i according to the proxy poverty measure.

- *ROC curves and Area Under the Curve (AUC)*: Following Hanna & Olken (2018)³, we trace Receiver Operator Characteristic (ROC) curves that describe the quality of a targeting method at counterfactual targeting thresholds (Extended Data Figure 4, left figures). At each counterfactual targeting threshold T we simulate targeting $T\%$ of observations according to the proxy poverty measure in question and calculate the true positive rate (TPR) and false positive rate (FPR) of the classifier with respect to reaching the $T\%$ poorest according to the ground-truth poverty measure. By varying T from 0% to 100%, we construct the ROC curves shown in Extended Data Figure 4. The area under the curve (AUC) is used to summarize the targeting quality, with a random targeting method achieving an AUC of 0.5 and perfect targeting an AUC of 1. For convenience, we also include “Coverage vs. Recall” figures (right figures of Extended Data Figure 4) that show how program recall varies as the eligible percentage of the population increases.

Note that since recall is another name for the true positive rate, panels b and d represent a rescaling of the ROC curves in panels a and c.

Our analysis focuses on analyzing the performance of a quota-based approach that ranks individuals from predicted poorest to predicted wealthiest, then targets the poorest 29% of individuals. We use the quota of 29% since the rural Novissi program had sufficient funding to provide benefits to the poorest 29% of registrants in eligible cantons. (This quota-based approach is not the only way that poverty scores could be used in targeting, though it is the only approach that we evaluate. For instance, a threshold-based approach might target everyone below a threshold poverty score; alternative approaches might provide cash transfers of different sizes depending on the poverty score of the beneficiary⁴³). The 29th percentile corresponds to a consumption threshold of USD 1.17/day in the 2018-2019 field survey dataset, and a PMT threshold of USD 1.18/day in the 2020 phone survey dataset. We calculate the following metrics to describe how accurately targeting the poorest 29% according to each targeting method reaches (1) the 29% truly poorest, (2) those below the international poverty line of USD 1.90/day (57% of observations in the 2018-2019 field survey, and 76% of observations in the 2020 phone survey), and (3) those below the extreme poverty line, which was defined as three-quarters of the poverty line, or USD 1.43/day (41% of observations in the 2018-2019 field survey, and 53% of observations in the 2020 phone survey):

- *Accuracy*: Classification accuracy measures the proportion of observations that are identified correctly (targeted observations that are poor according to the ground-truth poverty measure, and non-targeted observations that are not poor according to the ground-truth wealth measure). $Accuracy = \frac{TP+TN}{TP+FP+TN+FN}$.
- *Recall*: Recall measures the proportion of all poor observations that are reached by a given targeting method. $Recall = \frac{TP}{TP+FN}$. Recall is closely related to the concept of exclusion errors (i.e., the fraction of true poor who do not receive benefits, $\frac{FN}{TP+FN}$), since $Recall = 1 - Exclusion\ error$.
- *Precision*: Precision measures the proportion of targeted observations that are poor according to the ground-truth poverty measure. $Precision = \frac{TP}{TP+FP}$. Precision is closely related to the concept of inclusion errors (i.e., the fraction beneficiaries who are non-poor, $\frac{FP}{TP+FP}$), since $Precision = 1 - Inclusion\ error$.
- *Exclusion error*: The proportion of true poor excluded from benefits. Defined as $\frac{FN}{TP+FN}$.
- *Inclusion error*: The proportion of beneficiaries who are not poor, i.e., $\frac{FP}{TP+FP}$.

Note that the poverty lines are applied to consumption expenditure in the 2018-2019 field survey dataset, and to the proxy-means test estimates in the 2020 phone survey dataset.

Differences in Rural and National Evaluations

The results in Table 1 indicate that the phone-based targeting approach – as well as the counterfactual targeting approaches – was more effective in the actual rural Novissi program

(first four columns of Table 1) than it would have been in a hypothetical nationwide program (last four columns of Table 1). There are several factors that may account for these differences. Some of these factors are difficult for us to test empirically, for instance the fact that the surveys were conducted at different points in time, used different teams of enumerators, and different data collection modalities (phone vs. in person). We investigate two factors that we can explore empirically: the geographic concentration of each survey and the ground truth measure of poverty (consumption vs. PMT). We additionally explore whether targeting results are sensitive to the use of a nationwide PMT vs. a rural-specific PMT.

Geographic concentration: Whereas the rural Novissi evaluation focuses on Togo's 100 poorest cantons, the hypothetical national program is evaluated nationwide (397 cantons). We therefore present results in Supplementary Table 1 that restrict the simulation of the hypothetical national program to the 2,306 households in rural areas (out of 4,171 total). Comparing the results in Supplementary Table 1 to the last four columns of Table 1, we find that the performance of all methods drops, as would be expected when the beneficiary population is more homogeneous. The relative difficulty of estimating poverty among rural populations is also evident in Extended Data Table 6: the CDR-based method's performance at predicting both consumption and the PMT is lower when the analysis of the 2018-19 survey is restricted to the rural population (Panel A vs. Panel B). Importantly, we also observe that the relative performance of phone-based targeting increases: whereas the CDR-based method performed worse than the asset index and only slightly better than canton-based targeting in the full nationwide evaluation (last four columns of Table 1), the CDR-based method is on par with the asset index and substantially better than canton-based targeting when the nationwide survey is limited to rural areas (Supplementary Table 1).

Consumption vs. PMT: Whereas the national evaluation uses a measure of consumption as ground truth, the rural Novissi evaluation uses a PMT as ground truth. Supplementary Table 7 therefore simulates the hypothetical national program using a PMT as ground truth. Comparing the results in Supplementary Table 7 to the last four columns in Table 1, we find that using a PMT rather than consumption as ground truth increases targeting accuracy across all of the targeting methods. However, switching from consumption to the PMT does not substantially improve the performance of the phone-based method relative to the counterfactual approaches. This latter finding suggests that the use of the PMT is likely not a major source of the difference between the relative performance of the CDR-based method in the rural Novissi program (first four columns of Table 1) and the hypothetical nationwide program (last four columns of Table 1).

National PMT vs. Rural PMT: Since the best predictors of welfare differ for rural and urban populations, we explore whether targeting results change when the PMT is calibrated using a rural rather than national population. Specifically, we construct a rural-specific PMT using the same methodology described in Methods, 'Survey Data,' but restricting the training data to observations in the 2018-2019 field survey that are in rural areas. This rural PMT explains 17% of the variation in log-transformed consumption in rural areas, and is highly correlated (Pearson = 0.75) with the general PMT. We then produce rural PMT estimates for respondents to the 2020 phone survey, and retrain the phone-based poverty prediction model to predict the rural-specific PMT in that population. Supplementary Table 8 then presents results from simulating with the rural PMT as ground truth. Comparing Supplementary Table 8 to the first four columns of Table

1, we observe a noticeable improvement in the performance of the asset index, but other results are largely unchanged.

Relatedly, Extended Data Table 3 shows the feature importances for different phone-based prediction models. Panels A and B show the top-10 features for the main models presented in Table 1, i.e., for predicting a PMT in the 2020 rural phone survey, and predicting consumption in the 2018-19 nationwide household survey. Panels C and D shows the top-10 features for predicting a PMT in the 2018-19 survey, and predicting a PMT in the 2018-2019 household survey, restricted to rural areas. The feature importances for the two national-scale models are similar, suggesting the role of the ground truth poverty measure may not be as important as the role of geography in creating the poverty prediction models. The feature importances for the two rural-focused models are less similar, which may be due to the fact that the 2020 phone survey is concentrated in the 100 poorest cantons, while in Panel D we restrict to rural areas, but these rural areas still cover the entire country.

Taken together, the results in this subsection suggest that the benefits of phone-based targeting are likely to be greatest when the population under consideration is more homogeneous, and when there is less variation in other factors (such as place of residence) that are used in more traditional approaches to targeting.

Location-Based Targeting

Several results emphasize the importance of geographic information in effective targeting. In particular, we observe that basic geographic targeting performs nearly as well as phone-based targeting in specific simulations – in particular, in simulations of a nationwide program that can afford to target a large proportion of the total population (e.g., Extended Data Table 2). We also found that location-related features from the CDR are important in the phone-based prediction model (Methods, ‘Machine Learning Methods’).

For these reasons, Supplementary Table 9 explores the extent to which targeting could be based on a *CDR-location model* that only uses the CDR to infer an individual’s home location (see Supplementary Methods section 4). As with the phone (expenditures) model, the CDR-location model may be attractive to implementers since the data and technical requirements are reduced⁴⁴. In Supplementary Table 9, we observe that geographic targeting using phone-inferred home location is of slightly lower quality than geographic targeting using survey-recorded home location, and substantially worse than targeting using the machine learning approach.

We also investigate the correlation between different sources of information on an individual’s location. Supplementary Table 10 compares three different methods for identifying an individual’s location, using roughly 4,500 respondents to the 2020 phone survey. At the prefecture (admin-2) level, most people (90%) self-declare living in the same canton in which they are registered to vote; there is also strong overlap between the individual’s CDR-inferred location and self-declared location (70%). The accuracy is substantially lower at the canton level, which is likely due to error in the CDR-inference algorithm when spatial units are small, as well as to confusion among respondents as to which canton they live in (e.g., most respondents were confident in naming their village, but did not always know their canton).

Supplementary Table 11 presents additional analysis to compare the mobile phone activity of each subscriber with their home location, as recorded in the survey and as inferred from their CDR. We find that 62-85% of the average subscriber's activity occurs in their home prefecture, and that all of the modal subscriber's activity occurs in their home prefecture. These results are consistent with the importance of location-related features in the prediction algorithm (and the relatively low mobility of the rural Togolese population).

This analysis may also provide some context for the difference in the accuracy of the geographic targeting methods between the rural evaluation and the national evaluation in Table 1. While canton-based targeting performs better in the national evaluation, which is consistent with past work showing that finer-resolution geographic targeting is preferred to lower-resolution geographic targeting^{45,46}, prefecture-based targeting counter-intuitively performs better in the rural evaluation. We suspect this discrepancy is caused by three main factors: First, we expect that the estimates of average canton wealth are likely to be noisier than the estimates of average prefecture wealth, since the prefecture estimates aggregate over a larger population and the canton estimates rely on satellite-based inferences. Second, in the rural evaluation the prefecture is an important component of the PMT that is used as the ground truth measure of poverty (see Supplementary Table 3), so prefecture targeting relies on information that is structurally incorporated into the ground truth outcome (unlike in the national evaluation, where the ground truth outcome is consumption). The results in Supplementary Table 7 are consistent with this second hypothesis: the gap between prefecture and canton targeting in the national evaluation in Table 1 is smaller when switching the ground-truth poverty outcome from consumption to the PMT. Third, locations in the rural phone survey were self-reported, whereas locations were recorded on GPS devices by enumerators in the national survey; as noted, many respondents expressed confusion about their home canton. (The results in Supplementary Table 9, however, are not consistent with this third hypothesis: they indicate that targeting on canton inferred from mobile phone data is weaker than targeting on prefecture inferred from mobile phone data, suggesting that a difference in response quality between prefecture and canton in the survey is not a major factor in the difference in outcomes in the targeting simulations).

Temporal Stability of Results

When simulating the performance of phone-based targeting, our main analysis uses each survey dataset to both train the machine learning model and, via cross-validation, to evaluate its performance. These measures of targeting performance thus indicate what should be expected when training data (i.e., the ground truth measures of poverty and the matched CDR) are collected immediately prior to a program's deployment. This best-case scenario is what occurred in Togo in 2020: the phone survey was completed in October 2020 and Novissi was expanded beginning in November 2020. In other settings, however, it may not be possible to conduct a survey before launching a new program; it may likewise not be possible to access up-to-date mobile phone data.

To provide an indication of how long phone-based models and predictions remain accurate, Extended Data Table 4 compares (i) the best-case scenario to alternative regimes where (ii) the training data are old but the CDR are current, and (iii) the training data are old and the CDR are also old. In these simulations, the "old" data are from the 2018-19 national household survey and corresponding 2019 phone dataset; the "current" data are the subset of 2020 phone survey respondents for whom CDR are available in 2019 and 2020 ($N = 7,064$). In all simulations, the

2020 PMT is used as the ground truth measure of poverty. Predictions for (i) are generated over 10-fold cross validation; predictions for (ii) and (iii) are out-of-sample with respect to the training data, since the models are trained on the 2018-19 field survey. (An additional issue with (iii) is turnover on the mobile phone network: 1,851 (21%) of phone numbers collected in the 2020 survey were not on the mobile phone network in 2019, and therefore cannot be associated with a wealth prediction in (iii). See also Supplementary Fig. 6 for detailed information on rates of turnover on the mobile phone network).

The results in Extended Data Table 4 indicate that predictive performance decreases when the model is out of date, and decreases even further when the CDR are out of date. This is to be expected, since roughly two years elapsed between the “old” and “current” periods: in addition to changes in how people use their phones (which would disrupt the accuracy of the predictive model), the actual economic status of some individuals may have changed – for instance, due to the COVID-19 pandemic. There are also other important differences between the 2018-19 national household survey and the 2020 phone survey that could affect the extent to which a model trained on the former could accurately predict outcomes in the latter (such as the mode of data collection, the geographic concentration of the sample, and so forth – see Methods, ‘Differences in Rural and National Evaluations’).

For the main simulations focused on reaching the poorest 29%, Extended Data Table 4 suggests that accuracy decreases by 3-4 percentage points (4-6%) and precision decreases by 5-7 percentage points (10-14%) when out of date models and CDR are used for targeting. These losses are nearly as large as the gains of phone-based targeting over geographic targeting observed in Table 1, which emphasizes the importance of having current and representative training data for real-world deployment of phone-based targeting. However, in absolute levels, the phone-based predictions remain reasonably accurate despite the two-year gap between the training and test environments (i.e., the Spearman correlation with ground truth is $\rho = 0.35 - 0.36$).

Social Welfare

Using the two matched survey-CDR datasets, we calculate aggregate utility under each of the targeting methods using a social welfare function. Following Hanna & Olken (2018)³ we rely on constant relative risk-aversion (CRRA) utility, which models individual utility as a function of pre-transfer consumption and transfer size:

$$U = \frac{\sum_{i=0}^N (y_i + b_i)^{1-\rho}}{1-\rho}$$

Where N is the population size, y_i is the consumption of individual i , and b_i the benefits assigned to the individual. Following Hanna & Olken (2018)³, we use a coefficient of relative risk-aversion $\rho = 3$. To reflect the policy design of the Novissi program, we assume that all beneficiaries who receive a benefit receive the same value $b_i = b$. (In principle, the benefit b_i paid to i could depend on characteristics of i , such as i 's level of poverty. While such an approach would substantially increase total welfare, in practice it is much more difficult to implement). To construct the social welfare curves, we:

- Calculate a total budget available for each of the two datasets. We focus on programs that have a budget size analogous to that of rural Novissi, which aimed to distributed approximately USD 4 million among the 154,238 program registrants, or USD 25.93 per registrant. We therefore assign each dataset a total budget of USD 25.93 N , where N is the total size of the dataset.
- Simulate targeting $T\%$ of observations on the basis of each of our counterfactual targeting approaches.
- Assign equal benefits to each of the targeted observations, with the budget divided evenly among targeted observations (so lower targeting thresholds T correspond to more benefits for targeted individuals).
- Calculate aggregate utility by summing over benefits and consumption for each individual with the CRRA utility function. Note that non-targeted individuals are included in the welfare calculation; they are merely assigned 0 benefits. For the 2018-2019 field survey dataset we use consumption expenditure for y_i ; for the 2020 phone survey dataset we use the PMT estimates.
- By varying T between 0% and 100% of observations targeted, we trace out the social welfare curves shown in Figure 2.

Fairness

We are interested in auditing our targeting methods for fairness across sensitive subgroups. Note that notions of parity and fairness are debated in machine learning and policy communities: Barocas et al. (2018)⁴⁷ describe how the three most popular parity criteria --- demographic parity (benefits assigned to subgroups proportionally to their size), threshold parity (use of the same classification threshold for all subgroups), and error rate parity (equal classification error across subgroups) --- are in tension with one another. Moreover, Noriega et al. (2020)²⁶ describe how tensions over parity criteria, prioritized subgroups, and positive discrimination lead to complicated prioritization compromises in the administration of targeted social protection programs.

Here we focus on two targeting-specific parity criteria:

- *Demographic parity*: A targeting method satisfying demographic parity will assign benefits to a subgroup proportionally to the subgroup’s presence in the population of interest. We evaluate demographic parity among the poor: that is, we compare the proportion of each subgroup living in poverty (below the 29th percentile in terms of consumption) to the proportion of each subgroup that is targeted (below the 29th percentile in terms of the proxy poverty measure used for targeting).

$$DP = \frac{\text{True Positives} + \text{False Positives}}{N} - \frac{\text{True Positives} + \text{False Negatives}}{N}$$

- *Normalized rank residual*: We are interested in whether certain subgroups are consistently ranked higher or consistently ranked lower than they “should” be by the counterfactual targeting approaches. We therefore compare the distributions of rank residuals across subgroups and targeting methods:

$$RR_i = \frac{\hat{r}_i - r_i}{N}$$

where \hat{r}_i is the poverty rank of individual i according to the proxy poverty measure and r_i is the poverty rank of individual i according to the ground-truth poverty measure.

We focus on seven dimensions for parity: gender, ethnicity, religion, age group, disability status, number of children, and marital status. We also evaluate parity across whether an individual is “vulnerable,” where vulnerability is defined as one of the following traits: {female, over age 60, has a disability, has more than five children, is single}. We conduct this analysis using demographic information about the head of the household in the 2018-2019 field survey dataset, as these demographic variables were not all collected in the 2020 phone survey.

PROGRAM EXCLUSIONS

In Table 2, we present information on sources of exclusion from the Novissi program that are not inherently related to targeting. These estimates are drawn from diverse sources of administrative and survey data, specifically:

- *Voter ID penetration:* According to government administrative datasets, 3,633,898 individuals were registered to vote in Togo by late 2019. The electoral commission of Togo reports that this corresponds to 86.6% of eligible adults. While the total adult population in Togo is hard to pin down (the last census was in 2011), Togo’s national statistical agency (<https://inseed.tg/>) estimates that there are 3,715,318 adults in Togo, whereas the United Nations estimates 4.4 million adults in Togo⁴⁸, implying a voter ID penetration rates of 82.6% or 97.8%.
- *Phone penetration:* In the 2018-2019 field survey, 65% of individuals reported owning a mobile phone (Supplementary Figure 3a) and 85% of households included at least one individual who owns a phone (Supplementary Figure 3b). In rural areas, these rates drop to 50% of individuals and 77% of households. Rates of phone ownership are significantly lower among women (53%) than among men (79%), especially in rural areas (33% for women and 71% for men). These household survey-based estimates likely represent a lower bound, given the steady increase in phone penetration between 2018 and 2020. The Togolese government estimates 82% SIM card penetration in the country (though some people may have multiple SIM cards)⁴⁹. Based on data from the mobile phone companies, we observe 5.83 million unique active SIMs in Togo between March and September 2020.
- *Past phone use:* In order to construct a phone-based poverty estimate for a subscriber, they had to place at least one outgoing call or text on the mobile phone network in the period of mobile network observation prior to the program’s launch (March – September 2020, with program registrations in November-December 2020). In Togo, a lower bound on this source of exclusion is the typical monthly rate of mobile phone turnover, which we estimate to be roughly 2.5% (see Supplementary Figure 6). An upper bound is closer to 27%, which is the number of SIM cards that registered for Novissi November-December 2020 who did not make an outgoing transaction in the March-September. This discrepancy may be due to (i) individuals buying new SIM cards specifically to register

for Novissi; or (ii) individuals registering for Novissi using existing SIM cards that were not in active use, for instance the SIM cards in multi-SIM phones. Based on qualitative observation, multi-SIM phones are very common in Togo, and secondary or tertiary SIMs are infrequently used (or not used at all). It is possible that families registered one household member on a primary SIM and others on secondary or tertiary SIMs that may have had no previous network activity.

- *Program awareness:* Since individuals had to register for the Novissi program to receive benefits, program advertising and population awareness was a key goal. The program was advertised via radio, SMS, field teams, and direct communication with community leaders at the prefecture and canton level. In total, 245,454 subscribers attempted to register for the program. Although we do not observe the prefecture and canton of subscribers who attempt but do not succeed in registering in our administrative data, we know that 87% of successful registrants are in cantons eligible for benefits. Assuming the rate is approximately the same for attempters, we expect that around 213,545 of the attempters are in eligible cantons. The total voting population in eligible cantons is 528,562, for an estimated attempted registration rate of 40.40%.
- *Registration challenges:* Registration for the Novissi program required the completion of a short (5 question) USSD survey. Of the 245,454 subscribers that attempted to register for the program, 176,517 succeed, for a 71.91% rate of registration success.

Overlaps among sources of exclusion: The above sources of exclusion are not independent and are therefore not cumulative. For instance, individuals who are not registered to vote may also be systematically less likely to have a mobile phone. For this reason, Extended Data Table 5 uses the 2020 phone survey dataset – restricted to respondents who report living in an eligible canton – to calculate overlaps in sources of exclusion to the poor, including Voter ID possession, program awareness, registration challenges, and targeting errors using the phone-based targeting method. We cannot account for mobile phone ownership in this analysis since the 2020 survey was conducted over the phone, and sampled based on past CDR (see Supplementary Methods section 5).

