

A respondent-driven method for mapping small agricultural plots using tablets and high resolution imagery

Yuta J. Masuda<sup>1</sup>, Jonathan R.B. Fisher<sup>2</sup>, Wei Zhang<sup>3</sup>, Carolina Castilla<sup>4</sup>, Tim Boucher<sup>1</sup>, Genowefa Blundo<sup>5</sup>

<sup>1</sup> Global Science, The Nature Conservancy, 4245 Fairfax Dr #100, Arlington, VA 22203

<sup>2</sup> Center for Sustainability Science, The Nature Conservancy, 4245 Fairfax Dr #100, Arlington, VA 22203

<sup>3</sup> Environment, Production and Technology Division, International Food Policy Research Institute, 2033 K Street NW, Washington DC, 20006 USA

<sup>4</sup> Department of Economics, Colgate University, 13 Oak Drive Hamilton, NY 13346

<sup>5</sup> Département Environnements et sociétés, La recherche agronomique pour le développement, Campus international de Baillarguet TA C-DIR / B 34398 Montpellier, France

Keywords: agriculture, spatial data, household surveys, smallholder, methodology

**Abstract**

*Agricultural statistics on small farms are critical but still suffer from selection bias and are time consuming and costly. Less burdensome and reliable methods are needed. We report on a scalable method using a respondent's knowledge, high resolution imagery, and tablet devices to draw spatially explicit plot boundaries. We find the method may work best with respondents that own their plots and farmers, and for smaller plots (<1 hectare). We also find incongruence between survey questions and spatially-derived data, indicating the importance of incorporating spatial data to verify responses about plot characteristics.*

DRAFT

## Introduction

With over 375 million households relying on small plots for food and livelihoods (FAO, 2014), developing cost effective and reliable methods for collecting agricultural statistics is critical for tackling sustainable development challenges. Plot boundary measurements in particular are important because it enables integrating socioeconomic and productive activity data from household surveys with a variety of spatiotemporal remotely sensed data on land cover, species, climate, soil, hydrology, population density, and other biophysical features. Recent work by Donaldson and Storeygard (2016) highlight the uses and utility of incorporating spatial data into economics and for evaluating challenges in development. But current methods for gathering spatially explicit plot boundaries (herein referred to as just plot boundaries) for small agricultural plots can be time consuming and costly, resulting in poor data quality and narrow sectoral focus (Carletto *et al.*, 2015b). We report on a practical and scalable method for gathering small plot boundaries. Our method utilizes high-resolution satellite imagery and the respondent's knowledge about their own landscape to draw plot boundaries on tablet devices. We compare plot characteristics from survey data and satellite imagery data and find incongruence. We demonstrate the utility and limitations of the method and make recommendations on operationalizing this method for improved integration of spatially explicit information to household survey data.

Gathering spatially explicit agricultural survey data are critical for advancing knowledge on land governance, agricultural productivity, food security, environmental degradation, as well as the human-nature interactions. For instance, land, especially for the rural poor, is often the primary and most valuable asset. Accurate land area measurement is critical to many agricultural metrics, such as yield, productivity of inputs, landholding, and extent of land fragmentation (Carletto *et al.*, 2015b). It can also help identify land inequality, which is a critical factor in economic growth and civil conflict, and human development (Deininger and Squire, 1998; Macours, 2011; André and Platteau, 1998; Baten and Juif, 2013). Spatial data can augment household survey data with rich and potentially novel information. For instance, measuring multiple observations over vast areas is relatively simple and cost effective, and information on environmental changes may be timelier and more accurate compared to household surveys (e.g., rate and extent of deforestation (Hansen *et al.*, 2013)). Knowing where household surveys have a comparative advantage (e.g., demographics) and where they are limited is important when integrating these data. Technical advances, such as the availability, reliability, and cost of using Global Position Systems (GPS) devices have reduced the cost of using GPS devices to gather plot boundaries (Goldstein and Udry, 1999; Carletto *et al.*, 2015b).<sup>1</sup> For instance, Schøning *et al.* (2005) found land measurements via tape and