The final three columns of Extended Data Table 5 show, based on the 2020 phone survey dataset, average characteristics of the population “succeeding” at each step: average PMT, percent women, and average age. The first panel shows successive exclusions for the entire population; the second panel focuses on just the poorest 29% (i.e., those who “should” be receiving aid, were everyone to register for the program and were the targeting algorithm perfect). In Panel A, we observe that to a certain extent the “right” types of people are dropping out at each step, which would be consistent with self-targeting observed in other contexts³¹: in particular, those who attempt to register are poorer than the overall population (average PMT = 1.45 vs. 1.62). There are little differences in the share of the successful population who are women or average age, except in the targeting stage.

Comparing Panels A and B of Extended Data Table 5, we observe that the recall of the targeting algorithm is substantially higher among the population that owns a voter ID and succeeds in registration for the program (61%, as shown in Extended Data Table 5, last row) than the overall population surveyed in the 2020 phone survey (47%, as shown in Table 1, row 4). This may be due to self-selection (i.e., the type of poor people who register for Novissi tend to also have low

phone-based poverty scores). However, it could alternatively suggest that the phone-based targeting algorithm is best at identifying the poor among the types of subscribers who are aware of and register to the Novissi program.

References (Methods)

21. McBride, L. & Nichols, A. Retooling Poverty Targeting Using Out-of-Sample Validation and Machine Learning. *World Bank Economic Review* **32**, 531–550 (2018).
22. Brown, C., Ravallion, M. & van de Walle, D. A poor means test? Econometric targeting in Africa. *Journal of Development Economics* **134**, 109–124 (2018).
23. Skoufias, E., Diamond, A., Vinha, K., Gill, M. & Dellepiane, M. R. Estimating poverty rates in subnational populations of interest: An assessment of the Simple Poverty Scorecard. *World Development* **129**, 104887 (2020).
24. Eubanks, V. *Automating inequality: How high-tech tools profile, police, and punish the poor*. (St. Martin's Press, 2018).
25. Barocas, S., Hardt, M. & Narayanan, A. *Fairness and Machine Learning*. (fairmlbook.org, 2018).
26. Noriega-Campero, A. *et al.* Algorithmic targeting of social policies: fairness, accuracy, and distributed governance. in *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* 241–251 (ACM, 2020). doi:10.1145/3351095.3375784.
27. Kleinberg, J., Mullainathan, S. & Raghavan, M. Inherent Trade-Offs in the Fair Determination of Risk Scores. *arXiv:1609.05807 [cs, stat]* (2016).
28. Dwork, C., Hardt, M., Pitassi, T., Reingold, O. & Zemel, R. Fairness Through Awareness. *arXiv:1104.3913 [cs]* (2011).
29. Skoufias, E. Economic Crises and Natural Disasters: Coping Strategies and Policy Implications. *World Development* **31**, 1087–1102 (2003).
30. UNDP Regional Innovation Centre. A COVID cash transfer programme that gives more money to women in Togo. *The Innovation Dividend Podcast, EP 9* <https://undp-ric.medium.com/cina-lawson-a-covid-cash-transfer-programme-that-gives-more-money-to-women-in-togo-2386c5dff49> (2020). Accessed Jan. 1, 2022.
31. Alatas, V. *et al.* Self-Targeting: Evidence from a Field Experiment in Indonesia. *Journal of Political Economy* **124**, 371–427 (2016).
32. Gunnemann, J. *PMT Based Targeting in Burkina Faso*. <https://openknowledge.worldbank.org/handle/10986/33467> (2016) doi:10.1596/33467. Accessed Jan. 1, 2022.
33. Grosh, M. E. & Baker, J. L. *Proxy means tests for targeting social programs*. (The World Bank, 1995). doi:10.1596/0-8213-3313-5.
34. Tiecke, T. G. *et al.* Mapping the world population one building at a time. *arXiv:1712.05839 [cs]* (2017).
35. De Montjoye, Y.-A., Rocher, L. & Pentland, A. S. bandicoot: a python toolbox for mobile phone metadata. *J. Mach. Learn. Res.* **17**, 6100–6104 (2016).
36. de Montjoye, Y.-A., Hidalgo, C. A., Verleysen, M. & Blondel, V. D. Unique in the Crowd: The privacy bounds of human mobility. *Scientific Reports* **3**, 1376 (2013).
37. Cecaj, A., Mamei, M. & Biccocchi, N. Re-identification of anonymized CDR datasets using social network data. in *2014 IEEE International Conference on Pervasive Computing and Communication Workshops (PERCOM WORKSHOPS)* 237–242 (2014). doi:10.1109/PerComW.2014.6815210.
38. Alaggan, M., Gambs, S., Matwin, S. & Tuhin, M. Sanitization of Call Detail Records via Differentially-Private Bloom Filters. in *Data and Applications Security and Privacy XXIX* (ed. Samarati, P.) vol. 9149 223–230 (Springer International Publishing, 2015).

39. Mir, D., Isaacman, S., Caceres, R., Martonosi, M. & Wright, R. DP-WHERE: Differentially private modeling of human mobility. in 580–588 (2013). doi:10.1109/BigData.2013.6691626.
40. Steele, J. E. *et al.* Mapping poverty using mobile phone and satellite data. *Journal of The Royal Society Interface* **14**, 20160690 (2017).
41. Pokhriyal, N. & Jacques, D. C. Combining disparate data sources for improved poverty prediction and mapping. *PNAS* **114**, E9783–E9792 (2017).
42. Togo First. Togolese deputies approve biometric ID project. *Togo First* <https://www.togofirst.com/en/public-management/0409-6177-togolese-deputies-approve-biometric-id-project>.
43. Lindert, K., Karippacheril, T. G., Caillava, I. R. & Chávez, K. N. *Sourcebook on the Foundations of Social Protection Delivery Systems*. (World Bank Publications, 2020).
44. Kabemba, P. B., Laura Bermeo, and François. Cash and the city: Digital COVID-19 social response in Kinshasa. *Brookings Institute: Future Development Series* <https://www.brookings.edu/blog/future-development/2021/09/08/cash-and-the-city-digital-covid-19-social-response-in-kinshasa/> (2021).
45. Baker, J. L. & Grosh, M. E. Poverty reduction through geographic targeting: How well does it work? *World Development* **22**, 983–995 (1994).
46. Smythe, I. & Blumenstock, J. E. Geographic Micro-Targeting of Social Assistance with High-Resolution Poverty Maps. *arXiv:2004.03865* (2020).
47. Barocas, S., Hardt, M. & Narayanan, A. Fairness and Machine Learning Limitations and Opportunities. <https://www.semanticscholar.org/paper/Fairness-and-Machine-Learning-Limitations-and-Barocas-Hardt/bae7f0b3448a3eac77886f2a683c0cf9256bb8bf> (2018).
48. United Nations. 2019 Revision of World Population Prospects. *United Nations Department of Economic and Social Affairs* (2019).
49. République togolaise. Programme digital de transferts monétaires en réponse à la COVID-19. <http://www.fondation-farm.org/zoe/doc/colocnovissi.pdf> (2020). Accessed Jan. 1, 2022.
50. République togolaise. 86,6% des Togolais ont une carte d'électeur. *RepublicOfTogo.com* <https://www.republicoftogo.com/toutes-les-rubriques/politique/86-6-des-togolais-ont-une-carte-d-electeur>. Accessed Jan. 1, 2022.

Acknowledgements

Isabel Onate Falomir, Shikhar Mehra, Suraj Nair, Adrian Dar Serapio, Nathaniel Ver Steeg, and Rachel Warren provided invaluable research assistance on this project. This project would not have been possible without the dedication of our project partners in Togo, especially Minister Cina Lawson, Shegun Bakari, Leslie Mills, Kafui Ekouhoho, Morlé Koudeka, and Attia Byll. The data collection for the 2020 phone survey was made possible by the hard work of a dedicated team of enumerators from Institut National de la Statistique et des Études Economiques et Démographiques (INSEED), led by Stanislas Telou and Ruth Ogoumedi. The team at GiveDirectly was instrumental in implementing the Novissi expansion studied in this paper (especially Han Sheng Chia, Michael Cooke, Kristen Lee, Alex Nawar, and Daniel Quinn). We thank Esther Duflo, Luis Encinas, Tina George, Rema Hanna, Ethan Ligon, and Ben Olken for helpful feedback. We are grateful for financial support from Google.org, data.org, the Center for Effective Global Action, the Jameel Poverty Action Lab, the Global Poverty Research Lab at Northwestern University, and the World Bank, which financed the phone surveys and data collection under the WURI program. Blumenstock is supported by NSF award IIS – 1942702. Authors retained full intellectual freedom in conducting this research, and as such all opinions and errors are our own.

Author Contributions

E.A. and J.E.B. jointly supervised this project. E.A., S.B., and J.E.B analyzed data and prepared figures. E.A. and J.E.B. wrote the paper. All authors edited and revised the manuscript.

Competing interests

The authors declare no competing interests.

Additional Information

Supplementary information is available for this paper.

Correspondence and requests for materials should be addressed to jblumenstock@berkeley.edu.

Reprints and permissions information is available at www.nature.com/reprints.

Data Availability

The data used in this analysis include data that are available from public online repositories, data that are available upon request of the data provider, and data that are not publicly available because of restrictions by the data provider. The micro-estimates of wealth and population density used to derive satellite-based poverty maps are available from the Humanitarian Data Exchange (<https://data.humdata.org/dataset/relative-wealth-index> and <https://data.humdata.org/dataset/highresolutionpopulationdensitymaps-tgo>). The survey datasets are available upon request from the Institut National de la Statistique et des Études Economiques et Démographiques (<https://inseed.tg/> and inseed@inseed.tg). The mobile phone data and administrative data from the Novissi program contain proprietary and sensitive information, and

cannot be publicly released. Upon reasonable request, we can provide information to academic researchers on how to contact mobile network operators and the Togolese government to request these datasets.

Code Availability

The code used for these analyses is publicly available at the GitHub repository located at <https://github.com/emilylaiken/togo-targeting-replication/>.

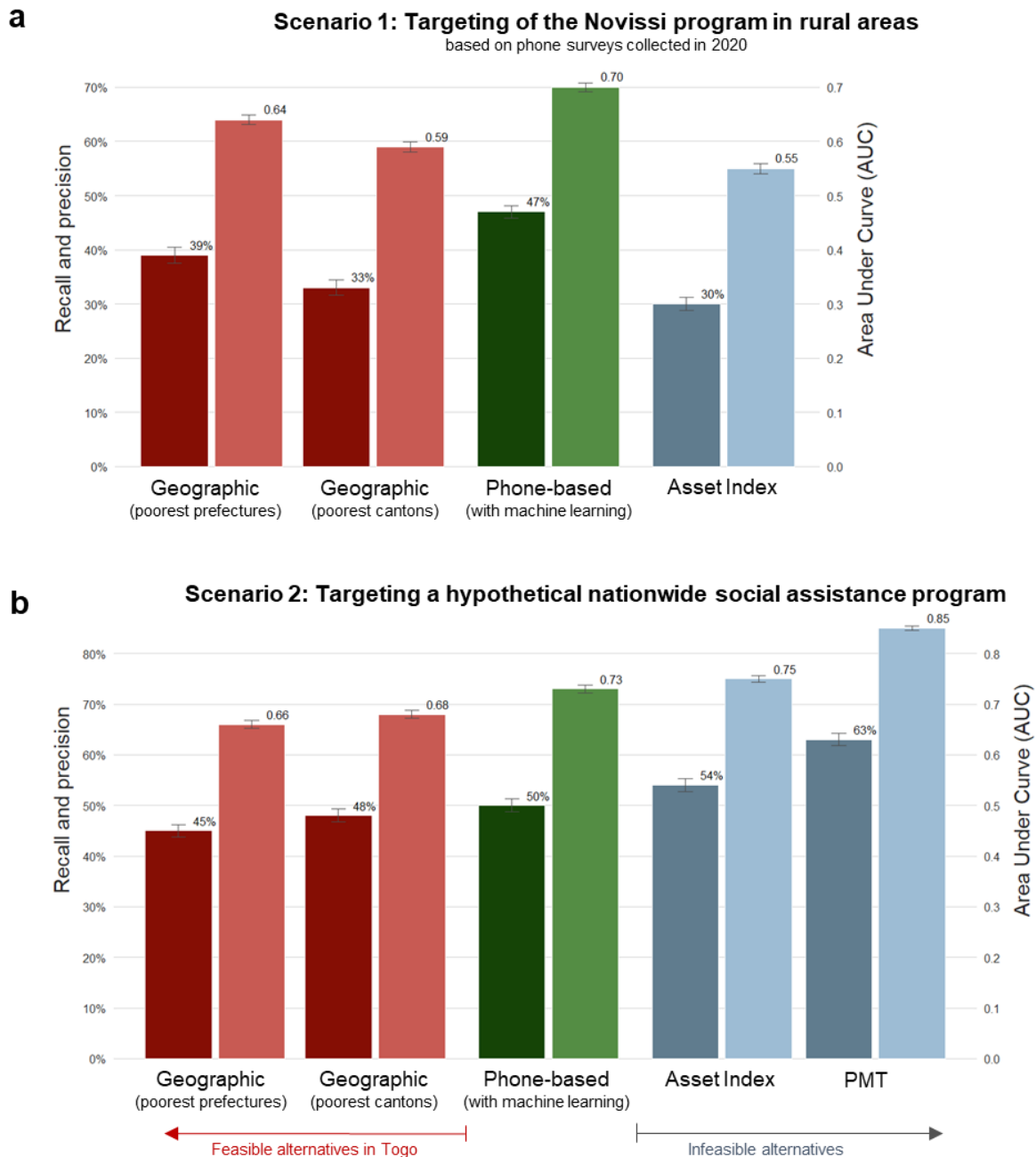
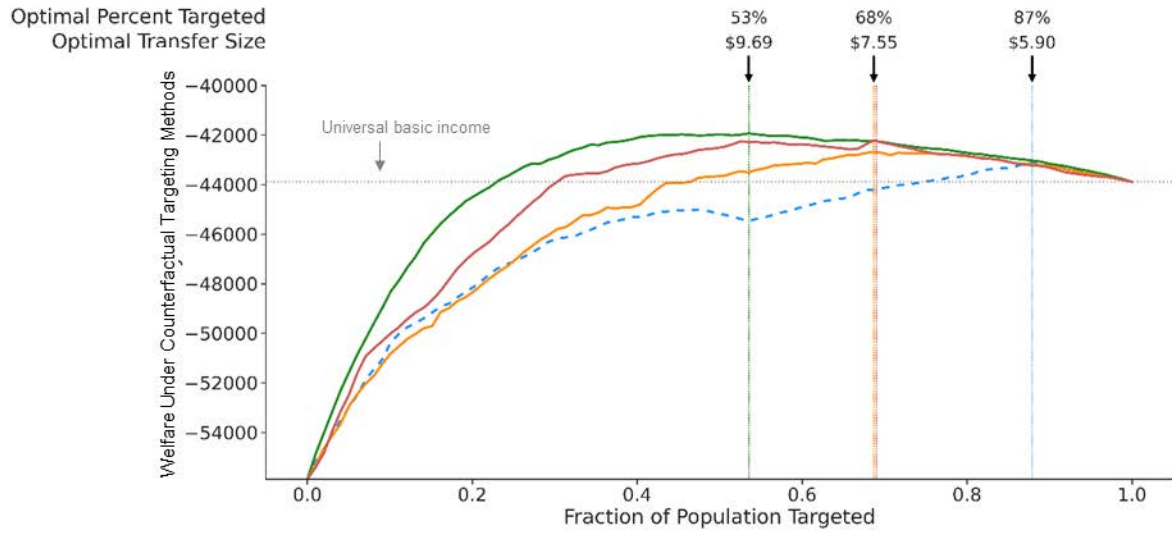


Figure 1 | Comparing Novissi targeting to alternatives. The performance of phone-based targeting (green) in comparison to alternative approaches that were feasible (red) and infeasible (grey) in Togo in 2020. Targeting is evaluated for **(a)** The actual rural Novissi program, which focused on Togo’s 100 poorest cantons (using a 2020 survey representative of mobile subscribers in the 100 cantons, where PMT is a ground truth for poverty since consumption data was not collected in the phone survey); and **(b)** a hypothetical nationwide anti-poverty program (using a national field survey conducted in 2018-2019, where consumption is a ground truth for poverty). Darker bars indicate recall and precision (left axis), which is equivalent to $1 - \text{exclusion error}$; lighter bars indicate Area Under Curve (right axis). The bar height represents the point estimate from the full simulation; standard deviations produced from $N=1,000$ bootstrap simulations are shown as whiskers. This figure highlights a subset of the results contained in Table 1.

a

Scenario 1: Targeting of the Novissi program in rural areas
based on phone surveys collected in 2020

**b**

Scenario 2: Targeting a hypothetical nationwide social assistance program
based on in-person surveys collected in 2018-2019

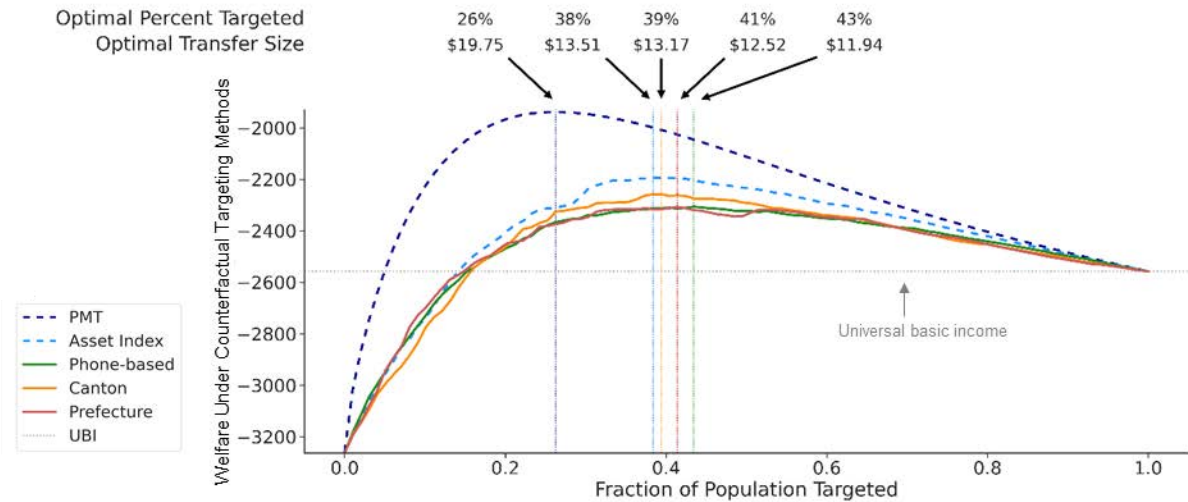


Figure 2 | Welfare analysis of different targeting mechanisms. Aggregate social welfare is calculated (assuming CRRA utility) under counterfactual targeting approaches. We assume a fixed budget of USD 4 million and a population of 154,238, with an equal transfer size for all beneficiaries. Utility curves for feasible targeting mechanisms are shown in solid lines; infeasible targeting mechanisms are shown in dashed lines. The horizontal dotted line indicates total social welfare for a universal basic income program that provides (very small) transfers to the entire population; vertical dotted lines indicate the targeting threshold and associated transfer size that maximizes social welfare for each targeting mechanism. Targeting is evaluated for **(a)** an anti-poverty program in Togo’s 100 poorest cantons; and **(b)** a hypothetical nationwide anti-poverty program.

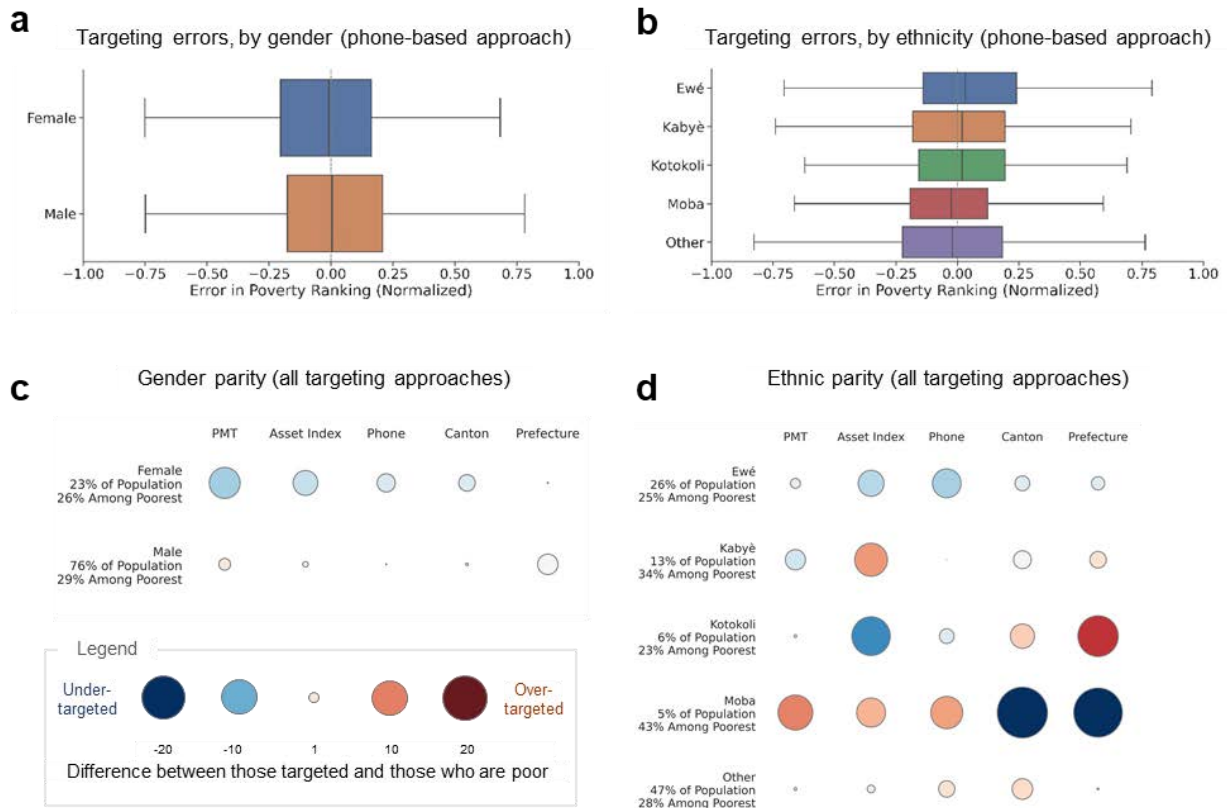


Figure 3 | Fairness of targeting for different demographic subgroups. Above: Distributions of differences between ranking according to predicted wealth from the ML approach and ranking according to true wealth (using the 2018-2019 field survey matched to CDR, $N = 4,171$), disaggregated by gender (a) and ethnicity (b). Boxes show the 25th and 75th percentiles, whiskers show the minimum and maximum, and the center line shows the median of the distribution. Left-skewed bars indicate groups that are consistently under-ranked; right-skewed bars indicate groups that are consistently over-ranked. Below: Evaluation of demographic parity across subgroups by comparing the proportion of a subgroup targeted under counterfactual approaches to the proportion of the subgroup that falls into the poorest 29% of the population (using data from the 2018-2019 field survey matched to CDR, $N = 4,171$). Bubbles show the percentage point difference between the proportion of the subgroup that is targeted and the proportion that is poor according to ground-truth data. Large red bubbles indicate groups that are over-targeted; large blue bubbles indicate groups that are under-targeted.

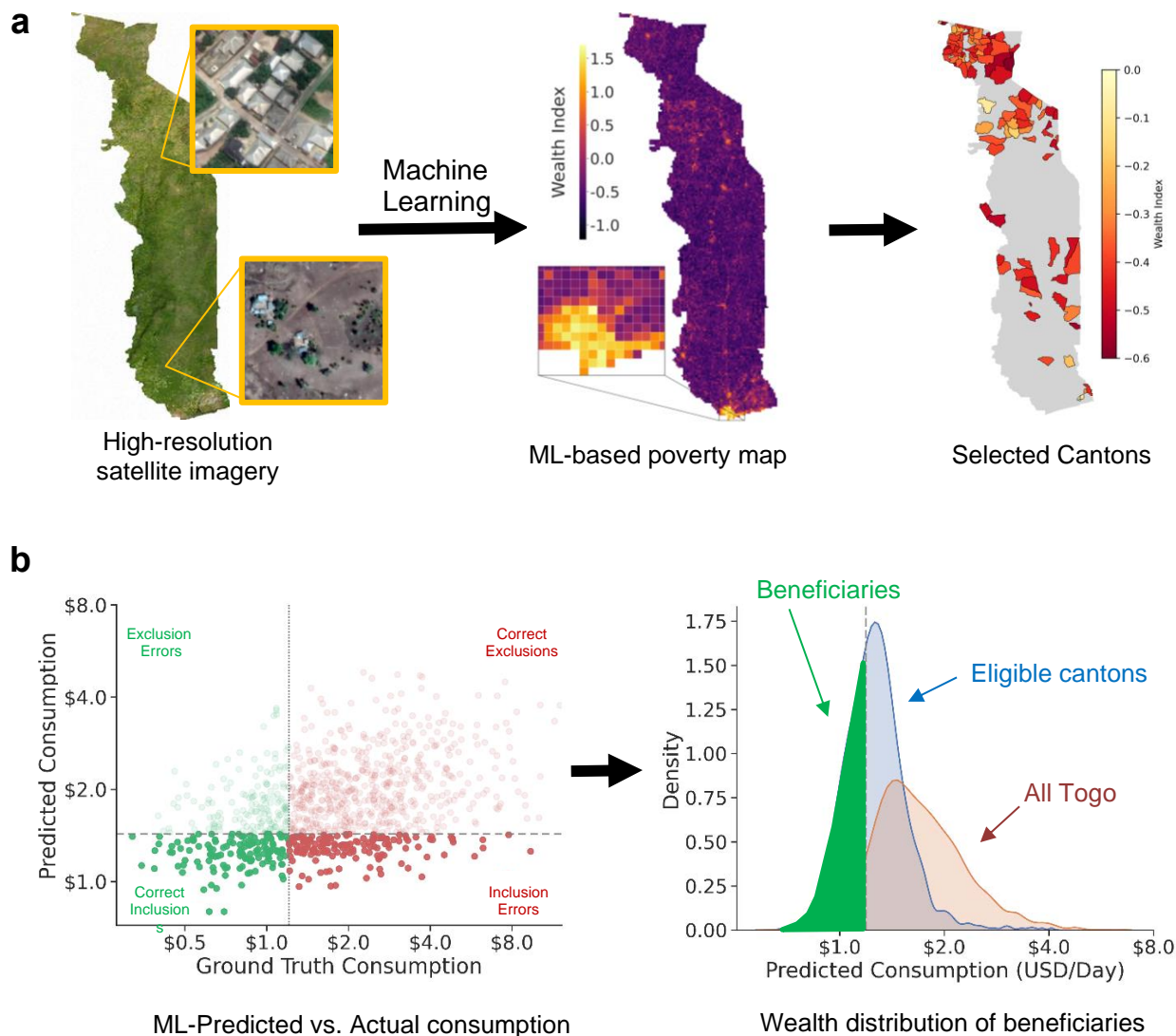
Targeting Novissi in rural Togo Based on 2020 Phone Survey ($N = 8,915$)					Hypothetical nationwide program Based on 2018-2019 Field Survey ($N = 4,171$)			
	Spearman	AUC	Accuracy	Precision & Recall	Spearman	AUC	Accuracy	Precision & Recall
<i>Panel A: Targeting methods considered by the Government of Togo in 2020</i>								
Prefecture (Admin-2 regions)	0.30 (0.017)	0.64 (0.008)	65% (0.87%)	39% (1.51%)	0.34 (0.017)	0.66 (0.008)	68% (0.74%)	45% (1.27%)
Canton (Admin-3 regions)	0.19 (0.019)	0.59 (0.009)	61% (0.78%)	33% (1.35%)	0.39 (0.016)	0.68 (0.008)	70% (0.71%)	48% (1.23%)
Phone (Expenditures)	0.13 (0.020)	0.57 (0.010)	60% (0.71%)	32% (1.23%)	0.26 (0.017)	0.63 (0.009)	65% (0.81%)	40% (1.40%)
Phone (Machine Learning)	0.38 (0.017)	0.70 (0.009)	69% (0.87%)	47% (1.18%)	0.45 (0.015)	0.73 (0.007)	71% (0.74%)	50% (1.28%)
<i>Panel B: Common alternative targeting methods that could not be implemented in Togo in 2020</i>								
Asset Index	0.10 (0.018)	0.55 (0.009)	60% (0.48%)	30% (0.83%)	0.51 (0.014)	0.75 (0.007)	74% (0.69%)	54% (1.19%)
PPI		[data not available]			0.63 (0.011)	0.81 (0.006)	77% (0.73%)	60% (1.25%)
PMT		[data not available]			0.72 (0.009)	0.85 (0.005)	78% (0.70%)	63% (1.20%)
<i>Panel C: Additional counterfactual targeting methods that were feasible in Togo in 2020</i>								
Random	0.00 (0.021)	0.50 (0.082)	59% (0.74%)	30% (0.26%)	0.00 (0.019)	0.50 (0.010)	59% (0.79%)	29% (1.36%)
Occupation (As implemented)	-0.11 (0.019)	0.45 (.007)	55% (0.62%)	22% (1.07%)	-0.09 (0.019)	0.46 (0.095)	56% (0.53%)	24% (0.91%)
Occupation (Optimally designed)	0.25 (0.016)	0.61 (0.008)	66% (0.58%)	41% (1.00%)	0.41 (0.016)	0.69 (0.008)	72% (0.72%)	52% (1.25%)

Table 1 | Performance of targeting mechanisms. Targeting performance using mobile phone data and machine learning (highlighted) in comparison to counterfactual targeting strategies. The “true poor” are those who, according to survey data, are in the poorest 29% of the population (the 29% threshold reflects the budget constraint of the rural Novissi expansion). The first four columns evaluate targeting with a 2020 phone survey representative of subscribers in Togo’s 100 poorest cantons, using a PMT as ground truth for poverty since consumption data were not collected. The last four columns evaluate targeting using nationally representative household survey data collected in 2018-2019, using consumption as a ground truth. Panel A compares the phone-based PMT (highlighted) to alternative targeting methods that the Government of Togo considered prior to expanding Novissi to rural areas. Panel B shows the performance of targeting methods that are commonly implemented but were infeasible in Togo at the time. Panel C indicates the performance of other targeting methods the government could have used. Accuracy, precision, and recall are evaluated by the extent to which they reach the poorest 29% (by construction, precision and recall are equal in this simulation, and are equal to 1 – exclusion error). Standard deviations, produced from 1,000 bootstrap simulations, shown in parentheses.

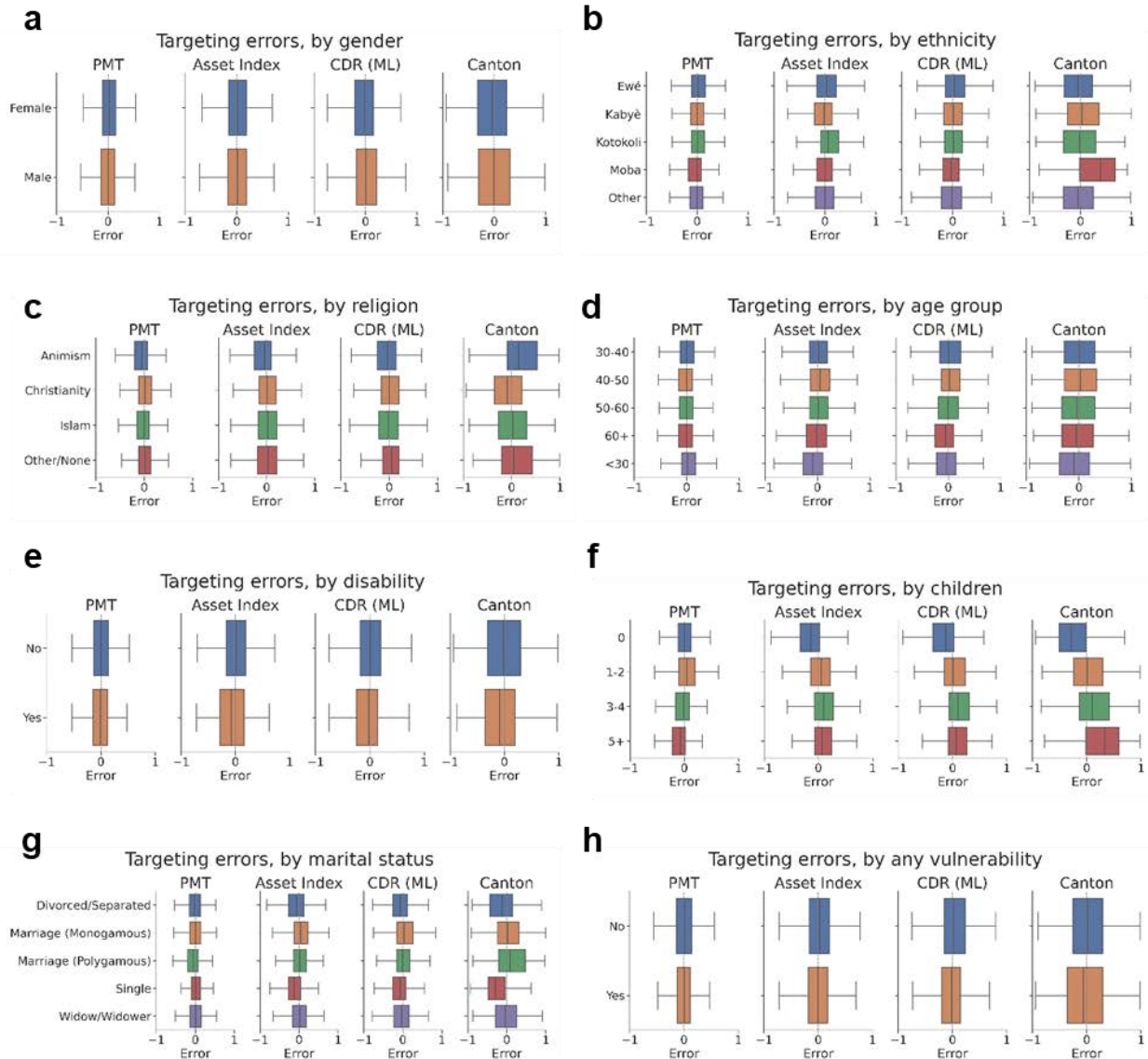
Exclusion Source	Proportion Included	Data and Calculations
Voter ID possession	83% - 98%	According to administrative data, 3,633,898 individuals are registered to vote in Togo. The electoral commission of Togo reports that this corresponds to 86.6% of eligible adults ⁵⁰ . The total adult population in Togo is not certain (the last census was in 2011), but Togo’s national statistical agency (https://inseed.tg/) estimates that there are 3,715,318 adults in Togo; the United Nations estimates 4.4 million adults ⁴⁸ . These imply a voter ID penetration rate of either 82.6% and 97.8%, respectively.
SIM card and mobile phone access	50% - 85%	65% of individuals interviewed in the 2018-2019 field survey ($N = 6,171$) reported owning a phone; 85% of individuals were in a household with one or more phones. Rural penetration is lower (50% of individuals and 77% of households), as is penetration among women (53% for women vs. 79% for men; in rural areas, it is 33% for women and 71% for men) – see Supplementary Fig. 3. Phone penetration in Togo likely increased between the field survey (2018-2019) and the Novissi expansion (October 2020); the Togolese government estimates 82% SIM card penetration ⁵⁰ .
Past mobile phone use	72% - 97%	Poverty estimates were only constructed for subscribers who placed at least one outgoing transaction between March and September 2020. In a typical month, 2.5% of all phone numbers are newly registered (Supplementary Fig. 6), so with a one-month gap between poverty inference and program registration we would expect 95-97% of registrations to be associated with a poverty score. However, 27% of all Novissi registrations (November-December 2020) did not match to CDR, likely due to new SIM purchases or registration on infrequently used SIMs (see Methods, ‘Program Exclusions’).
Program awareness	35% - 46%	245,454 unique subscribers attempted to register for the rural Novissi program. The total voting population of eligible areas is 528,562, implying a maximum registration rate of 46.44%. However, not all 245,454 registration attempts were made by people living in eligible areas; examining administrative data on home location from successful registrations we estimate that 87% of registration attempts came from eligible areas, implying an attempted registration rate of 40.40%. An alternative way to estimate attempted registration rates involves comparing the number of registration attempts made by phones below the poverty threshold (69,753) to our estimate of the number of voters in eligible cantons below the poverty threshold based on inferred home locations from mobile phone data (174,425, see Supplementary Methods section 4 for details), which implies an attempted registration rate of 34.79% after scaling by 87% (to account for registrations that came from outside of eligible areas).
Registration challenges	72%	Registration for the Novissi program requires entering basic information into a USSD (phone-based) platform. According to program administrative data, of the 245,454 subscribers who attempted registration, 176,517 (71.95%) eventually succeeded. The average registration required four attempts.
Targeting errors	47%	Based on the estimates from our targeting simulations using the 2020 phone survey (Table 1), the exclusion error rate of the phone-based targeting algorithm is 53%.

Table 2 | Sources of exclusion from rural Novissi benefits. We use multiple sources of administrative data, survey data, and government sources to estimate the extent to which different elements of the Novissi program’s design may have led to errors of exclusion. Novissi eligibility requirements included: a valid voter ID (as a unique identifier and for home location), access to a mobile phone (to fill register using the USSD platform), past mobile network transactions (to estimate poverty from mobile network behavior), program awareness (to know that the program exists and to attempt to register), ability to register via the USSD platform (which requires basic digital literacy), as well as targeting errors from the phone-based machine learning algorithm. While this table calculates sources of exclusion as though they were all independent, Extended Data Table 5 uses survey data to calculate overlaps in exclusions.

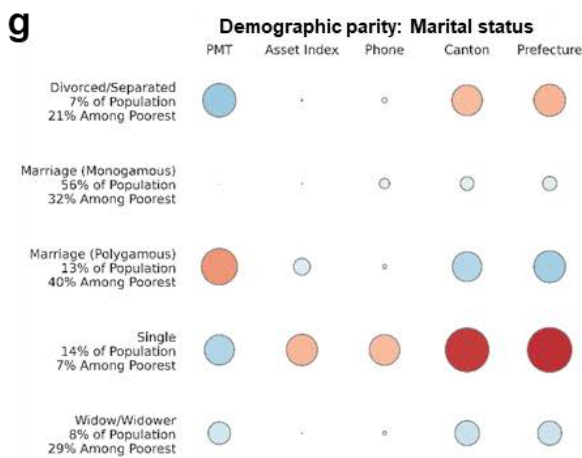
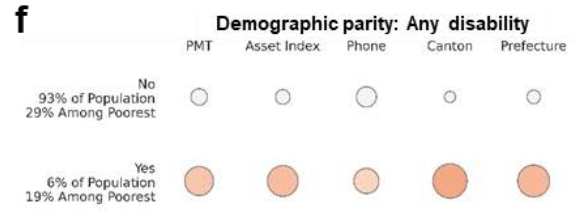
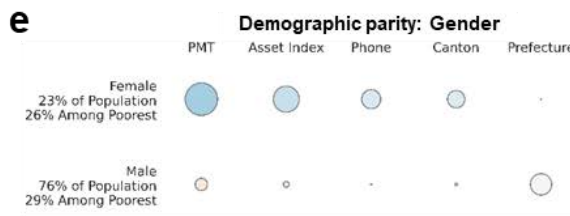
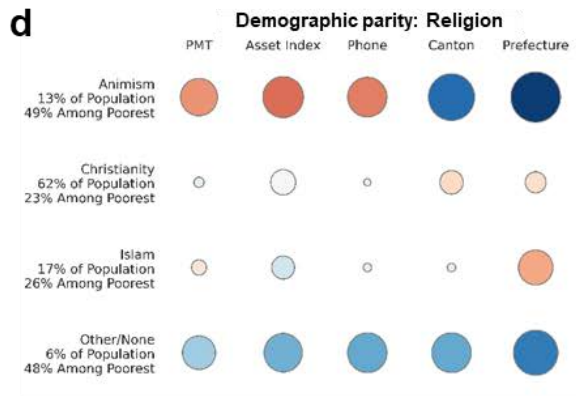
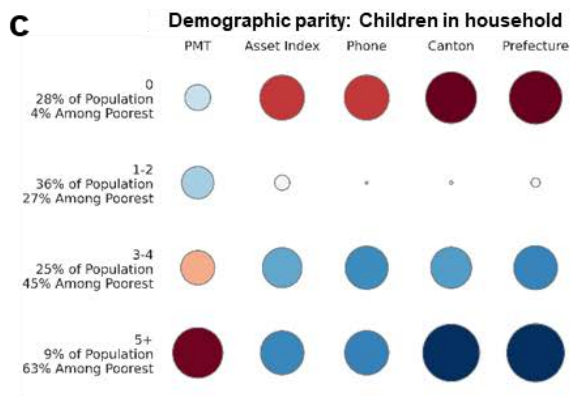
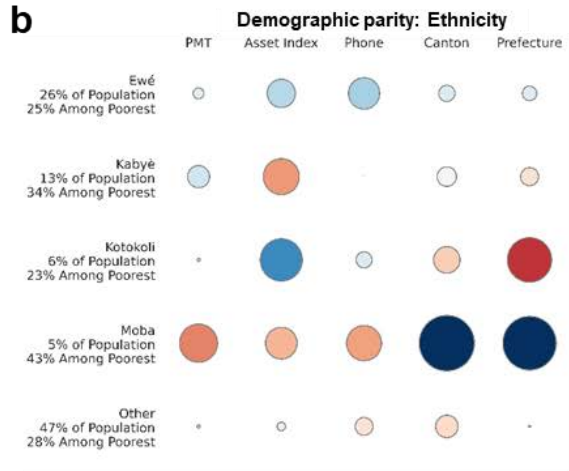
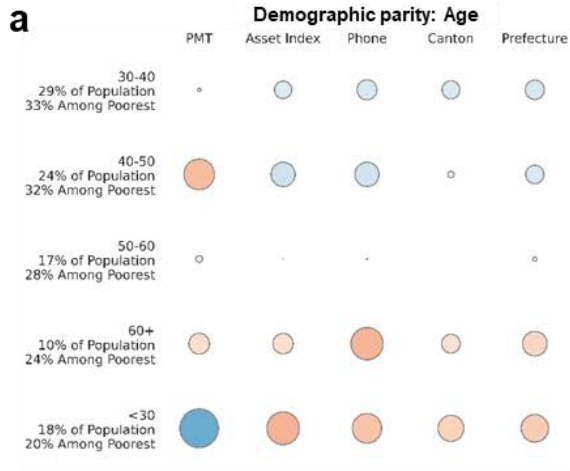
Extended Data (Figure Legends and Tables)