---

<sup>1</sup> For instance, the Tanzania Living Standards Measurement Study-Integrated Surveys on Agriculture have collected GPS measurements for 25 percent of sampled households, although plots that were not within an hour by any mode of transportation were excluded.

compass (the traditional method) took three times longer compared to GPS measurements (for a list of other methods see Appendix A). But GPS measurements are not a panacea for agricultural statistics. First, like the tape and compass method, GPS measurements still require enumerators to physically identify and walk around the plot perimeter. As a result, plots outside some predetermined minimum distance to the household are often excluded, and this can lead to systematic plot selection bias (Kilic *et al.*, 2017). GPS measurements also still suffer from technical limitations. Tree canopy cover, weather conditions, and plot size and slope can hinder accuracy of GPS measurements (Keita and Carfagna, 2009, 2010; Fermont and Benson, 2011). The accuracy of GPS devices also falls significantly for small agricultural plots (<0.5 hectares, or 1.2 acres) (Keita *et al.*, 2010; Fermont and Benson, 2011), which is often the primary population of interest to researchers and practitioners interested in sustainable development.

GPS devices have seen limited adoption despite its increasing affordability, and the use of survey questions to gather self-reported plot data is still common. But there is a clear cost to relying only on self-reported survey questions: lower reliability and accuracy (Goldstein and Udry, 1999; De Groote and Traorè, 2005; Carletto *et al.*, 2013; Carletto *et al.*, 2015a; Carletto *et al.*, 2015b; Arthi *et al.*, 2017). In particular, measurement error in plot size estimation may be driven by the difficult mental calculus respondents must perform when estimating plot size (Carletto *et al.*, 2013). Additionally, self-reported plot measurements often suffer from “heaping”, which occurs when respondents attempt to estimate plot sizes by rounding to discrete values (e.g., 0.31 hectares to 0.5 hectares). Enumerators also may be a source of measurement error, as they may be required to convert multiple land measurement units to a common unit in large-scale surveys (Fermont and Benson, 2011; Carletto *et al.*, 2015b). While measurement error may seem trivial, it has led to substantial debates in the agricultural development literature, such as the inverse farm size-productivity relationship (Benjamin, 1995; Bhalla and Roy, 1988; Lamb, 2003; Barrett *et al.*, 2010; Carletto *et al.*, 2013; Holden and Fisher, 2013; Ali and Deininger, 2015).

An alternative method to gathering data on plot boundaries, and by extension numerous plot characteristics, is to utilize high resolution satellite imagery. The potential of satellite imagery to identify plots is increasingly recognized by scholars (Keita and Carfagna, 2009; Nelson and Swindale., 2014; Carletto *et al.*, 2015b). Until recently, however, satellite imagery has lacked the requisite resolution and accessibility to researchers. Indeed, Bilsborow and Henry (2012) note that “very high resolution” satellite imagery (<5 meters) is needed to gather useful spatial information for small agricultural plots. Other scholars assumed analysts or algorithms would identify plots (Fermont and Benson, 2011; Nelson and Swindale, 2013), but this underutilizes local knowledge and creates strong dependencies on very high resolution imagery that ultimately still suffers from plot characteristic-driven measurement error (e.g., spatial markers for identifying plots). A method utilizing satellite imagery and local knowledge to

accurately and inexpensively gather plot boundary data could simultaneously reduce effort and increase accuracy, making it more practical to consistently integrate spatial and household survey data.

In this paper, we report on the enumerator-assisted respondent-driven plot mapping (RDPM) method - a method using high resolution satellite imagery on tablet devices with respondent input to delineate small agricultural plots. Our method employs an enumerator-assisted respondent-driven mapping exercise to identify plots using off the shelf software and hardware and low cost satellite imagery. We do so by pre-loading higher resolution imagery (2.5 m resolution) on tablet devices with GPS chips, and having survey respondents identify their plots with respect to familiar landmarks (e.g., their own home, local schools, roads, rivers) to draw plot boundaries. To date, we know of two efforts that utilized similar methods (Cadena, 2013; Vale and Stabile, 2016), although one focused on sampling and the potential selection bias introduced by their methodology (Vale and Stabile, 2016), while the other utilized rich cadastral data (Cadena, 2013). Our study setting is in a data poor environment, where land records are not digitized, land transactions are often informal, and information is localized. We describe the process and utility of our method, and also report on challenges and potential limitations of its use. We tested our method with 540 households from 30 villages in central Kenya covering 744 plots.