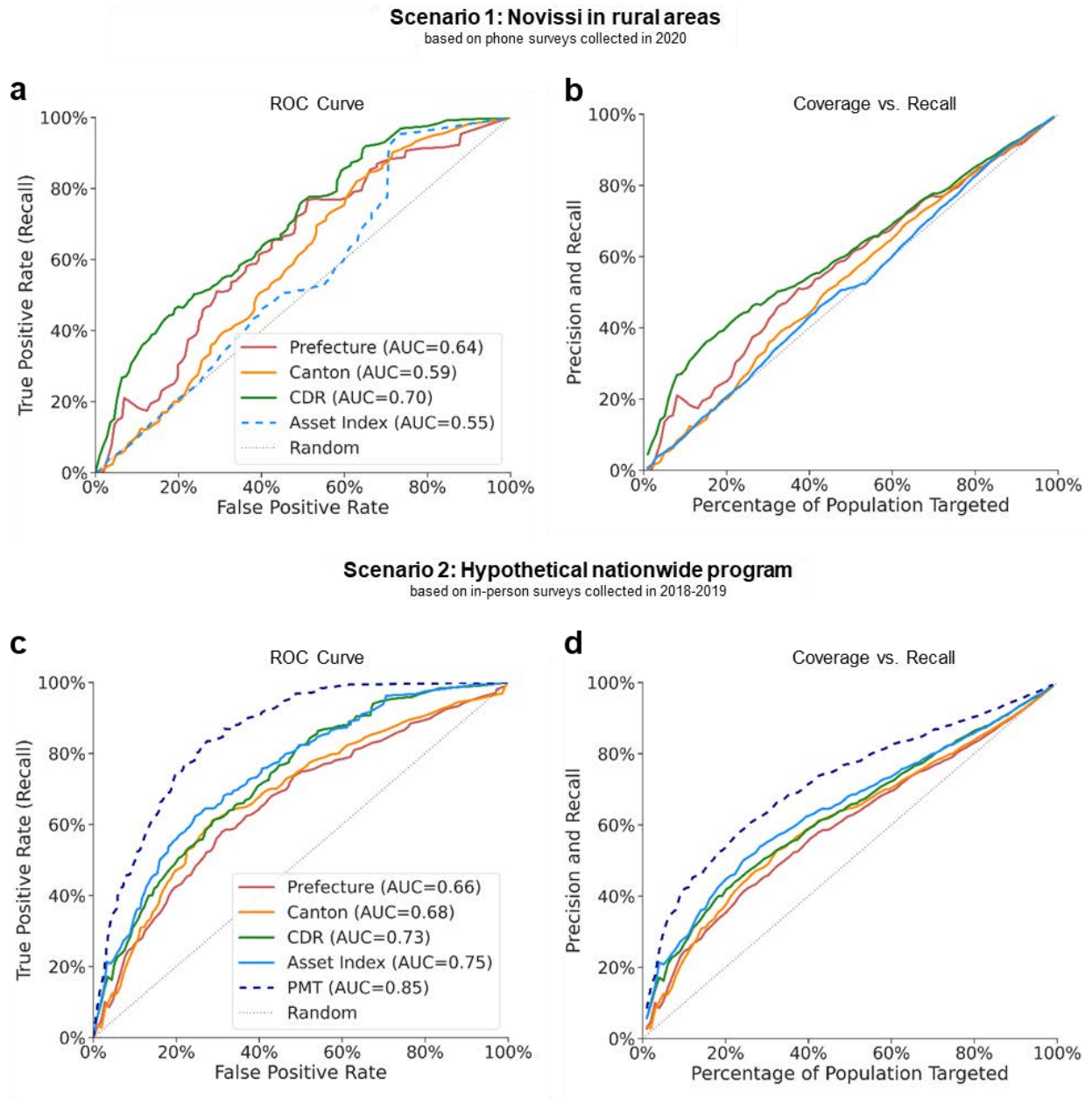
Extended Data Figure 1 | Overview of targeting methodology. a) Regional targeting. Micro-estimates of poverty (middle)¹³, are overlaid with population data to produce canton-level estimates of wealth. Individuals registered in the 100 poorest cantons (right) are eligible for benefits. **b) Individual targeting.** A machine learning algorithm is trained using representative survey data to predict consumption from features of phone use (Methods, ‘Machine Learning Methods’). The algorithm constructs poverty scores that are correlated with ground-truth measures of consumption (left). Subscribers who register for the program in targeted cantons with estimated consumption less than USD \$1.25/day are eligible for benefits (right). The red distribution shows the predicted wealth distribution of the entire population of Togo; the blue distribution shows the predicted wealth distribution in the 100 poorest cantons; and the green section indicates the predicted wealth distribution of Novissi beneficiaries.



Extended Data Figure 2 | Fairness with normalized rank residuals. Boxplots showing distributions of normalized rank residuals (see Methods, ‘Fairness’) aggregated by subgroup, using the 2018-2019 field survey dataset ($N = 4,171$). Boxes show the 25th and 75th percentiles, and the center line shows the median of the distribution. Left-shifted boxes indicate groups that are consistently under-ranked by a given targeting mechanism, right-shifted boxes indicate groups that are consistently over-ranked by a given targeting mechanism.



Extended Data Figure 3 | Fairness with demographic parity. We evaluate demographic parity across subgroups by comparing the proportion of a subgroup targeted under counterfactual approaches to the proportion of the subgroup that falls into the poorest 29% of the population (using data from the 2018-2019 field survey matched to CDR, $N = 4,171$). Bubbles show the percentage point difference between the proportion of the subgroup that is targeted and the proportion that is poor according to ground-truth data. Large red bubbles indicate groups that are over-targeted; large blue bubbles indicate groups that are under-targeted.



Extended Data Figure 4 | Targeting performance at different levels of program coverage.

Top figures (**a** and **b**) show performance for the rural Novissi program, evaluated using 2020 phone survey. Bottom figures (**c** and **d**) correspond to the hypothetical national program, evaluated using the 2018-2019 field survey. ROC curves on left (**a** and **c**) indicate the true positive and false positive rates at different targeting thresholds. Coverage vs. Recall figures on right (**b** and **d**) show how precision and recall vary as the percentage of the population receiving benefits increases, i.e., they indicate the precision and recall for reaching the poorest $k\%$ of the population in programs that target the poorest $k\%$. (Precision and recall are thus the same for each value of k by construction – see Methods, ‘Measures of Targeting Quality’).

	Targeting Novissi in rural Togo Based on 2020 Phone Survey (N = 8,915)			Hypothetical nationwide program Based on 2018-2019 Field Survey (N = 4,171)		
	Accuracy	Precision	Recall	Accuracy	Precision	Recall
<i>Panel A: Targeting methods considered by the Government of Togo in 2020</i>						
Prefecture (Admin-2 regions)	59% (0.94%)	61% (1.49%)	37% (0.99%)	67% (0.73%)	51% (1.26%)	44% (1.09%)
Canton (Admin-3 regions)	54% (0.86%)	53% (1.47%)	32% (0.91%)	69% (0.73%)	54% (1.26%)	47% (1.08%)
Phone (Expenditures)	53% (0.85%)	50% (1.32%)	31% (0.90%)	64% (0.85%)	45% (1.46%)	39% (1.25%)
Phone (Machine Learning)	61% (0.77%)	64% (0.94%)	39% (0.81%)	69% (0.73%)	55% (1.27%)	48% (1.09%)
<i>Panel B: Common alternative targeting methods that could not be implemented in Togo in 2020</i>						
Asset Index	53% (0.54%)	51% (0.009)	31% (0.57%)	72% (0.71%)	60% (1.23%)	51% (1.05%)
PPI	[data not available]			76% (0.73%)	67% (1.26%)	57% (1.09%)
PMT	[data not available]			78% (0.70%)	70% (1.20%)	60% (1.03%)
<i>Panel C: Additional counterfactual targeting methods that were feasible in Togo in 2020</i>						
Random	53% (0.84%)	51% (1.31%)	31% (0.88%)	56% (0.81%)	33% (1.39%)	28% (1.20%)
Occupation (As implemented)	47% (0.76%)	41% (1.17%)	25% (0.80%)	54% (0.55%)	29% (0.96%)	25% (0.82%)
Occupation (Optimally designed)	59% (0.68%)	61% (1.61%)	37% (0.71%)	71% (0.74%)	58% (1.28%)	50% (1.10%)

Extended Data Table 1 | Performance of targeting households below the *extreme* poverty line. Analysis is similar to that presented in Table 1, but targeting is evaluated on the extent to which each method (still targeting the poorest 29%) provides benefits to individuals consuming less than the international *extreme* poverty line, set at 75% of the international poverty line or USD \$1.43 per person per day (53% of observations in the 2020 phone survey dataset and 41% of observations in the 2018-2019 field survey). Spearman correlation and AUC are not reported here as they do not depend on the classification threshold, and are thus identical to the values reported in Table 1.

	Targeting Novissi in rural Togo Based on 2020 Phone Survey (N = 8,915)			Hypothetical nationwide program Based on 2018-2019 Field Survey (N = 4,171)		
	Accuracy	Precision	Recall	Accuracy	Precision	Recall
<i>Panel A: Targeting methods considered by the Government of Togo in 2020</i>						
Prefecture (Admin-2 regions)	47% (0.67%)	86% (1.16%)	34% (0.46%)	60% (0.67%)	68% (1.15%)	39% (0.67%)
Canton (Admin-3 regions)	44% (0.87%)	80% (1.51%)	31% (0.59%)	62% (0.69%)	71% (1.19%)	41% (0.68%)
Phone (Expenditures)	41% (0.77%)	76% (1.32%)	30% (0.52%)	57% (0.91%)	63% (1.56%)	36% (0.90%)
Phone (Machine Learning)	48% (0.76%)	87% (1.30%)	34% (0.51%)	63% (0.69%)	72% (1.19%)	42% (0.69%)
<i>Panel B: Common alternative targeting methods that could not be implemented in Togo in 2020</i>						
Asset Index	42% (0.52%)	77% (0.89%)	30% (0.35%)	65% (0.69%)	76% (1.19%)	44% (0.68%)
PPI	[data not available]			69% (0.66%)	83% (1.14%)	48% (0.66%)
PMT	[data not available]			71% (0.56%)	87% (0.97%)	50% (0.56%)
<i>Panel C: Additional counterfactual targeting methods that were feasible in Togo in 2020</i>						
Random	39% (0.76%)	73% (1.31%)	29% (0.51%)	49% (0.88%)	49% (1.51%)	28% (0.87%)
Occupation (As implemented)	38% (0.77%)	71% (1.33%)	28% (0.52%)	48% (0.61%)	46% (1.05%)	27% (0.60%)
Occupation (Optimally designed)	46% (0.61%)	84% (1.06%)	33% (0.42%)	64% (0.68%)	74% (1.18%)	43% (0.68%)

Extended Data Table 2 | Performance of targeting households below the poverty line.

Analysis is similar to that presented in Table 1, but targeting is evaluated on the extent to which each method (still targeting the poorest 29%) provides benefits to individuals consuming less than the international poverty line of USD \$1.90 per person per day (76% of observations in the 2020 phone survey dataset and 57% of observations in the 2018-2019 field survey). Spearman correlation and AUC are not reported here as they do not depend on the classification threshold, and are thus identical to the values reported in Table 1.

Feature	Importance	Feature	Importance
<i>Panel A: Predicting consumption, using 2018-2019 national household survey</i>		<i>Panel B: Predicting a PMT, using 2020 phone survey</i>	
% in Tone	20	% in Tandjoare	31
% nocturnal calls	18	% in Doufelgou	21
# in Lome Commune	17	% in Cinkasse	18
% in Tandjoare	15	Mean data volume	15
% in Tchamba	14	% in Kpendjal-Ouest	14
% in Lome Commune	13	% in Agoe-Nyive	14
% in Agoe-Nyive	13	# in Kpendjal	14
SD call duration (weekends)	12	Median call duration (night)	13
Min time between calls (weekdays)	11	# in Golfe	11
Radius of gyration (night)	11	% in Keran	11
<i>Panel C: Predicting a PMT, using 2018-2019 national household survey</i>		<i>Panel D: Predicting a PMT, using 2018-2019 survey restricted to rural areas</i>	
% in Tchamba	25	% in Tchamba	29
% in Tandjoare	24	% in Tandjoare	22
% in Doufelgou	22	% in Doufelgou	22
% in Agoe-Nyive	20	% in Agoe-Nyive	22
% in Lome Commune	19	% in Kloto	21
# in Lome Commune	19	% in Tone	16
% in Tone	17	Radius of gyration	15
Radius of gyration (night)	16	% in Kpendjal-Ouest	12
Entropy of text contacts (day)	14	# in Dankpen	11
% in Tchaoudjo	13	SD churn rate	11

Extended Data Table 3 | Feature importances. Feature importances for the 10 most important features selected by machine learning models trained to predict (a) Proxy Means Test from CDR, using a 2020 phone survey of mobile subscribers in Togo’s 100 poorest cantons ($N = 8,915$); (b) consumption from CDR in the 2018-2019 field survey dataset ($N = 4,171$); (c) PMT from CDR in the 2018-2019 field survey dataset ($N = 4,171$), and (d) PMT from CDR in the 2018-2019 field survey dataset restricted to rural areas ($N = 2,306$). Feature importance is calculated based on the total number of times a feature is split upon in the prediction ensemble. Features are color-coded as follows: CDR features are shown in blue, location features in green, mobile money features in purple, and mobile data features in red.

Temporal stability of phone-based targeting (rural Novissi program)

Based on 2020 Phone Survey ($N = 7,064$)

	Model	Phone data	Spearman	AUC	Accuracy	Precision	Recall
<i>Panel A: Reaching the 29% Poorest</i>							
(1) Best case	Current	Current	0.42 (0.019)	0.72 (0.010)	72% (0.73%)	51% (1.27%)	51% (1.27%)
(2) Old model	Old	Current	0.35 (0.019)	0.68 (0.010)	69% (0.75%)	46% (1.29%)	46% (1.29%)
(3) Old model and data	Old	Old	0.36 (0.019)	0.68 (0.010)	68% (0.74%)	44% (1.27%)	44% (1.27%)
(4) Geographic (Prefecture)			0.31 (0.020)	0.65 (0.009)	67% (0.89%)	43% (1.53%)	43% (1.53%)
(5) Geographic (Canton)			0.20 (0.023)	0.59 (0.011)	62% (0.83%)	34% (1.44%)	34% (1.44%)
<i>Panel B: Reaching the extreme poor (48% of observations)</i>							
(1) Best case	Current	Current	0.42 (0.019)	0.72 (0.01-)	62% (0.75%)	68% (1.30%)	41% (0.79%)
(2) Old model	Old	Current	0.35 (0.019)	0.68 (0.010)	60% (0.69%)	64% (1.30%)	38% (0.82%)
(3) Old model and data	Old	Old	0.36 (0.019)	0.68 (0.010)	60% (0.75%)	63% (1.29%)	38% (0.78%)
(4) Geographic (Prefecture)			0.31 (0.020)	0.65 (0.009)	59% (0.94%)	62% (1.61%)	38% (0.98%)
(5) Geographic (Canton)			0.20 (0.023)	0.59 (0.011)	54% (0.96%)	53% (1.65%)	32% (1.00%)
<i>Panel C: Reaching the poor (74% of observations)</i>							
(1) Best case	Current	Current	0.42 (0.019)	0.72 (0.01-)	49% (0.56%)	90% (0.97%)	35% (0.38%)
(2) Old model	Old	Current	0.35 (0.019)	0.68 (0.010)	47% (0.67%)	86% (1.16%)	34% (0.46%)
(3) Old model and data	Old	Old	0.36 (0.019)	0.68 (0.010)	47% (0.62%)	86% (1.08%)	34% (0.42%)
(4) Geographic (Prefecture)			0.31 (0.020)	0.65 (0.009)	48% (0.69%)	87% (1.18%)	34% (0.46%)
(5) Geographic (Canton)			0.20 (0.023)	0.59 (0.011)	44% (0.98%)	80% (1.69%)	31% (0.66%)

Extended Data Table 4 | How quickly does the accuracy of a phone-based targeting model degrade? Table compares three scenarios: (1) “Best case”: when the model is calibrated using survey data and phone data gathered just before deployment – these results are comparable to the paper’s main analysis (slight differences are due to the sample restrictions described below); (2) “Old model”: when the model is trained using a survey conducted two years before deployment, but the phone data are collected just before deployment; and (3) “Old model and data”: when the phone-based wealth estimates are generated using survey and phone data from two years prior. Rows (4) and (5) show geographic targeting results using the same sample as in rows (1) – (3). In the simulations, the “old” data are from the 2018-19 national household survey and corresponding 2019 phone dataset; the 2020 phone survey PMT is used as the ground truth measure of poverty (restricted to respondents for whom CDR are available in 2019 and 2020, $N = 7,064$).

Exclusion Source	<i>N</i>	Succeed	Drop Out	% Remaining	PMT	% Women	Age
<i>Panel A: Attrition among overall population</i>							
All	8,915	--	--	100.00%	1.62 (0.72)	23% (42%)	33.21 (11.91)
Own a voter ID	8,898	99.70%	0.30%	99.70%	1.62 (0.71)	23% (42%)	33.17 (11.87)
Attempt to register	5,145	45.48%	54.52%	45.34%	1.45 (0.57)	23% (42%)	33.30 (12.00)
Succeed in registration	4,092	76.84%	23.16%	34.84%	1.43 (0.54)	23% (42%)	33.05 (11.87)
Targeted by phone PMT	2,277	46.99%	53.01%	16.37%	1.28 (0.44)	21% (40%)	35.79 (11.96)
<i>Panel B: Attrition among the poorest 29%</i>							
All poor	3,209	--	--	100.00%	1.00 (0.15)	19% (39%)	36.22 (10.99)
Own a voter ID	3,207	99.77%	0.23%	99.77%	1.00 (0.15)	19% (39%)	36.16 (10.89)
Attempt to register	2,253	60.55%	39.45%	60.41%	0.99 (0.15)	20% (40%)	36.94 (11.19)
Succeed in registration	1,845	78.61%	21.39%	47.49%	0.99 (0.15)	19% (40%)	35.37 (11.03)
Targeted by phone PMT	1,257	60.56%	39.44%	28.76%	0.96 (0.15)	17% (37%)	36.67 (10.83)

Extended Data Table 5 | Overlapping sources of exclusion from rural Novissi. Progressive sources of attrition from the rural Novissi program, where each row shows exclusion conditional on exclusions from preceding rows. The final three columns show characteristics of the population “succeeding” at each step. Panel A: Results estimated using the 2020 phone survey ($N = 8,915$). Panel B: Results estimated for just the poorest 29% from the 2020 survey ($N = 3,209$). There is no attrition based on mobile phone ownership or past phone use in this sample (in contrast to Table 2) since only active phone users were sampled for the phone survey. Values reweighted using sample weights. (In some cases, sample weights create large differences in the weighted and raw percentages. For instance, 5,145 out of 8,898 voters (57.8%) attempt to register (Panel A), but the weighted percentage is 45.5%. The importance of sample weights is consistent with the wide distribution of sample weights shown in Supplementary Fig. 10).

	Consumption	PMT	Asset Index
<i>Panel A: 2018-2019 Field Survey (N = 4,171)</i>			
ML	0.46	0.62	0.74
Single Feature	0.13	0.16	0.11
<i>Panel B: 2018-2019 Field Survey, Rural Only (N = 2,306)</i>			
ML	0.31	0.44	0.51
Single Feature	0.09	0.12	0.08
<i>Panel C: 2020 Phone Survey (N = 8,915)</i>			
ML	--	0.41	0.40
Single Feature	--	0.13	0.14

Extended Data Table 6 | Performance of phone-based approach to predicting wealth and consumption. Accuracy (Pearson correlation coefficients) for predicting poverty measures from CDR. ML predictions are produced over 5-fold cross validation and evaluated for pooled correlation. The “single feature” model estimates wealth and consumption based on the individual’s total expenditures on calling and texting.

Supplementary Materials for

Machine Learning and Phone Data Can Improve Targeting of Humanitarian Aid

Emily Aiken, Suzanne Bellue, Dean Karlan, Chris Udry, Joshua Blumenstock

Correspondence to: jblumenstock@berkeley.edu

Contents

Supplementary Discussion.....	2
1. Related work.....	2
2. Limitations and Concerns.....	2
Supplementary Methods	6
3. Selection of Variables for Proxy-Means Test.....	6
4. Home Location Inference from Mobile Phone Data	7
5. Design of the 2020 Phone Survey	9
References.....	13
Supplementary Figures	16
Supplementary Tables.....	27

Supplementary Discussion

1. Related work

There is a rich history of theoretical and empirical work that compares and evaluates methods for targeting social transfer programs. While there is increasing interest in “universal basic income”, in which everyone is eligible for transfers, most countries use one or more targeting mechanisms to determine eligibility¹. Typically, the goal of targeting is to ensure that the poorest individuals receive transfers.¹

Many programs include some degree of *self-targeting*, in which beneficiaries are required to take some action in order to receive benefits³⁻⁵. If the benefits of the program, relative to the costs associated with that action, are higher for poorer people, self-targeting can direct a greater share of benefits to the poor. *Geographic targeting* is also common, whereby benefits are restricted to individuals who live in specific regions^{6,7}. Empirical evidence on geographic targeting indicates that more granularly targeted programs can be more effective at prioritizing the poor, but the implementation of such programs requires fine-grained poverty maps and distribution mechanisms that can be deployed in small regions⁸⁻¹⁰. With *proxy means tests (PMT)*, a number of variables are collected for each household, which are then used to impute an approximate measure of consumption or wealth for that household.^{11,12} Likewise, a simple poverty scorecard or *poverty probability index (PPI)* uses a small number of variables to impute a poverty score.^{13,14} PMTs and PPIs are frequently used in LMICs, but do require that the government collect and maintain a comprehensive social registry that records the information of each household. Finally, *community-based targeting (CBT)* approaches rely on members of the community to identify the poorest households in the area^{15,16}. CBT-based approaches do not always target the lowest-consumption households, but allow the community to define their own notion of poverty, which can lead to higher satisfaction among community members⁴ but may also lower perceptions of program legitimacy¹⁷.