## **Materials and Methods**

We use data from household survey, satellite imagery, and data collected via RDPM. Data come from a broader study on livelihoods, farming practices, soil erosion, and water quality in the Upper Tana watershed. We first describe the study site, the household survey, and the methodology. We then briefly discuss how these data were merged and analyzed.

### *Study site*

Our study site covers an area approximately 120 square kilometers consisting of 68 villages in the Upper Tana basin north of Nairobi in central Kenya. Data were collected from September-October 2015. We employed a two-stage simple random selection process. In the first stage 30 villages were randomly selected. In the second stage, approximately 18 households per village were interviewed from each village. Our final household survey sample consists of 540 households.

### *Household survey*

The gender disaggregated household survey collected information on household demographics, agricultural practices, household and productive assets, water treatment, conservation practices, access to markets and credit, community group membership, intrahousehold decision-making, and other topics. Importantly, the survey collected detailed information on agricultural plots, such as size, input use,

agricultural practices, crop types, and land cover. The survey was administered with tablet devices using the CSPro 6.1 computer-assisted personal interviewing (CAPI) software (CSPro 2015). A pair of male and female enumerators interviewed each household. For married households, male enumerators interviewed the husbands and the female enumerators interviewed the wives.

#### *Enumerator-assisted respondent-driven plot mapping*

RDPM provides a method for collecting plot boundaries from respondents at the residence by drawing the plot boundaries on GPS-enabled tablet devices pre-loaded with higher resolution imagery. By marrying high resolution imagery with respondent's knowledge about their own plots and landscape, the RDPM allows us to collect spatial data for each plot, even for irregularly shaped plots and plots that are far away from the household that may be systematically excluded using other methods (Kilic *et al.*, 2017). Because polygons are drawn directly on spatially explicit high resolution images, measurement error is solely driven by the enumerator or respondent, and not by technical errors caused by plot characteristics (e.g., tree cover and plot slope) and satellite position, signal propagation, and receivers (Hoffman-Wellenhof *et al.*, 2008; Keita and Carfagna, 2009, 2010; Fermont and Benson, 2011). The RDPM arguably saves significant time and effort required by both the tape and compass method and GPS measurement by allowing data collection to occur at the household's residence.

There are two primary components for the RDPM: high resolution imagery (e.g., 2.5m resolution) loaded within mobile GIS software and a tablet device with a GPS chip. We used ArcGIS Collector software to collect the spatial extent of each plot during the household surveys. Spatial data are linked to the household survey data via unique household and plot identifiers, both of which were assigned during the survey. A tablet device (either Android, iOS, or Windows) and a subscription to ArcGIS Online (AGO) (an online mapping service from ESRI, <http://arcgis.com>) are required to implement ArcGIS collector. Using Esri's World Imagery (2.5 meter resolution) map service as the base map made setup and implementation simple to cache the imagery to avoid the need for an internet connection in the field.

The first step in gathering plot boundaries was setting up the AGO server component. A feature service (a type of map layer) was published showing the study area boundary, and a second blank feature service was created from a shape file to capture desired fields, such as household and plot IDs. Finally, AGO accounts were created for each enumerator, and a group was created in AGO to control access.

We then configured the android tablets for data collection. We used 16 GB micro-SD cards in each tablet to store the cached satellite imagery since we did not have reliable data connectivity in the field. We then installed ArcGIS Collector and moved the application to the SD card. Each enumerator was assigned a unique account to ensure all edits could be tied to a particular enumerator and device. The web map was then downloaded to the tablet, choosing the study area as the "work area" and the

maximum available scale level (1:4,514). Edits made offline synced with AGO once an internet connection was available.

Enumerators were trained to assist respondents in identifying the location of the agricultural plot. Enumerators were instructed to first orient the respondent by centering the map to the respondent's home. Other landmarks in the community, such as schools, also helped orient the respondent. Respondents were then asked to identify plots, one at a time, to complement survey data for each plot (e.g., tenure status). Enumerators were instructed to not let the respondent actually draw the plot. Field testing revealed respondents were often uncomfortable articulating plot boundaries on tablets on their own. As a result, respondents were instead asked to clearly articulate the boundaries of the plot while the enumerator drew the polygon while being observed by the respondent. Once enumerators drew the plot boundaries, they were asked to input the household and plot IDs and repeat the process for each plot. If enumerators or respondents made errors, enumerators were instructed to delete and redraw the polygon. At the end of each day, enumerators were instructed to wirelessly upload polygon data to the AGO server. This allowed timely review of spatial data by the project supervisor so any potential issues were quickly identified and fixed.

#### *Merging household survey data and polygon data*

In total, the household survey collected data on 1,160 plots from 540 households. 80 percent of these plots were reported as being drawn in the survey (n=932), but in the spatial data we identified 778 complete agricultural plots. The discrepancy is due to data entry error by respondents in the survey (accidentally inputting a new entry for a plot when there were no more plots). In the spatial data, we dropped polygons that were clearly erroneous due to enumerator error (e.g., one large polygon covering multiple smaller polygons, duplicate polygons) (n=34), leaving 744 polygons. All 744 plots had matching household survey data and spatial data, which leaves 376 households total (70 percent of the original sample).

#### **Analysis and empirical strategy**

We report three sets of analyses to investigate the feasibility and utility of the RDPM for gathering small plot agricultural statistics. All covariates are consistent across each set of analysis and are estimated using robust standard errors clustered at the household-level. First, we investigate the determinants of successfully drawing plots using the RDPM, which reveal possible limitations and selection bias of the RDPM. We estimate a logistic regression such that:

$$\log\left(\frac{p_{ijk}}{1 - p_{ijk}}\right) = \beta_0 + \mathbf{X}_{ijk}\beta_1 + \mathbf{Z}_{ijk}\beta_2 + \mathbf{W}_{jk}\beta_3 + \mathbf{V}_k\beta_4 \quad (1)$$

where  $p_i = P(Y = 1 | \mathbf{X}_{ijk}, \mathbf{Z}_{ijk}, \mathbf{W}_{jk}, \mathbf{V}_k)$  if the enumerator successfully drew a plot with the respondent's guidance and zero otherwise.  $i, j,$  and  $k$  index respondent, household, and enumerator, respectively.  $\mathbf{X}_{ijk}$  is a vector of individual (respondent) covariates,  $\mathbf{Z}_{ijk}$  is a vector of household covariates,  $\mathbf{W}_{jk}$  is a vector of plot covariates, and  $\mathbf{V}_k$  is a vector of enumerator indicators. Individual-level covariates include information on sex, age, marital status, education, indicator for being the household head and being a farmer as their primary occupation. Household-level covariates include data on household wealth, household gender composition (percent female), household size, and the number of plots the household manages. Finally, plot covariates capture information on whether the household owns the plot, indicators for the plot's primary use, self-reported plot size, and distance from the household to the plot.

Our second set of analyses examines the quality of polygons. We rate polygon quality based on two criteria: 1) whether a polygon has clear overlap with structures that indicate measurement error (e.g., a polygon cutting through a building or road), and 2) a binary subjective rating of the polygon's fit to remotely sensed imagery (e.g., whether or not the boundary appears to match visible land cover features). We create three variables using the two criteria: one variable for each criterion that is coded one for poor quality and zero otherwise, and a variable coded as one if a plot had either clear or subjective overlap and zero otherwise. We estimate equation (1) above using the polygon quality variables and the same respondent, household, and plot covariates.