2. Limitations and Concerns

While mobile phone data can create new options for the accurate targeting of humanitarian aid, there are several important limitations. A full discussion of the social, political, and ethical implications of these issues has been the focus of prior work and is beyond the scope of this article¹⁸⁻²²; we nonetheless highlight a few key issues that we believe require careful consideration before these methods can be implemented in a policy environment:

Phone ownership and access: As discussed in Methods, ‘Program Exclusions’, many individuals in LMICs do not own mobile phones. Thus, any targeting system based on mobile phone data may exclude those without phones from receiving program benefits. In the case of the Novissi program, the government used the mobile money system to disburse the cash transfers as a way to minimize human contact during the pandemic. Thus, in Togo, the use of phone data for

¹ How a program defines “poverty” is also a source of considerable debate². In this paper, we use the term “socioeconomic status” somewhat loosely to refer to an individual’s access to resources. By contrast, we use “consumption” to refer to how much an individual spends or consumes, and “wealth” to refer to an individual’s assets. “Poverty” is a condition in which an individual’s access to resources falls below a minimal level, based on consumption or wealth, as described in Methods, ‘Survey Data’.

targeting only created additional exclusions by requiring that program registrants had made at least one transaction on their SIM card in the months prior to registration. In general, incomplete mobile phone access highlights the need to allow for alternative pathways for individuals to register and receive benefits, and to create additional mechanisms for appeals, grievance redress mechanisms, and manual enrollment.

Data privacy: Mobile phone metadata, even when pseudonymized, contains sensitive information. Methods, ‘Data Privacy Concerns’ describes several steps taken to protect the confidentiality of the data used in this project. More generally, special considerations arise when using personal data from vulnerable populations²³, and human rights doctrine emphasizes that any form of communications surveillance should be “necessary and proportionate”²⁴.

In implementing the approach described in this paper, we developed an IRB protocol, as well as a data management plan, that was approved by U.C. Berkeley’s Committee for the Protection of Human Subjects. We followed principles of data minimization to limit the data collected and stored, and implemented organizational safeguards to restrict access to data. As an example, only IRB-approved researchers ever received access to CDR; data from the phone companies were shared with neither the Government of Togo nor GiveDirectly. Even the poverty scores derived from the phone data were restricted to IRB-approved researchers; the only data the government received was the list of SIM cards belonging to eligible beneficiaries below the targeted poverty threshold.

Future projects using mobile phone data for targeting should ensure that principles of data minimization and data sunsetting restrict the use of sensitive data to social protection objectives and limit the potential for “function creep.”²⁵ Further research on applying the guarantees of differential privacy to mobile phone metadata^{26,27} or implementing federated learning systems²⁸ could reduce the risk of data misuse or central data breaches.

Data access and consent: The fact that our approach requires access to mobile phone data owned by private companies poses an obstacle to the immediate and widespread use of such data for targeting humanitarian aid. There now exist several general frameworks and recommendations to facilitate the use of CDR in humanitarian applications^{19,29}. Yet such frameworks are still nascent, and without careful consideration may exclude important stakeholders and perspectives²²; they also widen the scope for private companies to influence humanitarian and development decisions³⁰. There also exist many ethical frameworks that rely on informed consent from participants for the use of personal data, including digital data such as CDR^{31,32}. Future programs should consider how consent pathways can be integrated with phone-based targeting, including opt-in (calculating poverty scores only after consent is provided) and opt-out (scrubbing data if consent is not provided at the time of registration) options.

Data representativity: To train the machine learning models, ground truth measures of consumption and wealth were collected using in-person and phone surveys. Since response rates were imperfect in the phone survey, we reweighted survey observations to make the training data more representative of all mobile subscribers (Methods, ‘Survey Data’). However, there are limits to the representativity of our training data, as dynamics of phone ownership and phone sharing vary across population subgroups (Supplementary Figure 3), and reweighting is an imperfect proxy.

To test for systematic bias based on data representativity, we perform ex-post audits to limit the likelihood that the trained models systematically disadvantage specific subgroups of the population (Methods, ‘Fairness’), and find that the phone-based targeting method is no more biased than counterfactual targeting approaches. We believe such audits are essential to future work on wealth prediction and targeting based on nontraditional data. Audits could be improved with additional context-specific research about which sub-populations are at the greatest risk for systematic exclusion (for example, in this paper we test for bias across age groups, genders, ethnicities, and more), and on considering alternative definitions for bias and fairness.^{33,34}

Unit of analysis: As noted in Methods, ‘Experimental Design’, our analysis focuses on *individuals* rather than *households* as the unit of analysis, partly reflecting the design of the Novissi program, and partly because there are no data in Togo that associate individuals with households. This limitation is important, since many real-world programs are targeted at the household level, but CDR are more naturally linked to individual subscribers. An important area for future work will thus be to explore the extent to which CDR can facilitate household-level targeting. Such work must account for the fact that a single SIM card is often shared across multiple members of the same household (and occasionally between households), and that some individuals use multiple SIM cards. Ideally, such an analysis would leverage authoritative data that uniquely identifies and links households, individuals, SIM cards, and phones.

Method of evaluation: Our main results are based on simulations of targeting methodologies using survey data collected prior to expansion of Novissi. An alternative approach to evaluating targeting performance would rely on survey data after program implementation, which would make it possible to more directly verify who did and didn’t receive program benefits, address issues related to the unit of analysis described above, and better attribute exclusion errors to different aspects of program design. While public health considerations in Togo prevented us from conducting a post-program survey, we hope future implementations of phone-based targeting can use post-program surveys to provide complementary evidence to what is described in this paper.

Poverty dynamics: The phone-based approach we describe uses machine learning algorithms to predict which individuals are “poor”, based on ground-truth assessment of poverty collected in surveys prior to program implementation. In the actual rural Novissi program, the ground truth measure of poverty was based on a proxy means test; in the hypothetical national program, ground truth is based on consumption (Methods, ‘Survey Data’). However, particularly in the context of a crisis, an individual’s poverty status can change; in such settings, pre-program poverty assessments may not accurately capture the population with the greatest need for support. Our data do not permit us to test whether phone data and machine learning can be used to determine if an individual has experienced a sudden fall in income or consumption, but we believe this is a promising area for future work.

Manipulation and gaming: When mobile phone data are used to determine eligibility for social benefits, individuals have incentives to strategically alter their behavior in order to “game” the system. This dilemma is not unique to phone-based targeting; it is a key consideration in the design of any targeting mechanism^{35,36}, and one that affects traditional proxy means tests and poverty scorecards^{37,38}. However, recent evidence suggests that such distortionary effects may be limited³⁹, and complex eligibility criteria (such as the gradient boosting procedure described in Methods, ‘Machine Learning Methods’) should limit the scope for such gaming⁴⁰. With Novissi

in Togo, the one-off nature of the program likely eliminated most scope for strategic behavior; however, if such an approach were used continuously over time, alternative “manipulation-proof” approaches to machine learning may be more appropriate⁴¹.

General equilibrium considerations: Our analysis of targeting effectiveness assumes there are no general equilibrium effects of the program on prices, wages, or interactions with informal transfers or insurance. For example, geographic targeting of transfers might lead to localized inflows of cash transfers that are large relative to the local economy, leading to changes in local demand for goods or supply of labor and therefore prices, wages or profits of local businesses^{42,43}. Similarly, since individuals are embedded in family and broader networks of informal transfers for redistribution, patronage and insurance and different targeting choices could have different effects on these existing informal arrangements^{44,45}. Equilibrium effects such as these may have important implications for the eventual distribution of impacts from the transfers. However, to cause a reversal of the policy implication of our analysis, general equilibrium effects would need to be more nuanced than merely present – for example, it would need to be that the false negatives under one method are more likely to share resources than the false negatives on another method.

Supplementary Methods

3. Selection of Variables for Proxy-Means Test

Our proxy-means test is used in analysis for both the 2018-2019 field survey (where we evaluate the PMT’s accuracy as a targeting mechanism) and the 2020 phone survey (where we use the PMT as a measure of ground-truth poverty in the absence of a consumption measure). We construct the PMT using all observations from the 2018-2019 field survey ($N = 6,171$). We begin by identifying all information on demographics and asset ownership collected in the field survey that may correlated with poverty. In total, we identify 56 variables, including information on household assets and housing quality, education, marital status, age, ethnicity, health, location, and more.

Our goal is to identify a small subset of variables that are most predictive of household consumption. We use stepwise forward selection to identify the most predictive feature subsets of size K , for K ranging from 1 to 30. Specifically, we randomly divide our survey observations into a training set (75%) and test set (25%). For $K=1$, we train a machine learning model to predict household consumptionⁱⁱ from each feature individually, and select the feature associated with the best model. For $K=2$, we test adding each remaining feature to our model, and select the feature that adds the most predictive power. We continue the process for all K up to 30.

We perform the stepwise forward selection process first for a Ridge regression (where the optimal L2 penalty is selected via a wide grid) and second for a random forest (where the optimal ensemble size is chosen from $\{50, 100\}$ via 3-fold cross validation and the optimal tree depth is chosen from $\{2, 4, 6, 8\}$). Supplementary Figure 4 plots the predictive accuracy (measured with R^2) for each value of K for the two models.

We observe that the random forest is not significantly more accurate than the regression, and note a greater degree of overfitting with the random forest. We therefore select the Ridge regression, as the resulting coefficients are easier to interpret. We identify an “elbow” in the accuracy progression at $K=12$ features, so we use the feature subset of size $K=12$ in our PMT. These features and the weights associated with them are recorded in Supplementary Table 3.

ⁱⁱ While in the rest of this paper we use price-index adjusted per capita household consumption, in this exercise our outcome variable is raw household consumption (because the data necessary to construct price index adjusted consumption was not available to us prior to the 2020 phone survey when this analysis was performed).

4. Home Location Inference from Mobile Phone Data

We use home locations for mobile network subscribers inferred from mobile network data for a set of supplementary analyses (Supplementary Fig. 8, Supplementary Tables 9-11) and for sampling the 2020 phone survey. For the supplementary analyses, which require assigning a home prefecture and canton to each mobile network subscriber, we use standard frequency-based approaches to home location inference using the locations of cell phone towers through which subscribers place calls. These frequency-based methods have been developed in past work⁴⁶⁻⁴⁸ and are described in more detail in Section (i) below. For sampling the 2020 phone survey, which required identifying which subscribers were likely to live in rural Novissi-eligible cantons, we developed a new approach to home location inference from mobile phone metadata using supervised learning, which is described in Section (ii) below.

i. Frequency-based home location inference

“Frequency-based” methods of home location inference, based on the locations of cell towers used by subscribers, are used widely in the literature.⁴⁶⁻⁴⁸ Chi et al. (2020)⁴⁸ validate a set of different approaches to home location inference in comparison to ground truth location data, including the location (in our case, prefecture or canton) with the maximum phone transactions, the location with the maximum number of phone transactions in a given time frame (for example, daily between 8pm and 6am), and the location with the maximum number of unique days with phone transactions. Chi et al. (2020) find that the third method -- the maximum number of unique days with phone transactions -- is most accurate on their validation set of mobile phone metadata from Rwanda; we therefore select this approach to frequency-based home location inference. As displayed in Supplementary Table 10, this method is highly correlated with both the home prefecture and home canton recorded in voter data and with the home prefecture and home canton reported in surveys.

ii. Home location inference using machine learning

For sampling the 2020 phone survey, we were not interested in identifying the canton or prefecture each subscriber lived in; rather, we were interested in identifying which of the 5.83 million mobile network subscribers active between March and September 2020 lived in any of the 100 poorest cantons that were eligible for rural Novissi aid. This binary classification task is better suited to machine learning than the multiclass classification task of assigning subscribers to home locations; we therefore adopted a new approach to home location inference using machine learning for identifying subscribers likely to be living in the 100 poorest cantons for survey sampling.

Specifically, we trained our machine learning model on the dataset of all subscribers that registered for Novissi when it was first available in the Greater Lomé region (while only residents of Greater Lomé were eligible for this program, any registered voter in Togo could sign up for the platform for immediate eligibility in future programs). In total, this dataset includes 1.1 million subscribers with Novissi registration data matched to CDR. These registration data includes the canton in which each subscriber is registered to vote (we refer to this as the ‘ground-truth’ home canton).

The raw training dataset is not representative of all mobile network subscribers in Togo, as a nonrandom subset of subscribers registered for Novissi (for example, more than half of the registered subscribers are in the Greater Lomé region). To make the training data more representative, we calculated the expected share of subscribers in each canton based on the total number of voters registered in each canton and the mobile phone penetration rate in the prefecture (based on the 2018-2019 field survey). We “balanced” the training dataset by sampling observations at random from cantons with a disproportionately high number of registrants until the proportions in the training dataset reflected the expected proportion of mobile network subscribers in each canton.

Finally, we trained a machine learning model to predict whether each subscriber lived within the 100 eligible cantons. As in poverty prediction, we use a gradient boosting model with optimal hyperparameters chosen via cross-validation. The model uses the same “features” that we use for statistical home location inference – specifically the (normalized) number of unique days on which each subscriber places a transaction in each canton of Togo. The model obtains an AUC score of 0.90 and cross-validated accuracy of 93%. We then use the trained machine learning model to produce estimates of the likelihood that all 5.83 million mobile network subscribers live in an eligible canton.

5. Design of the 2020 Phone Survey

This section describes the design and implementation of the 2020 phone survey, which took place in the last week of September and the first week of October 2020.

i. Sampling

The 2020 phone survey was designed to be representative of active mobile phone subscribers living in Togo's 100 poorest cantons. The sample frame for the survey was all mobile phone subscribers active on one of the two mobile networks in Togo between March 1 and September 30, 2020 ($N = 5.83$ million). Sampling was based on four metrics associated with each mobile phone subscriber: inferred probability of living in a rural Novissi-eligible area, registration to a previous Novissi program, inferred wealth based on phone data, and total mobile phone expenditure.

- *Inferred probability of living in a rural Novissi-eligible canton:* We used the machine learning model described in Appendix B section (ii) to assign each subscriber a probability of living in a rural Novissi-eligible canton.
- *Registration to a previous Novissi program:* At the time of the survey, 22% of mobile network subscribers in Togo were already registered in the Novissi system, and therefore were associated with a ground-truth home canton based on the canton in which they are registered to vote. In our dataset of inferred home location likelihoods, we assigned any subscriber registered to vote in one of the 100 targeted cantons a 100% likelihood of geographical eligibility ($N = 86,856$). We assigned any subscriber registered to vote outside of these cantons a 0% likelihood of geographical eligibility ($N = 1,046,905$).
- *Inferred poverty based on mobile phone data:* We used ground-truth poverty data collected in a previous nationally-representative phone survey conducted in June 2020 to train a machine learning model to predict poverty from CDR. We followed the methods described in Methods, 'Machine Learning Methods' using the PMT as ground truth and CDR features from March 1 to September 30, 2020. We used the machine learning model to predict the poverty of each of the 5.83 million mobile phone subscribers in Togo.
- *Mobile phone expenditure:* We constructed the measure of total phone expenditure for each subscriber described in Methods, 'Parsimonious Phone Expenditure Method'.

Based on the total number of voters registered in targeted cantons and individual mobile phone penetration in each canton (based on the 2018-2019 field survey, measured at the prefecture level), we estimated that around 240,000 subscribers live in eligible cantons. We identified the 240,000 subscribers most likely to be living in a targeted canton (including all 86,856 subscribers registered in targeted cantons). Only these 240,000 subscribers were eligible to be surveyed.

We oversampled survey respondents based on two counterfactual targeting methods that we simulated pre-survey: predicted poverty based on phone data, and mobile phone expenditures, as described in Methods, 'Predicting Poverty from Phone Data'. We divided the 240,000 subscribers into four quartiles based on phone-inferred poverty and mobile phone expenditures. We overlapped the quartiles to form eight "cells", based on the combination of the two targeting

methods (for example, cell AA represents being in the lowest quartile by both targeting methods, while cell AD represents being in the lowest quartile by one method and the lowest quartile by the other, and cell BC represents being in the second-lowest quartile by one method and the second-highest quartile by the other). We assigned a cell weight of 0.20 to cells AD and BC (where the two methods disagree the most), a cell weight of 0.15 to cells AC and BD, a cell weight of 0.10 to cells AB and BC, and a cell weight of 0.05 to cells CD and DD (where the two methods disagree least).

Our sampling probabilities for the 240,000 survey-eligible subscribers were constructed as the product of a subscriber’s cell weight and their probability of residing in a targeted canton (so subscribers likely to be living in targeted cantons are oversampled within each cell). The distributions of these draw probabilities are shown in Supplementary Figure 10 Panel A. We use the inverse of these draw probabilities as sample weights in our downstream analysis, in combination with response weights - see Section (iv) below. We drew 40,000 phone numbers at random from the 240,000 survey-eligible subscribers, with assigned draw probabilities. We provided these 40,000 phone numbers in a random order to enumerators with the expectation that not all of them would be called in order to reach a goal interview quota of 10,000; indeed, only 30,244 phone numbers were called before the quota was reached – see Section (ii) below.

ii. Response Rates

In total, enumerators conducted 10,701 interviews out of 30,244 phone numbers that were called (overall response rate of 35.38%). Phone numbers were called in a random order, and were assigned to enumerators by language (with random assignment with groups of enumerators speaking the same language). While we have little information on subscribers pre-survey, we can examine differential nonresponse by (1) inferred geography based on CDR, (2) registration to a previous Novissi program, and (3) pre-survey mobile phone use (we focus on the phone-predicted measure of poverty and measure of daily expenditures on calls and texts that are used in the rest of the paper). Supplementary Table 12 displays response rates disaggregated along these dimensions. We find that response rates are higher for those registered to a prior Novissi program, those inferred to be living in the regions of Lomé Commune, Maritime, or Savanes, and those with a high daily phone expenditure. Section (iv) describes how we reweight survey observations to account for differential nonresponse.

iii. Removing Low-Quality Surveys

We identified unreliable enumerators by comparing the data collected in the survey with the information contained in the Novissi registry for the subset of survey respondents who had registered to a previous Novissi program. We begin our analysis by constructing “value-added” (VA) estimates for the enumerators in our data. We predict the VA of each enumerator on the basis of the correct answers to three questions for which we obtained ground-truth information from the Novissi database (canton, age and sex), and on the frequency of surveys with a single head of household (which avoids the roster part of the survey and simplifies the enumerator’s work). We control for interviewee characteristics such as region and interview language to

separate the enumerator's impact from observable interviewee selection.ⁱⁱⁱ Our approach to estimating enumerators' VA parallels the parametric empirical Bayes estimator of teacher's VA in past work.⁴⁹⁻⁵¹ We then normalize the VAs for each of the four dimensions (canton, age, gender, and number of surveys with only one adult), and take the average for enumerators who conducted more than twenty interviews. The bottom ten percent of enumerators have an average VA one standard deviation below the mean VA across all enumerators; we classify their surveys as "poor quality." The interviews of the three interviewers with an average VA lower two standard deviations below the total average VA are classified as "very poor quality."

1,180 surveys associated with enumerators who are ranked "poor quality" or "very poor quality" are removed from the dataset. We drop a further 606 surveys with missing data for the PMT or one or more of the counterfactual targeting methods, for a final survey dataset size of 8,915.

iv. Reweighting for Nonresponse

As noted in section (ii), certain groups are more likely to respond to the survey than others. To make the final analysis representative of the initial sample frame (i.e., active mobile subscribers in the 100 poorest cantons) rather than just survey respondents, we reweight survey observations by likelihood of response based on pre-survey covariates.^{52,53} In our case, we train a machine learning model (using an LGBM and the same set of hyperparameters used for wealth prediction from phone data) to predict response from our usual set of CDR features, along with whether a subscriber registered to a previous Novissi program. This model is trained on all 30,244 numbers that were called, with "response" defined as responding to the survey, including all questions necessary to construct the PMT and counterfactual targeting outcomes, consenting to matching between survey responses and mobile phone data, and that survey passing the quality assessment step (see Section iii), for a total "responded" population of 8,915 (29%). As in other machine learning models described in this paper, we tune hyperparameters over 5-fold cross validation and produce predictions for each observation over 10-fold cross validation. The model achieves a cross-validated AUC score of 0.71; feature importances for the model are shown in Supplementary Table 13. To assess the model's accuracy, Supplementary Figure 11 compares binned estimates of response probability with true rates of response, and indicates that the response prediction model is well-calibrated. Supplementary Figure 10 Panel B displays the distribution of response probabilities for observations included in the final survey dataset.

The final survey weights used in the paper are the product of the inverse of the response probability and the inverse of the sampling probability described in Section (i); the distribution of survey weights are shown in Supplementary Figure 10 Panel C.

v. Survey Content

Surveys lasted 30 minutes on average, and included questions on the demographics of the respondent and household members, assets owned by the household, subjective wellbeing of the

ⁱⁱⁱ As the phone number list was randomized and then distributed to the enumerators, we believe there is little room for sorting.

respondent, the social services available to the household, and the impacts of COVID-19 on the household. The full survey instrument is publicly available online.^{iv}

^{iv} <https://jblumenstock.com/files/papers/TogoInstrument2020.pdf>

References

1. Hanna, R. & Olken, B. A. Universal Basic Incomes versus Targeted Transfers: Anti-Poverty Programs in Developing Countries. *Journal of Economic Perspectives* **32**, 201–226 (2018).
2. Ravallion, M. *The Economics of Poverty: History, Measurement, and Policy*. (Oxford University Press, 2016).
3. Nichols, A. L. & Zeckhauser, R. J. Targeting Transfers through Restrictions on Recipients. *The American Economic Review* **72**, 372–377 (1982).
4. Alatas, V. *et al.* Self-Targeting: Evidence from a Field Experiment in Indonesia. *Journal of Political Economy* **124**, 371–427 (2016).
5. Kleven, H. J. & Kopczuk, W. Transfer Program Complexity and the Take-Up of Social Benefits. *American Economic Journal: Economic Policy* **3**, 54–90 (2011).
6. Baker, J. L. & Grosh, M. E. Poverty reduction through geographic targeting: How well does it work? *World Development* **22**, 983–995 (1994).
7. Schady, N. R. Picking the Poor: Indicators for Geographic Targeting in Peru. *Review of Income and Wealth* **48**, 417–433 (2002).
8. Coady, D. P. The Welfare Returns to Finer Targeting: The Case of The Progresa Program in Mexico. *Int Tax Public Finan* **13**, 217–239 (2006).
9. Elbers, C., Fujii, T., Lanjouw, P., Özler, B. & Yin, W. Poverty alleviation through geographic targeting: How much does disaggregation help? *Journal of Development Economics* **83**, 198–213 (2007).
10. Smythe, I. & Blumenstock, J. E. Geographic Micro-Targeting of Social Assistance with High-Resolution Poverty Maps. in *In Submission (KDD)* (2021).
11. Grosh, M. E. & Baker, J. L. *Proxy means tests for targeting social programs*. (The World Bank, 1995). doi:10.1596/0-8213-3313-5.
12. Filmer, D. & Pritchett, L. H. Estimating Wealth Effects Without Expenditure Data—Or Tears: An Application To Educational Enrollments In States Of India*. *Demography* **38**, 115–132 (2001).
13. Desiere, S., Vellema, W. & D’Haese, M. A validity assessment of the Progress out of Poverty Index (PPI)TM. *Evaluation and Program Planning* **49**, 10–18 (2015).
14. Brown, C., Ravallion, M. & van de Walle, D. A poor means test? Econometric targeting in Africa. *Journal of Development Economics* **134**, 109–124 (2018).
15. Alderman, H. Do local officials know something we don’t? Decentralization of targeted transfers in Albania. *Journal of Public Economics* **83**, 375–404 (2002).
16. Galasso, E. & Ravallion, M. Decentralized targeting of an antipoverty program. *Journal of Public Economics* **89**, 705–727 (2005).
17. Premand, P. & Schnitzer, P. Efficiency, Legitimacy, and Impacts of Targeting Methods: Evidence from an Experiment in Niger. *The World Bank Economic Review* (2020) doi:10.1093/wber/lhaa019.
18. Mann, L. Left to Other Peoples’ Devices? A Political Economy Perspective on the Big Data Revolution in Development. *Development and Change* **49**, 3–36 (2018).
19. Kerry, C. F., Kendall, J. & de Montjoye, Y.-A. Enabling Humanitarian Use of Mobile Phone Data. *Brookings Issues in Technology Innovation* (2014).
20. Blumenstock, J. E. Don’t forget people in the use of big data for development. *Nature* **561**, 170–172 (2018).

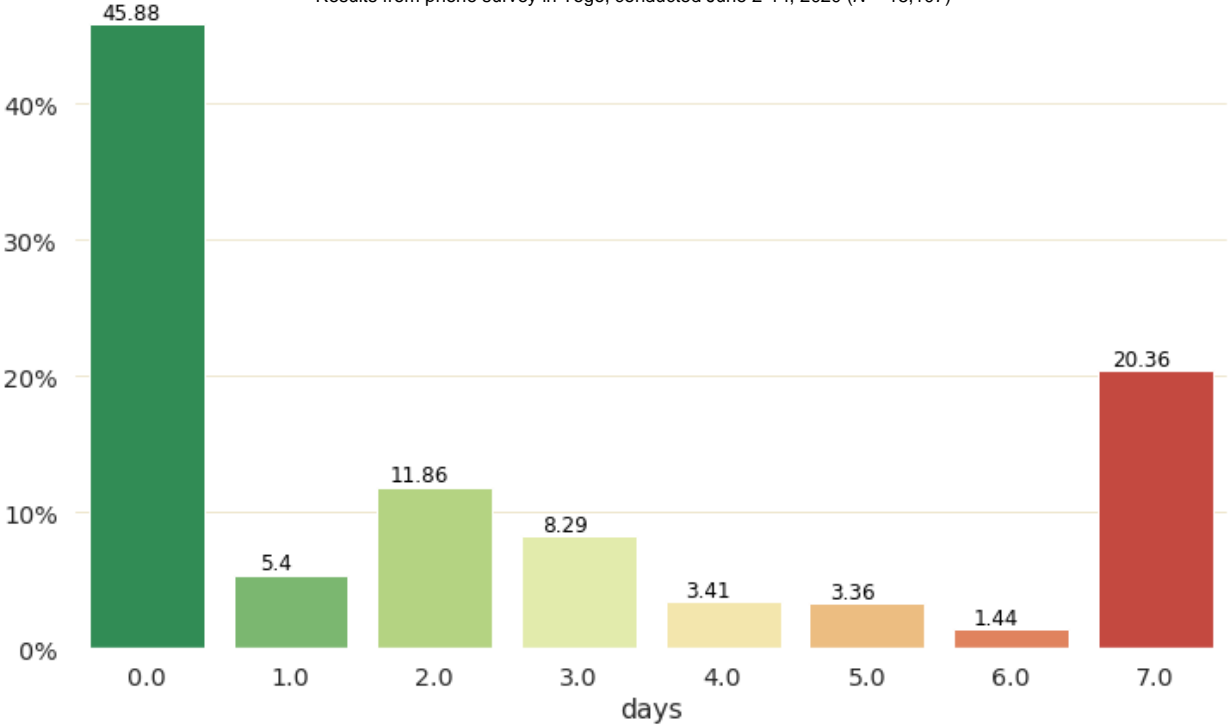
21. Ohlenberg, T. *AI in Social Protection – Exploring Opportunities and Mitigating Risks*. (Deutsche Gesellschaft für Internationale Zusammenarbeit (GIZ), 2020).
22. Abebe, R. *et al.* Narratives and Counternarratives on Data Sharing in Africa. in *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* 329–341 (Association for Computing Machinery, 2021). doi:10.1145/3442188.3445897.
23. Taylor, L. No place to hide? The ethics and analytics of tracking mobility using mobile phone data. *Environ Plan D* **34**, 319–336 (2016).
24. Necessary and Proportionate. International Principles on the Application of Human Rights to Communications Surveillance. (2013).
25. Taylor, L. No place to hide? The ethics and analytics of tracking mobility using mobile phone data. *Environ Plan D* **34**, 319–336 (2016).
26. Alaggan, M., Gambs, S., Matwin, S. & Tuhin, M. Sanitization of Call Detail Records via Differentially-Private Bloom Filters. in *Data and Applications Security and Privacy XXIX* (ed. Samarati, P.) vol. 9149 223–230 (Springer International Publishing, 2015).
27. Mir, D., Isaacman, S., Caceres, R., Martonosi, M. & Wright, R. DP-WHERE: Differentially private modeling of human mobility. in 580–588 (2013). doi:10.1109/BigData.2013.6691626.
28. Li, T., Sahu, A. K., Talwalkar, A. & Smith, V. Federated Learning: Challenges, Methods, and Future Directions. *IEEE Signal Processing Magazine* **37**, 50–60 (2020).
29. Oliver, N. *et al.* Mobile phone data for informing public health actions across the COVID-19 pandemic life cycle. *Science Advances* eabc0764 (2020) doi:10.1126/sciadv.abc0764.
30. Taylor, L. & Broeders, D. In the name of Development: Power, profit and the datafication of the global South. *Geoforum* **64**, 229–237 (2015).
31. Cate, F. H. & Mayer-Schönberger, V. Notice and consent in a world of Big Data. *International Data Privacy Law* **3**, 67–73 (2013).
32. Ioannidis, J. P. A. Informed Consent, Big Data, and the Oxymoron of Research That Is Not Research. *The American Journal of Bioethics* **13**, 40–42 (2013).
33. Noriega-Campero, A. *et al.* Algorithmic targeting of social policies: fairness, accuracy, and distributed governance. in *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency* 241–251 (ACM, 2020). doi:10.1145/3351095.3375784.
34. Kleinberg, J., Mullainathan, S. & Raghavan, M. Inherent Trade-Offs in the Fair Determination of Risk Scores. *arXiv:1609.05807 [cs, stat]* (2016).
35. Goodhart, C. *Monetary Relationships: A View from Threadneedle Street*. (University of Warwick, 1975).
36. Akerlof, G. A. The economics of ‘tagging’ as applied to the optimal income tax, welfare programs, and manpower planning. *The American economic review* **68**, 8–19 (1978).
37. Martinelli, C. & Parker, S. W. Deception and Misreporting in a Social Program. *Journal of the European Economic Association* **7**, 886–908 (2009).
38. Camacho, A. & Conover, E. Manipulation of Social Program Eligibility. *American Economic Journal: Economic Policy* **3**, 41–65 (2011).
39. Banerjee, A., Hanna, R., Olken, B. A. & Sumarto, S. *The (lack of) Distortionary Effects of Proxy-Means Tests: Results from a Nationwide Experiment in Indonesia*. <http://www.nber.org/papers/w25362> (2018) doi:10.3386/w25362.
40. Finkelstein, A. E-Z Tax: Tax Salience and Tax Rates. *The Quarterly Journal of Economics* **124**, 969–1010 (2009).

41. Björkegren, D., Blumenstock, J. E. & Knight, S. Manipulation-Proof Machine Learning. *arXiv:2004.03865 [cs, econ]* (2020).
42. Egger, D., Haushofer, J., Miguel, E., Niehaus, P. & Walker, M. W. *General Equilibrium Effects of Cash Transfers: Experimental Evidence from Kenya*. <https://www.nber.org/papers/w26600> (2019) doi:10.3386/w26600.
43. Cunha, J. M., De Giorgi, G. & Jayachandran, S. The Price Effects of Cash Versus In-Kind Transfers. *The Review of Economic Studies* **86**, 240–281 (2019).
44. Geng, X., Janssens, W., Kramer, B. & List, M. Health insurance, a friend in need? Impacts of formal insurance and crowding out of informal insurance. *World Development* **111**, 196–210 (2018).
45. Mobarak, A. M. & Rosenzweig, M. *Selling formal insurance to the informally insured*. <https://www.econstor.eu/handle/10419/59144> (2012).
46. Blumenstock, J. E., Cadamuro, G. & On, R. Predicting poverty and wealth from mobile phone metadata. *Science* **350**, 1073–1076 (2015).
47. Vanhoof, M., Reis, F., Ploetz, T. & Smoreda, Z. Assessing the Quality of Home Detection from Mobile Phone Data for Official Statistics. *Journal of Official Statistics* **34**, 935–960 (2018).
48. Chi, G., Lin, F., Chi, G. & Blumenstock, J. A general approach to detecting migration events in digital trace data. *PLoS One* **15**, e0239408 (2020).
49. Kane, T. J. & Staiger, D. O. *Estimating Teacher Impacts on Student Achievement: An Experimental Evaluation*. <https://www.nber.org/papers/w14607> (2008) doi:10.3386/w14607.
50. Chetty, R., Friedman, J. N. & Rockoff, J. E. Measuring the Impacts of Teachers I: Evaluating Bias in Teacher Value-Added Estimates. *American Economic Review* **104**, 2593–2632 (2014).
51. Gilraine, M., Gu, J. & McMillan, R. *A New Method for Estimating Teacher Value-Added*. <https://www.nber.org/papers/w27094> (2020) doi:10.3386/w27094.
52. Kalton, G. & Flores Cervantes, I. Weighting Methods. [http://lst-iiiep.iiiep-unesco.org/cgi-bin/wwwi32.exe/\[in=epidoc1.in\]/?t2000=019622/\(100\)](http://lst-iiiep.iiiep-unesco.org/cgi-bin/wwwi32.exe/[in=epidoc1.in]/?t2000=019622/(100)) **19**, (2003).
53. Buskirk, T. D. & Kolenikov, S. Finding Respondents in the Forest: A Comparison of Logistic Regression and Random Forest Models for Response Propensity Weighting and Stratification. *Survey Methods: Insights from the Field (SMIF)* (2015) doi:10.13094/SMIF-2015-00003.
54. Tiecke, T. G. *et al.* Mapping the world population one building at a time. *arXiv:1712.05839 [cs]* (2017).

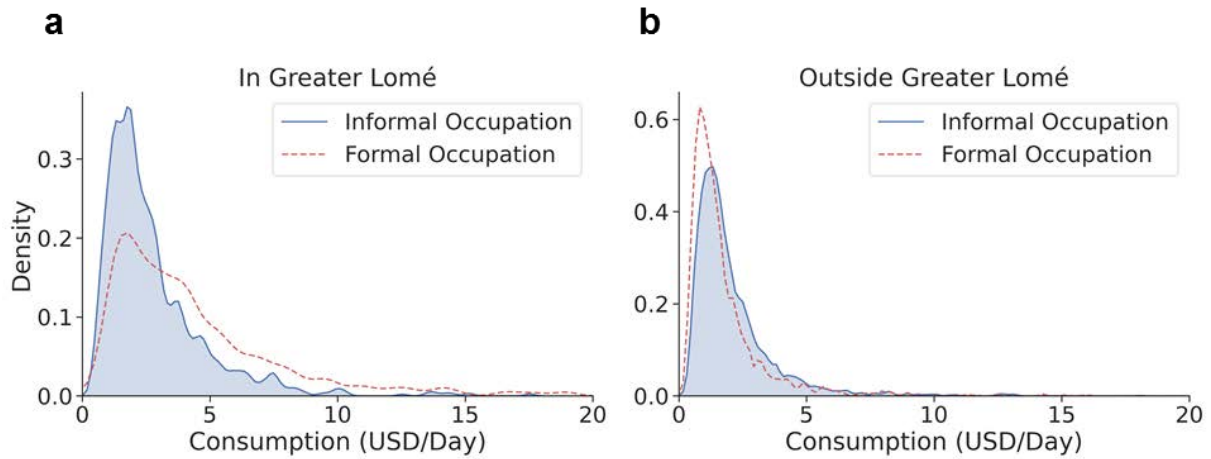
Supplementary Figures

*“In the past week, on how many days did you or someone in your household have to **reduce the number of meals eaten in a day?**”*

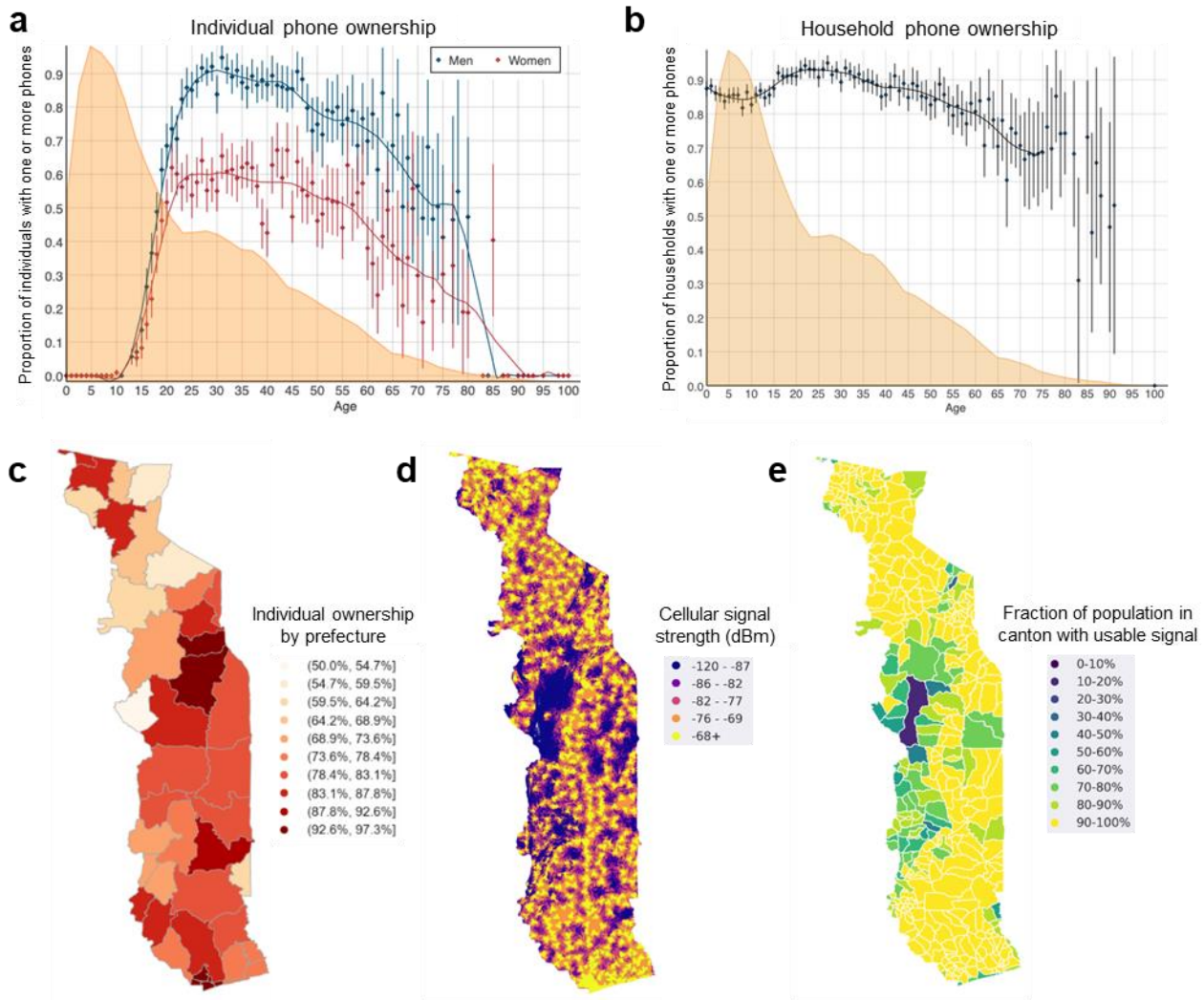
Results from phone survey in Togo, conducted June 2-14, 2020 (N = 15,107)



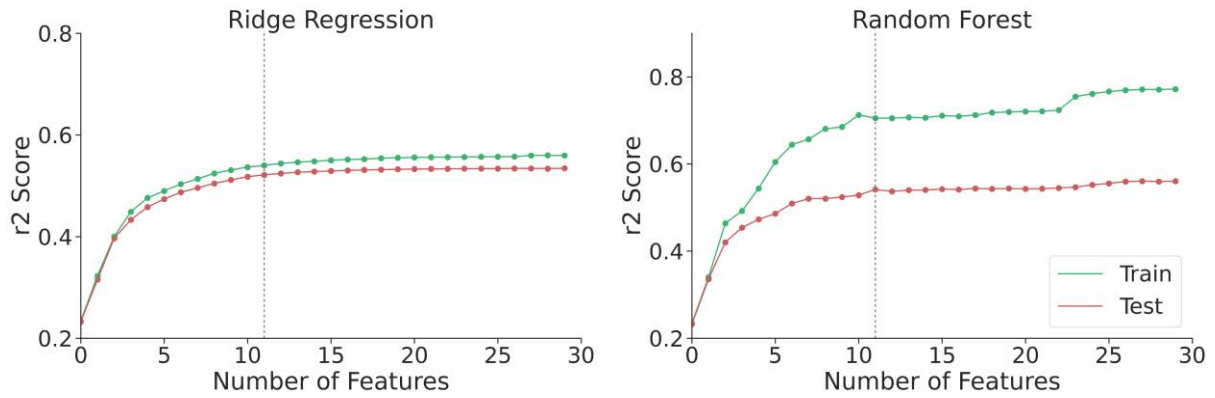
Supplementary Figure 1 | Food insecurity in Togo. In June 2020, we conducted a phone survey of 15,107 mobile phone owners in Togo. Survey weights are used to make responses representative of the population of mobile phone owners in Togo.