In our final analysis we investigate whether there are differences between survey responses and data extracted from spatial data. We primarily examine differences in measured plot size from polygons and self-reported survey data.<sup>2</sup> For land area, we operationalize two variables following Carletto *et al.* (2013): the absolute and relative difference between self-reported land areas and plot boundary land areas. The relative difference is calculated as  $(\text{Self reported plot size} - \text{Estimated plot size}) / \text{Estimated plot size}$ . We estimate the following linear regression model:

$$y_{ijk} = \beta_0 + \mathbf{X}_{ijk}\beta_1 + \mathbf{Z}_{ijk}\beta_2 + \mathbf{W}_{jk}\beta_3 + \mathbf{V}_k\beta_4 + \varepsilon_{ijk} \quad (2)$$

---

<sup>2</sup> We also assessed correlations on reported and remotely sensed land cover to investigate the congruence between spatial and survey data, which is reported on Appendix C.

where the dependent variable  $y_{ijk}$  is the absolute and relative difference in plot size. For all models, we estimated pooled models because we are interested how individual characteristics are associated with each dependent variable, but we also estimated household fixed effects model for robustness checks and present results in Appendix B. Results are largely similar.

## Results

### *a. Descriptive statistics*

Survey respondents consist mainly of male household heads (first column in Table 1). Male respondents have significantly higher educational attainment and are significantly more likely to be able to read, be married, be village leaders, and have occupations outside of farming. 94 percent of households own at least one plot, and households manage, on average, 2.2 plots. The vast majority of plots are used for agriculture. Potatoes (67%), cabbage (55%), and maize (40%) are the most common crops grown by households. Approximately 90 percent of plots are within 10 minutes walking distance from homes. Nearly 90 percent of plots are reportedly less than one acre.

### *b. Who was able to draw plots successfully?*

Approximately 70 percent of households were able to successfully draw at least one agricultural plot (second column in Table 1), and respondents and the households they belong to are largely similar to the overall sample population (third column in Table 1 for comparisons). Logistic regression results reveal that most observable characteristics are not significantly associated with being drawing plot boundaries successfully (Figure 1; Table B1 in Appendix B for full model results and robustness checks). Respondents were more likely to assist enumerators in drawing plot boundaries if they were farmers (12 percent), from larger households (2 percent), and owned the plot they were drawing (20 percent). Interestingly, many plot characteristics were significantly associated with successfully drawing plots. Plots used as a livestock field or cropland increase the probability of being drawn by approximately 12 percent and 8 percent, respectively. As self-reported plot size and the plot's distance from the household increase, respondents were significantly less likely to successfully draw plots by 2 percent and 1 percent, respectively.

Drawing plots requires familiarity with its boundary and location. For instance, farmers frequently visit plots used for livestock fields and crops, so it is unsurprising that these plots are significantly and positively associated with complete plot boundary drawings. Further, rented plots may be leased on a seasonal basis, and in general renters may be less familiar with a plot's boundary and location. Larger plots are less likely to be successfully drawn (Figure 2). For plot owners, plots are more likely to be drawn for plots smaller than 11 acres, and the probability was higher for farmers that owned

their plots. The probability significantly decreases once farmers do not own their plot, as plots smaller than 6 acres and 4 acres are likely to not be drawn for farmers and non-farmers, respectively. The negative relationship between the distance from the household to the plot is unsurprising. The RDPM centers maps on the respondent's home, and plots far from the home may be hard to identify without familiar or obvious landmarks.

### *c. Plot boundary quality*

744 plots were successfully drawn, but the quality of plot boundaries vary. For instance, 15 percent of the 744 polygons overlap with other polygons, although 10 percent of these polygons overlapped less than 5 percent with another polygon. Less than one percent (n=5) of successfully drawn plot boundaries had both subjective and objective overlap, and 2.7 percent (n=20) had clear overlap (without subjective overlap) and 10 percent (n=79) of plot boundaries had just subjective overlap. Less than one percent (n=5) had both subjective and objective overlap. We find that no individual or household characteristics are determinants of plot boundaries quality. We do find, however, that distance to the plot increases the likelihood of having poor quality plot boundaries ( $\beta_{distance\ to\ plot} = 0.048$ ,  $p < 0.10$ . Table 2), suggesting familiarity with the area around the plot is a determinant of plot boundary quality.