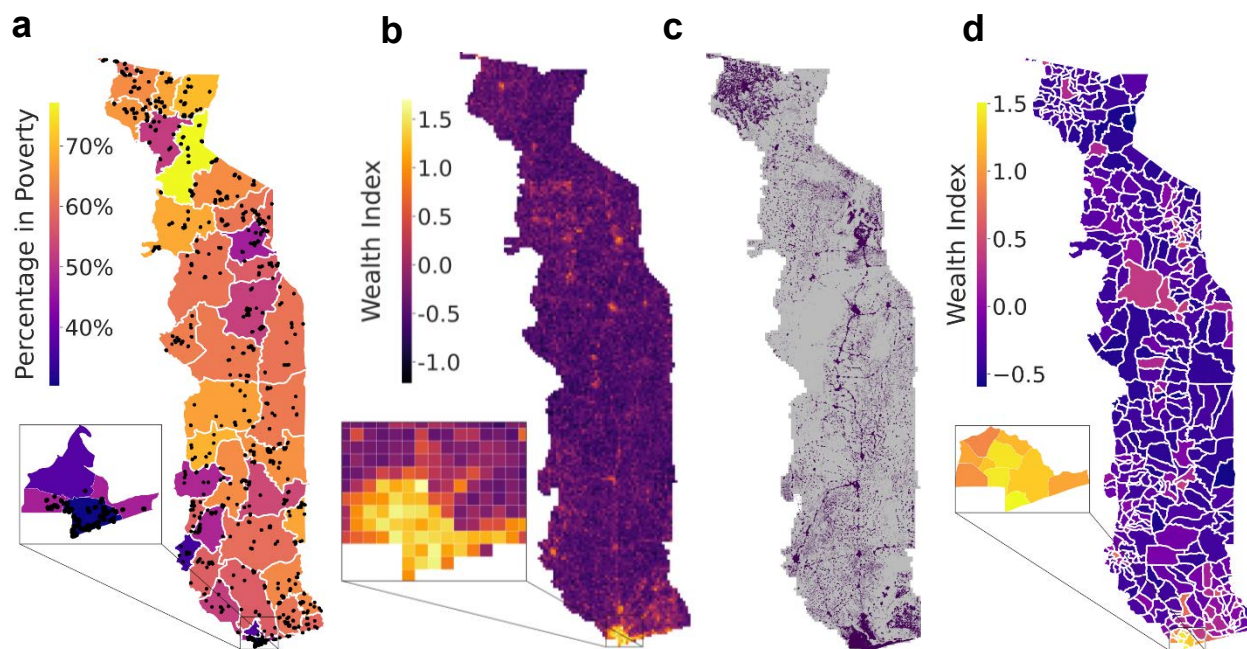
Supplementary Figure 2 | Wealth of formal vs. informal workers. Results based on analysis of nationally-representative household survey data collected by the Government of Togo in 2018-2019 ($N = 6,171$). Data is collected at the household-level, we assign a household-level informal occupation indicator if at least one of the adult household members is unemployed or employed in an informal occupation. See Methods, ‘Data Sources’.



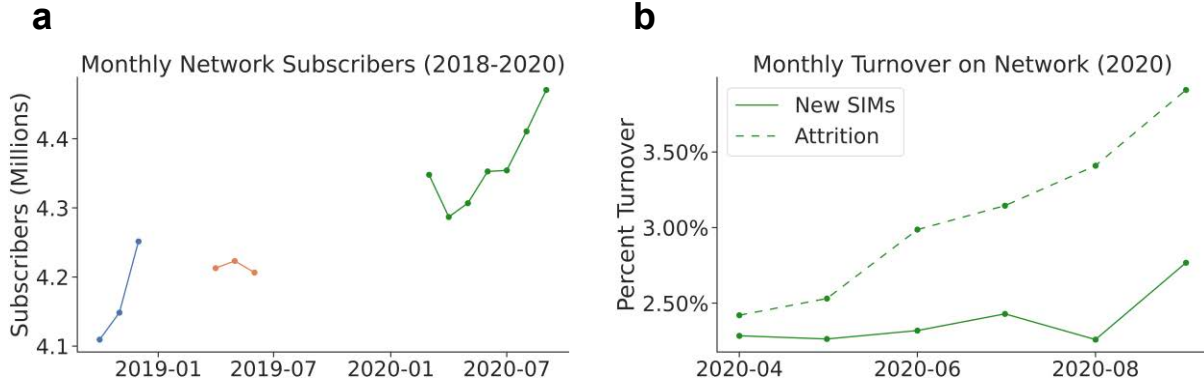
Supplementary Figure 3 | Mobile phone penetration and coverage in Togo. Based on nationally-representative household survey data collected in 2018-2019, we estimate **a)** the percentage of adults in Togo with one or more mobile phone, disaggregated by age and gender (the dots indicate the sample mean, while vertical bands indicate 95% confidence intervals derived from $N=27,483$ total individual survey responses); **b)** the percentage of households in Togo with one or more mobile phones, disaggregated by the age of the head of household (the dots indicate the sample mean, while vertical bands indicate 95% confidence intervals derived from $N=27,483$ total individual survey responses); and **c)** the percentage of individuals in each prefecture with one or more mobile phones. Using data on the location and signal strength of all cell towers in Togo, made available by Togocel (one of the two phone companies in Togo), we calculate **d)** the signal strength across Togo; and **e)** the fraction of the population in each canton with access to a usable signal, where signal greater than -86 dBm is generally considered usable, and sub-canton estimates of population density are derived from satellite imagery and downloaded from the Humanitarian Data Exchange⁵⁴.



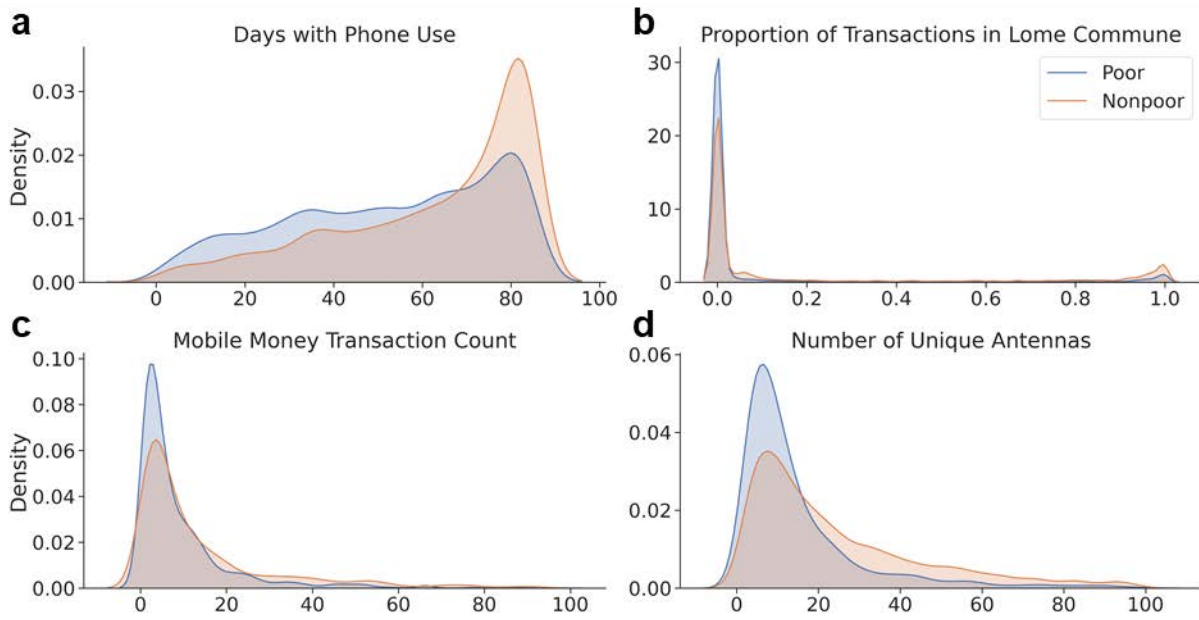
Supplementary Figure 4 | Selection of variables for proxy-means test. Each plot shows the accuracy (measured by r^2 score) of a proxy-means test using the most predictive feature subset of size K , where K is plotted on the x-axis. The left plot shows the accuracy obtained by a Ridge regression; the right plot shows the accuracy obtained by a random forest. Feature subsets are selected via stepwise forward selection.



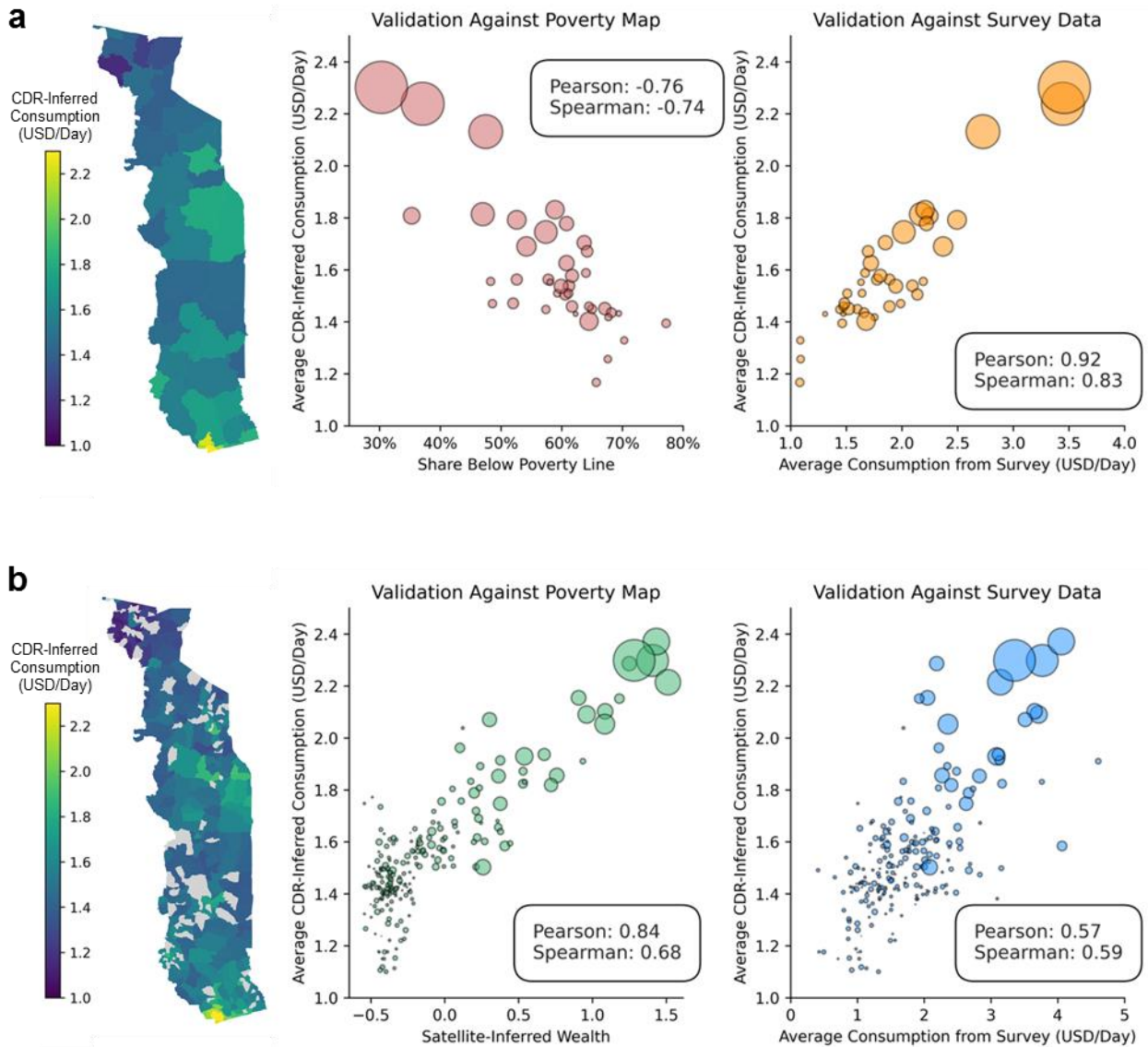
Supplementary Figure 5 | Poverty maps. (a) Prefecture (admin-2) poverty map inferred from 2017 field survey ($N = 26,902$), showing the percent of the population living under the poverty line by prefecture. Overlaid with locations of survey observations in black points. **(b)** High-resolution estimates of consumption derived from satellite imagery. **(c)** High-resolution estimates of population density derived from satellite imagery. **(d)** Canton (admin-3) poverty map inferred from satellite imagery by combining high-resolution consumption estimates and population density estimates to calculate weighted average consumption per canton.



Supplementary Figure 6 | Mobile phone network activity. (a) Count of unique subscribers making at least one outgoing transaction (call or text) on the mobile network in each month. October-December 2018 shown in blue, April-June 2019 in orange, and March-September 2020 in green. **(b)** Monthly turnover from the network in April-September 2020. New SIMs are quantified as the proportion of subscribers in each month whose first observed transaction is in that month. Attrition is quantified as the proportion of subscribers in each month who make no further outgoing transactions after that month. Note that we do not observe CDR in the months prior to March 2020, so we show results starting in April 2020 in Panel B; nonetheless a small proportion of the new SIMs in Panel B are inevitably due to sparsity in the CDR (that is, subscribers who placed a transaction prior to March 2020 that is not recorded in our dataset). Likewise, we do not observe CDR past December 2020, so a small part of the attrition measured in Panel B is due to sparsity in CDR transactions.



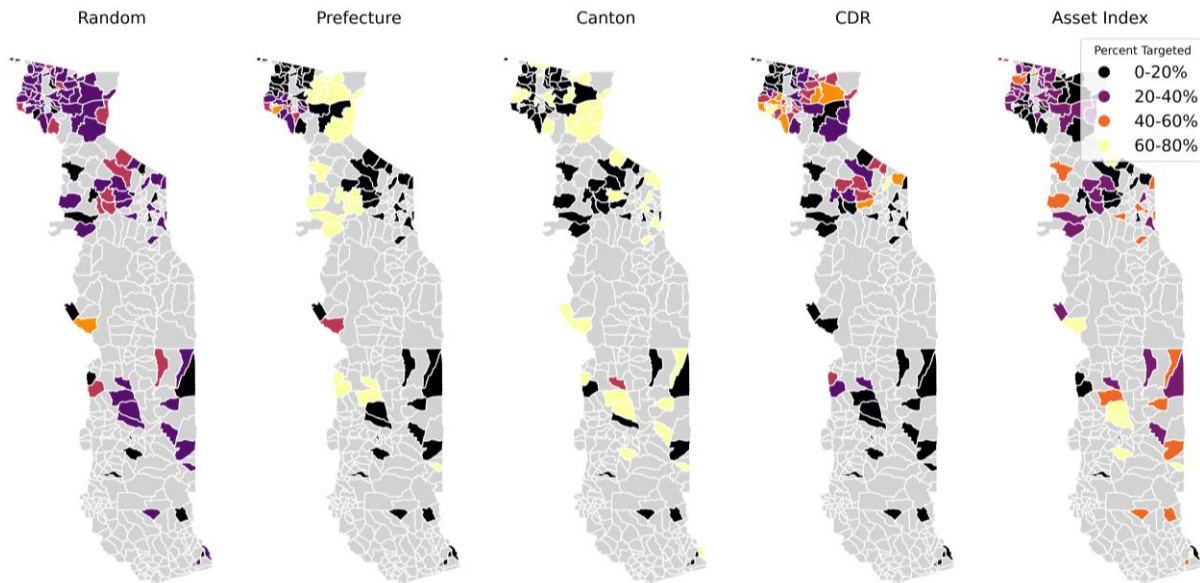
Supplementary Figure 7 | CDR features. Comparing the distribution for CDR features for those above and below the international poverty line (USD 1.90/day) in the 2018-2019 field survey dataset.



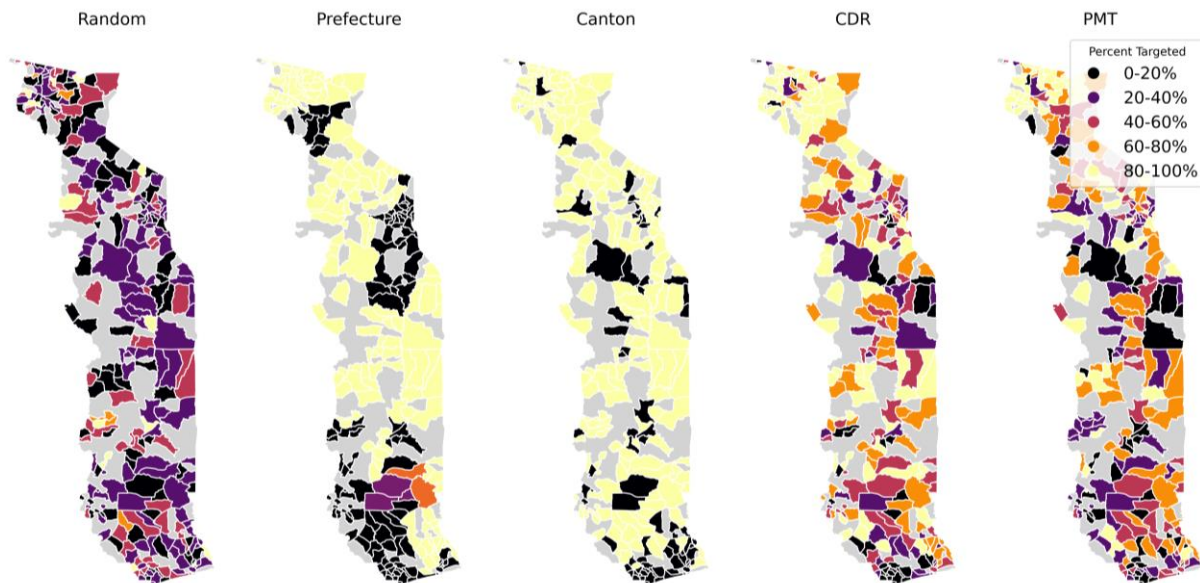
Supplementary Figure 8 | Spatial validation of phone-based poverty predictions. a) Map shows average phone-inferred consumption of subscribers in each *prefecture* (using CDR-based predictions trained on the 2018-19 in-person survey). Scatter plots compare average prefecture consumption, as derived from CDR (shown on y-axis), against two measures of poverty derived from the 2018-19 in-person survey (shown on x-axis): the share of people in the prefecture below the poverty line (middle plot), and the average consumption of households in the prefecture (right plot). **b)** Map shows average phone-inferred consumption of subscribers in each *canton* (cantons with no associated subscribers are shown in grey). Scatter plots compare average consumption per canton from the 2018-19 phone survey (evaluated across the 75% of all cantons in which there are observations in the 2018-19 field survey). Bubbles are sized by the number of subscribers assigned to each prefecture/canton.

a**Scenario 1: Targeting of the Novissi program in rural areas**

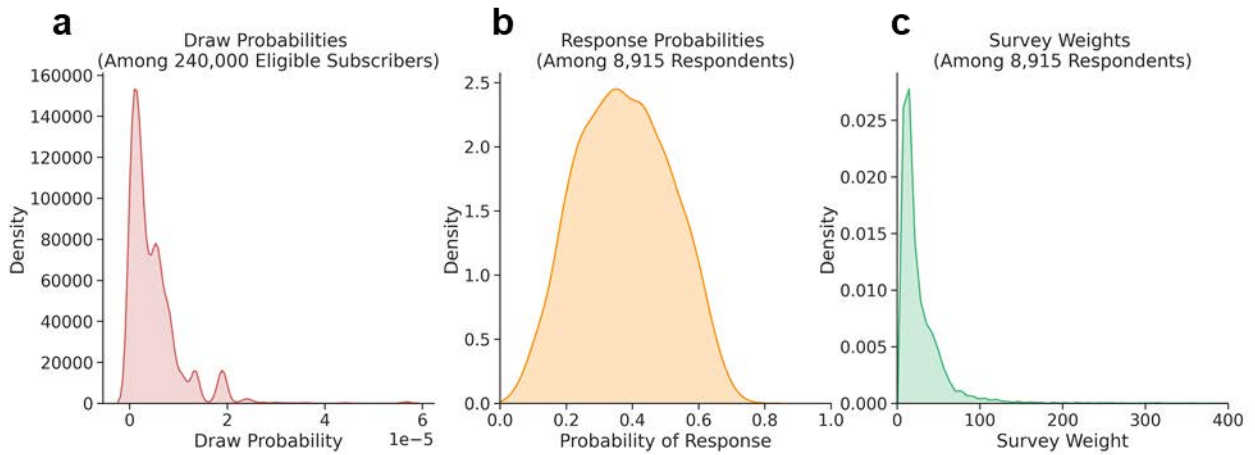
based on phone surveys collected in 2020

**b****Scenario 2: Targeting a hypothetical nationwide social assistance program**

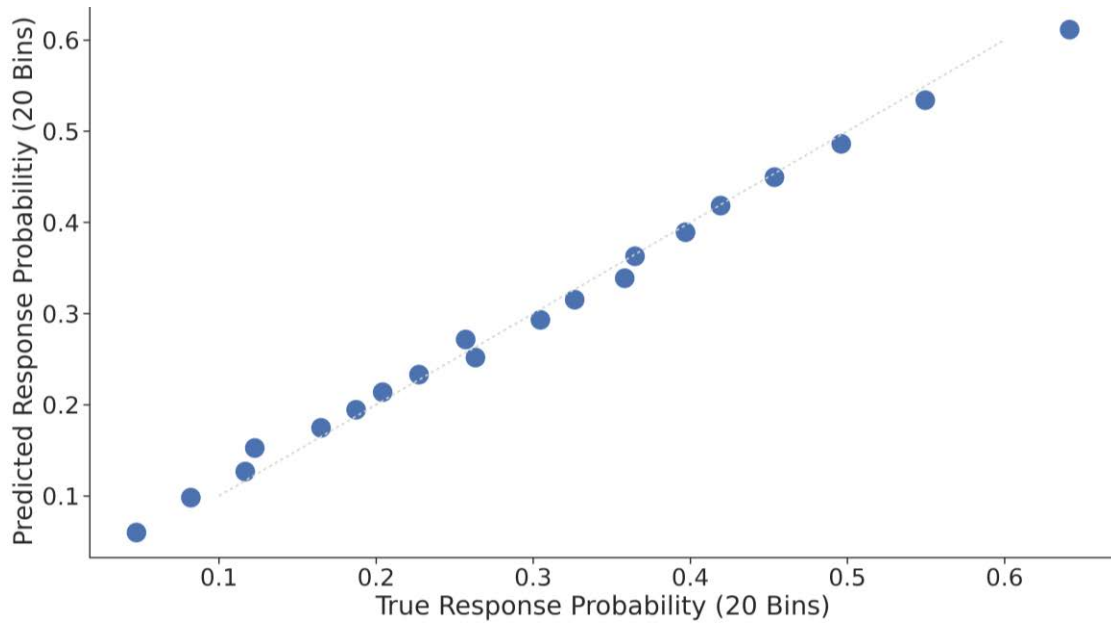
based on in-person surveys collected in 2018-2019



Supplementary Figure 9 | Share targeted by canton by different targeting methods. Panel A: Targeting share for the Novissi program in rural Togo, evaluated using individuals from the 2020 phone survey who report living in one of the 100 eligible cantons ($N = 6,745$). The respondent's self-reported canton and prefecture are used to color the map. Panel B: Targeting share for the hypothetical nationwide program, using data from the 2018-19 national household survey. Note that certain cantons have no observations in the 2018-2019 survey; these are shown in grey in Panel B. Cantons outside of the 100 poorest are shown in grey in Panel A.