Panel A in Figure 3 presents a sample of small agricultural plots in the study area. Examining the high quality plot boundaries in Figure 3 demonstrates the benefits of using high resolution imagery to identify plots. For instance, Panel B shows the RDPM was able to delineate small adjacent agricultural plots. Each plot is approximately 0.25 acres (0.1 hectares), well below the threshold of small agricultural plots (Keita *et al.* 2010; Fermont and Benson 2011). The accuracy of the RDPM will likely improve as higher resolution data become available. Panel C demonstrates the RDPM's ability to draw irregularly shaped polygons, which correctly excludes non-cropland, such as houses and roads. The plot in Panel C is approximately 2.1 acres (0.87 hectares).

### *d. Comparison of spatial and survey data*

We now turn to examine the subset of plots that did not have clear or subjective overlaps (n=650). Self-reported and estimated plot sizes are significantly positively correlated at 0.69 ( $p < 0.01$ ); however, the correlation decreases to 0.27 for plots smaller than 1.2 acres. The weaker correlation for smaller plots is consistent with what others have found. For instance, Shøning *et al.* (2005) found correlation to be 0.65 overall between GPS measurements and tape and compass measurements, but the correlation was 0.89 for plots between larger than 1.2 acres (0.5-0.7 hectares), decreasing to 0.12 for plots under approximately 1.2 acres (0.5 hectares). Our data also reveal dispersion as plot size increases (Figure 4), as there is greater variation between self-reported and estimated plot size as the plot size increases

( $\beta_{plot\ size} = 0.5$ ,  $p < 0.001$ . Table 2). A possible explanation for this increased difference may stem from respondents being unable to correctly identify large plot boundaries. We also find respondents underestimate larger plots relative to smaller plots ( $\beta_{plot\ size} = -0.28$ ,  $p < 0.10$ . Table 2), which mirror findings by De Groote and Traorè (2005). Further, consistent with previous literature (Carletto *et al.*, 2015c), we find “heaping” in the distribution of self-reported plots compared to estimated plot sizes (blue line in Figure 4) for plots between 0.5-3 acres. The Kolmogorov-Smirnov test for equality of distributions also reveal significant differences in self-reported and estimated plot size distributions ( $p < 0.01$ ). Finally, we find the quality of plot boundaries collected via the RDPM does not systematically differ by observable individual and household characteristics. This suggests the accuracy of the RDPM is robust to diverse populations, as it is not contingent on factors such as gender, wealth, or educational attainment.

## Discussion and conclusion

Nearly two decades ago the National Research Council (1998) released a report arguing the utility, potential, and importance of integrating spatial data with social science research on socio-ecological systems and sustainability, and recent work has called attention has echoed this report (Donaldson and Storeygard, 2016). Despite advances in technology, accessibility, and affordability, agricultural household surveys still rarely gather spatially explicit plot boundaries. Our method for collecting spatial boundaries for small agricultural plots addresses shortcomings of other methods by utilizing high resolution satellite imagery (2.5 meters) and respondent knowledge to identify their own plots on GPS-enabled tablet devices. We demonstrate that there may be distinct advantages to the RDPM. For instance, the RDPM appears to perform better for smaller agricultural plots than larger plots. Because data were collected inside the household, the RDPM requires less effort from enumerators and respondents compared to other methods. Further, the RDPM utilizes spatial data to avoid errors caused by, for instance, unit conversion (Fermont and Benson, 2011; Carletto *et al.*, 2015b) by directly estimating the plot boundary. As a result, we believe that the RDPM provides a scalable method for gathering plot boundaries. Although we compared spatial data with household survey data and found reason to believe RDPM data provide utility to practitioners and researchers, further work should assess the accuracy of RDPM data via field results and address likely flaws in household survey. If accuracy is acceptable, the RDPM provides a new tool for combining household survey data with rich spatial datasets.

Our model estimates on respondents able to draw their plot boundaries suggest the RDPM should be conducted with household members most familiar with the household’s plots, such as plot owners and farmers. Not doing so may create systematic plot selection bias. While plot selection bias in GPS measurements is driven by exclusion criteria (Kilic *et al.*, 2017), such as the distance and accessibility of

the plot from the home, the RDPM relies more heavily on the respondent's knowledge of their local geography. Further, plots covering more land area may require particular attention. Although this is a minor concern for our study population because the average self-reported plot size is smaller than one acre, this may especially be an issue for plots that are infrequently used because respondents may not be familiar with plot boundaries, and unfamiliarity may increase with plot size. Having multiple household members help identify and draw plots when employing the RDPM may increase the likelihood of successful plot drawings. This would utilize multiple household members' knowledge about their landscape by, for instance, having the household member responsible for animal husbandry identify plot boundaries for livestock fields, or the household member responsible for firewood collection identify plot boundaries for forested plots.