Supplementary Figure 10 | Distribution of sample weights for 2020 phone survey. Panel A: Distribution of draw probabilities among subscribers eligible for the survey. Panel B: Distribution of response probabilities for observations included in the final survey dataset, based on the response prediction model. Panel C: Distribution of sample weights (product of the inverse of the draw probability and the inverse of the response probability) for observations included in the final survey dataset.



Supplementary Figure 11 | Calibration of response probabilities for 2020 phone survey. We compare the predicted probability of response (y-axis, binned into 20 quantiles) to the realized probability of response (x-axis, again binned into 20 quantiles) to confirm that the response prediction model is well-calibrated.

Supplementary Tables

Targeting a hypothetical nationwide program – but only in rural areas

Based on 2018 Phone Survey Restricted to Rural Areas ($N = 2,306$)

	Spearman	AUC	Accuracy	Precision & Recall
<i>Panel A: Targeting methods considered by the Government of Togo in 2020</i>				
Prefecture (Admin-2 regions)	0.16 (0.023)	0.57 (0.011)	64% (0.97%)	37% (1.67%)
Canton (Admin-3 regions)	0.19 (0.025)	0.59 (0.013)	63% (0.98%)	36% (1.69%)
Phone (Expenditures)	0.15 (0.024)	0.59 (0.012)	63% (1.05%)	36% (1.81%)
Phone (Machine Learning)	0.30 (0.023)	0.65 (0.012)	67% (1.00%)	43% (1.73%)
<i>Panel B: Common alternative targeting methods that could not be implemented in Togo in 2020</i>				
Asset Index	0.36 (0.023)	0.68 (0.011)	67% (1.01%)	44% (1.74%)
PPI	0.55 (0.017)	0.77 (0.009)	72% (1.07%)	52% (1.84%)
PMT	0.61 (0.016)	0.80 (0.007)	73% (1.06%)	54% (1.84%)
Rural PMT	0.52 (0.018)	0.75 (0.008)	72% (1.02%)	51% (1.75%)
<i>Panel C: Additional counterfactual targeting methods that were feasible in Togo in 2020</i>				
Random	0.00 (0.024)	0.50 (0.012)	59% (1.04%)	29% (1.79%)
Occupation (Novissi)	-0.13 (0.024)	0.44 (0.011)	54% (0.89%)	21% (1.53%)
Occupation (Optimal)	0.31 (0.023)	0.63 (0.010)	63% (0.62%)	37% (1.06%)

Supplementary Table 1 | Performance of targeting the hypothetical national program, when restricted to rural areas. Analysis is similar to that presented in the last four columns of Table 1, but analysis is restricted to the 2,306 survey respondents (of the 4,171 total) who live in rural areas.

Asset	Magnitude (2018-2019 Field Survey)	Magnitude (2020 Phone Survey)
Electricity access	0.38	
Toilet	0.37	0.41
TV	0.35	
Electricity grid	0.35	
Garbage disposal	0.33	
Waste disposal	0.33	
Iron	0.26	0.06
Radio	0.20	0.23
Clean water (wet season)	0.16	
Clean water (dry season)	0.16	
Refrigerator	0.12	0.02
Walls	0.12	
Floor	0.11	
Mobile phone	0.11	
Water disposal	0.10	
Motorcycle	0.10	0.88
Computer	0.09	0.02
Roof	0.08	
Stove	0.07	0.06
Car	0.06	0.00
Tablet	0.01	0.00
Air conditioner	0.01	0.00
House	0.00	
Electricity (offgrid)	0.00	

Supplementary Table 2 | Asset-based wealth index. Magnitude of first principal component for 2018-2019 field survey and 2020 phone survey.

Feature	β	Feature (continued)	β
Car	2.77	HHW Education 4	-0.18
Stove	1.77	Pref. Lacs	-0.18
Refrigerator	1.32	Pref. Sotouboua	-0.18
HHH Education 8	1.12	Pref. Kloto	-0.21
HHH Education 9	0.91	HHW Education 6	-0.21
HHH Hospitalization	0.81	Pref. Kpele	-0.21
Iron	0.63	Pref. Bas-Mono	-0.23
HHH Education 3	0.55	Pref. Lome Commune	-0.23
TV	0.50	Pref. Danyi	-0.24
All children in school	0.48	Pref. Yoto	-0.26
Pref. Cinkasse	0.39	Pref. Agoe-Nyive	-0.27
Pref. Tchamba	0.33	HHH Education 5	-0.27
Toilet	0.26	No children in school	-0.31
HHH Education 7	0.17	Pref. Assoli	-0.32
Pref. Est-Mono	0.14	Pref. Kpendjal-Ouest	-0.33
HHW Education 0	0.12	Pref. Zio	-0.33
Pref. Tchaoudjo	0.09	Pref. Amou	-0.34
Pref. Bassar	0.09	HHW Education 3	-0.34
Pref. Haho	0.07	Pref. Plaine du Mo	-0.34
Pref. Dankpen	0.04	Pref. Anie	-0.34
Pref. Moyen-Mono	-0.03	Pref. Tandjoare	-0.35
Pref. Oti-Sud	-0.06	Pref. Binah	-0.37
Pref. Oti	-0.08	Pref. Ave	-0.39
Pref. Wawa	-0.11	Pref. Keran	-0.41
Pref. Vo	-0.11	Pref. Kpendjal	-0.46
Pref. Ogou	-0.12	HHW Education 2	-0.50
Pref. Tone	-0.14	Pref. Kozah	-0.51
Pref. Agou	-0.15	HHH Education 2	-0.57
Pref. Akebou	-0.17	Pref. Blitta	-0.61
HHW Education 1	-0.17	HHH Education 1	-0.63
Some children in school	-0.17	Pref. Golfe	-0.68
Number of children	-0.17	Pref. Doufelgou	-0.75

Supplementary Table 3 | Proxy means test. Weights for linear model, trained on 2018-2019 field survey ($N = 6,171$).

Feature	β	Feature (continued)	β
Refrigerator	0.38	Pref. Akebou	0.03
HHH Hospitalization	0.32	Pref. Ogou	0.02
Motorcycle	0.31	Pref. Ave	-0.01
TV	0.28	Pref. Moyen-Mono	-0.03
Pref. Vo	0.26	Number of children	-0.08
Computer	0.24	Pref. Plaine du Mo	-0.08
Pref. Tchamba	0.21	Pref. Est-Mono	-0.10
Garbage removal	0.17	Pref. Dankpen	-0.12
Pref. Wawa	0.17	Pref. Binah	-0.13
Toilet	0.16	Pref. Tchaoudjo	-0.13
Pref. Kloto	0.16	Pref. Cinkasse	-0.14
Pref. Haho	0.16	Pref. Oti-Sud	-0.15
Pref. Yoto	0.14	Pref. Anie	-0.15
Pref. Bas-Mono	0.14	Pref. Oti	-0.18
Pref. Golfe	0.14	Pref. Kozah	-0.19
Pref. Kpele	0.14	Pref. Tone	-0.22
Pref. Lacs	0.13	Pref. Assoli	-0.22
Floor of solid materials	0.10	Pref. Blitta	-0.23
Pref. Zio	0.10	No children in school	-0.24
Pref. Lome Commune	0.09	Some children in school	-0.28
Pref. Agou	0.09	Pref. Doufelgou	-0.29
Roof of solid materials	0.08	Pref. Kpendjal-Ouest	-0.32
Pref. Bassar	0.08	Pref. Keran	-0.33
Pref. Amou	0.06	Pref. Kpendjal	-0.35
Pref. Danyi	0.06	Pref. Tandjoare	-0.40
Pref. Soutoubua	0.05		

Supplementary Table 4 | Rural-specific proxy means test. Weights for linear model, trained on 2018-2019 phone survey restricted to rural areas ($N = 3,895$).

	<i>2018-2019 Field Survey (N=6,171)</i>			<i>2020 Phone Survey (N=8,915)</i>	
	Consumption	Proportion	N	Proportion	N
Intellectual Professions	\$4.11 (3.55)	7%	277	7%	577
Intermediate Professions	\$3.95 (3.56)	5%	197	3%	264
Administrators	\$3.89 (3.57)	1%	32	0%	16
Managers and Directors	\$3.70 (3.03)	3%	106	0%	36
Unemployed/Unknown	\$3.19 (2.44)	8%	339	3%	275
Direct Services and Merchants	\$2.75 (2.11)	23%	940	28%	2,111
Industry/Artisans	\$2.47 (1.83)	15%	587	12%	1,026
Military Professions	\$2.45 (1.25)	0%	17	1%	26
Elementary Professions	\$2.21 (1.83)	2%	65	3%	249
Factory Workers	\$2.17 (1.44)	7%	267	2%	165
Agricultural Professions	\$1.53 (0.94)	29%	1,744	41%	4,170

Supplementary Table 5 | Occupation categories. Average daily per capita consumption per occupation category, with counts by category, separately for the 2018-2019 field survey and 2020 phone survey. Occupation categories for the 2018-2019 survey are for the household head, for the 2020 survey are for the individual respondent.

	2018-2019 National Household Survey					2020 Phone Survey
	Full Survey	Phone Number	No Phone Number	Phone Number, Matched	Phone Number, Unmatched	Full Survey
Consumption	2.39 (2.41)	2.56 (2.38)	1.75 (2.41)	2.59 (2.42)	2.21 (1.78)	<i>[data not available]</i>
PMT	2.10 (1.43)	2.22 (1.47)	1.65 (1.16)	2.23 (1.47)	2.03 (1.38)	1.62 (0.72)
Occupation (% Formal)	56.42% (49.59%)	51.98% (49.96%)	72.94% (44.43%)	51.28% (49.99%)	59.63% (49.08%)	66.54% (47.19%)
% Rural	51.93% (49.96%)	45.17% (49.77%)	77.12% (42.01%)	43.79% (49.61%)	60.17% (48.97%)	96.19% (19.15%)
% Women	28.15% (44.98%)	23.61% (42.47%)	45.07% (49.76%)	23.43% *42.36%	25.63% (43.68%)	23.27% (42.25%)
Age	43.97 (14.43)	41.96 (13.19)	51.26 (16.28%)	41.97 (13.15%)	41.84 *(13.71%)	33.20 (11.90)
<i>N</i>	6,089	4,571	1,518	4,171	400	8,915

Supplementary Table 6 | Summary statistics for two survey datasets. Means and standard deviations for key outcomes in the 2018-2019 national household survey ($N = 6,089$) and 2020 phone survey concentrated in the 100 poorest cantons ($N = 8,915$). For the 2018-2019 national household survey, we break down the sample into two groups: households that provided enumerators with a phone numbers ($N = 4,571$) and those that do not ($N = 1,518$). We further break down the sample providing a phone number into two groups: households for which the phone number appears in data obtained from the mobile network operators ($N = 4,171$) and those for which it does not ($N = 400$). For the 2018-19 phone survey, occupation, gender, and age are assigned based on the head of household; for the 2020 phone survey they are assigned based on the respondent.

Targeting a hypothetical nationwide program – with PMT as ground truth
Based on 2018-2019 National Household Survey ($N = 4,171$)

	Spearman	AUC	Accuracy	Precision & Recall
<i>Panel A: Targeting methods considered by the Government of Togo in 2020</i>				
Prefecture (Admin-2 regions)	0.50 (0.014)	0.73 (0.006)	72% (0.73%)	51% (1.25%)
Canton (Admin-3 regions)	0.54 (0.013)	0.73 (0.006)	74% (0.70%)	55% (1.22%)
Phone (Expenditures)	0.31 (0.017)	0.64 (0.008)	66% (0.75%)	41% (1.30%)
Phone (Machine Learning)	0.56 (0.014)	0.78 (0.006)	73% (0.73%)	54% (1.25%)
<i>Panel B: Common alternative targeting methods that could not be implemented in Togo in 2020</i>				
Asset Index	0.68 (0.010)	0.82 (0.005)	77% (0.73%)	60% (1.26%)
PPI	0.74 (0.009)	0.86 (0.005)	81% (0.68%)	67% (1.18%)
<i>Panel C: Additional counterfactual targeting methods that were feasible in Togo in 2020</i>				
Random	0.00 (0.020)	0.50 (0.011)	59% (0.78%)	30% (1.34%)
Occupation (Novissi)	-0.13 (0.019)	0.44 (0.009)	54% (0.51%)	21% (0.88%)
Occupation (Optimal)	0.50 (0.015)	0.73 (0.006)	76% (0.71%)	59% (1.22%)

Supplementary Table 7 | Performance of targeting the hypothetical national program, with PMT as ground truth. Analysis is similar to that presented in the last four columns of Table 1, but with the PMT as ground truth instead of consumption.

Targeting Novissi in rural Togo – with rural PMT as ground truth

Based on 2020 Phone Survey ($N = 8,915$)

	Spearman	AUC	Accuracy	Precision & Recall
<i>Panel A: Targeting methods considered by the Government of Togo in 2020</i>				
Prefecture (Admin-2 regions)	0.31 (0.023)	0.65 (0.011)	65% (1.02%)	40% (1.76%)
Canton (Admin-3 regions)	0.19 (0.025)	0.60 (0.012)	62% (1.03%)	34% (1.78%)
Phone (Expenditures)	0.16 (0.023)	0.58 (0.012)	61% (1.02%)	33% (1.76%)
Phone (Machine Learning)	0.41 (0.022)	0.69 (0.012)	68% (0.99%)	46% (1.70%)
<i>Panel B: Common alternative targeting methods that could not be implemented in Togo in 2020</i>				
Asset Index	0.46 (0.021)	0.71 (0.011)	68% (0.99%)	46% (1.71%)
<i>Panel C: Additional counterfactual targeting methods that were feasible in Togo in 2020</i>				
Random	0.00 (0.023)	0.50 (0.012)	59% (0.99%)	29% (1.70%)
Occupation (Novissi)	-0.12 (0.024)	0.45 (0.011)	55% (0.93%)	23% (1.60%)
Occupation (Optimal)	0.26 (0.022)	0.61 (0.010)	65% (0.66%)	40% (1.14%)

Supplementary Table 8 | Performance of targeting Novissi in rural Togo, when using the rural-specific PMT as ground truth. Analysis is similar to that presented in the first four columns of Table 1, but with the rural-specific PMT (as described in Methods, ‘Survey Data’) as ground truth.

	Targeting Novissi in rural Togo Based on 2020 Phone Survey (<i>N</i> = 8,915)				Hypothetical nationwide program Based on 2018-2019 Field Survey (<i>N</i> = 4,171)			
	Spearman	AUC	Accuracy	Precision & Recall	Spearman	AUC	Accuracy	Precision & Recall
Prefecture (Survey-recorded)	0.30 (0.017)	0.64 (0.008)	65% (0.87%)	39% (1.51%)	0.34 (0.017)	0.66 (0.008)	68% (0.74%)	45% (1.27%)
Canton (Survey-recorded)	0.19 (0.019)	0.59 (0.009)	61% (0.78%)	33% (1.35%)	0.39 (0.016)	0.68 (0.008)	70% (0.71%)	48% (1.23%)
CDR Prefecture (Phone-inferred)	0.23 (0.016)	0.61 (0.008)	63% (0.76%)	36% (1.32%)	0.27 (0.017)	0.63 (0.008)	67% (0.74%)	44% (1.40%)
CDR Canton (Phone-inferred)	0.12 (0.021)	0.56 (0.011)	58% (0.83%)	28% (1.43%)	0.31 (0.017)	0.65 (0.008)	69% (0.73%)	47% (1.27%)
Phone (Machine Learning)	0.38 (0.017)	0.70 (0.009)	69% (0.87%)	47% (1.18%)	0.45 (0.015)	0.73 (0.007)	71% (0.74%)	50% (1.26%)

Supplementary Table 9 | Geographic targeting with phone-inferred location. First two rows and final row replicate the results shown in Table 1. We add two additional counterfactual geographic targeting approaches based on location information derived from mobile phone data: targeting based on the average wealth of their home prefecture (row 3) or of their home canton (row 4). Home prefectures and cantons are inferred from outgoing mobile phone transactions as described in Supplementary Methods section 4; the poverty of associated with each prefecture and canton is taken from the poverty maps shown in Supplementary Figure 5.

	Prefecture-level	Canton-level
Survey \leftarrow \rightarrow Voter	90.08%	69.77%
Survey \leftarrow \rightarrow Phone Data	70.08%	46.56%
Voter \leftarrow \rightarrow Phone Data	67.48%	44.89%

Supplementary Table 10 | Correlation between sources of location data in 2020 phone survey. Correlation between the three sources of home location data available for observations in the 2020 phone survey: self-reported location collected in a survey, voter location recorded at the time of voter registration, and home location inferred from phone data. Each entry represents the percentage of observations (without sample weights applied) for which the two datasets agree on the respondent's location. Percentages are taken among the population ($N = 4,515$) for whom all three data sources are available (that is, individuals who were surveyed, whose phone numbers were registered for the rural Novissi program so that the canton and prefecture associated with their voter ID are included in Novissi administrative data, and who place at least one outgoing call between March to September 2020 so that their phone number is tied to a home prefecture and canton). This analysis cannot be carried out for the 2018-2019 field survey as fewer than 15% of the phone numbers collected in the survey registered for the rural Novissi program.

	Share of phone transactions made from home prefecture (inferred from CDR)	Share of phone transactions made from home prefecture (self-reported in survey)
<i>Panel A: 2018-2019 national household survey and April-June 2019 CDR</i>		
Mean (and standard deviation)	75.46% (31.90%)	62.00% (40.05%)
Median	91.18%	81.84%
Mode	100.00%	100.00%
N	3,459,308	3,992
<i>Panel A: 2020 phone survey and March-September 2020 CDR</i>		
Mean (and standard deviation)	85.32% (18.78%)	68.00% (36.79%)
Median	94.00%	87.16%
Mode	100.00%	100.00%
N	5,615,393	8,183

Supplementary Table 11 | Percentage of mobile phone activity initiated from a subscriber’s home prefecture. Table indicates the fraction of outgoing calls and text messages that are routed through a cell tower in the subscriber’s home prefecture. In the first column, “home location” is inferred from the subscriber’s CDR as described in Appendix B; in the second column, “home location” is recorded during a survey with the respondent. Panel A: results based on analysis from 2019, using CDR from three months in 2019 in the first column ($N = 3,459,308$), and survey respondents with known GPS coordinates from the 2018-2019 field survey in the second column ($N = 3,992$). Panel B: results based on analysis from 2020, using CDR from 7 months in 2020 in the left column ($N = 5,615,393$), and survey respondents with self-reported prefectures in the 2020 phone survey in the right column ($N = 8,183$).

Group	Response Rate	N
<i>Panel A: Previous Novissi registration</i>		
Registered	37.82%	15,402
Unregistered	25.61%	14,085
<i>Panel B: Phone-inferred region</i>		
Lomé Commune	35.45%	189
Maritime	40.83%	1,254
Plateaux	30.17%	3,627
Centrale	31.91%	702
Kara	35.31%	6,582
Savanes	30.42%	17,034
<i>Panel C: Phone-predicted poverty (USD/day)</i>		
<\$1.32	33.50%	7,372
\$1.32-\$1.42	33.55%	7,372
\$1.42-\$1.57	30.10%	7,371
\$1.57+	30.79%	7,372
<i>Panel D: Phone expenditures (USD/day)</i>		
<\$0.03	22.56%	7,372
\$0.03-\$0.08	28.15%	7,372
\$0.08-\$0.21	34.82%	7,371
\$0.21+	42.66%	7,372

Supplementary Table 12 | Response rates for 2020 phone survey. Response rate disaggregated by four dimensions: registration to a previous Novissi program (Panel A), region of Togo inferred from location of mobile phone transactions (Panel B), daily consumption inferred from mobile phone activity and machine learning (Panel C), and daily phone expenditures (Panel D).

Feature	Importance
Registered to previous Novissi program	15
Togocom subscriber	13
% nocturnal calls	10
% in Kpendjal	9
Active days	9
Mean balance of contacts	8
Median interactions per contact	7
Median time between calls (weekdays)	7
Active days (weekend)	6
Minimum time between calls	6

Table S13 | Feature importances for response reweighting model for 2020 phone survey. As described in Supplementary Methods section 5, the gradient boosting ensemble model is trained to predict the probability of response for a phone number drawn for the 2020 phone survey on the basis of pre-survey observable covariates (from CDR and previous Novissi registrations). Feature importance is calculated based on the total number of times a feature is split upon in the prediction ensemble.