We tested the RDPM in an area that is approximately 120 km<sup>2</sup>, but we believe the RDPM provides a scalable alternative to existing methods, thus enabling more widespread adoption of the practice of collecting spatial boundaries for surveyed plots. For instance, survey teams may be deployed to specific areas of the country, and high resolution imagery can be uploaded onto tablets for that region to efficiently use memory. Alternatively, in countries with reliable and broad mobile phone data coverage, imagery can be retrieved during the interview. The increasing quality of freely available satellite imagery also creates new possibilities and greater accuracy for identifying small agricultural plots using RDPM. Cloud cover may still be an impediment for areas with persistent cloud cover throughout the year (Carletto *et al.*, 2015c). While Billsborrow and Henry (2012) did not have a single cloud free image covering their entire study area in the 8 year period they examined, this is increasingly uncommon and only limits a remote sensing approach. By contrast, the RDPM relies on a patchwork of the latest available imagery for each tile, meaning that while a given study area is likely made up of imagery from different dates, it is very unlikely there will be areas with no imagery capture in the last year or two.

Like all methods, however, there are tradeoffs and limitations. While the sample that drew plots are largely similar to the overall respondent sample, regression results suggest the RDPM may systematically exclude agricultural plots from some subpopulations. We find respondents that are farmers and landowners are more likely to be able to identify and draw plot boundaries, and larger plots and plots further from the household are less likely to be drawn. We believe having multiple household members identify plot boundaries, however, can minimize or eliminate this bias. The RDPM is also subject to enumerator error. Our results consistently found some enumerators were associated with lower likelihood of drawing plots and poorer plot quality, highlighting the importance of carefully training enumerators. In addition, the RDPM may be better suited for some plot or crop types. For instance, identifying specific boundaries within an intact forest may be challenging for respondents. Identifying plot boundaries for areas that have clear spatial markers, such as rice fields with drainage ditches or non-cropped areas

between plots, may be simpler. Despite these limitations in certain contexts, we believe that the overall RDPM offers significantly lower cost than GPS measurements. Informal enumerator feedback suggests a plot can be drawn in under ten minutes. Further, although we did not formally test the RDPM's relative accuracy to other methods, we believe the RDPM potentially has equal or greater accuracy than GPS measurements for small plots because GPS measurements still suffer from errors caused by plot characteristics and technical issues (Hoffman-Wellenhof *et al.*, 2008; Keita and Carfagna, 2009, 2010; Fermont and Benson, 2011).

Finally, we believe the RDPM opens new opportunities for merging spatial and survey data analysis. For instance, preliminary analyses using this study's plot boundaries with very high resolution spatial data (0.5 meter resolution) identified past sustainable agricultural practices not identified in the survey. Ayana *et al.* (2017) used remote sensing to identify the presence of drainage ditches and ridge-tillage furrows in this study area and found 69 percent of successfully drawn survey plots contained ditches and furrows in the remote sensing analysis even though only 2.6 percent of plots in the survey reported having either drainage ditches or ridge-tillage. Further, the RDPM has applications for other spatial information. Future researchers may ask respondents to trace the route from their house to the market, providing a new way for gathering information on informal road networks. For rural households that rely heavily on footpaths, this type of spatial information can only be obtained by relying on local knowledge.

Future work should carefully match survey questions on plots by, for instance, taking into account seasonal land use and land cover patterns (e.g., cropland may be temporarily bare after harvest, but remain cropland as a land use). This also provides an opportunity to improve analysis on crop yields by linking plot boundaries with spatial time series data to determine the proportion of plots used for agriculture. Further, the RDPM opens a way to test respondent recall on land cover change or land degradation (e.g., overgrazing of fields) for small agricultural plots by comparing survey questions about past events to spatial time series data. As new technologies (e.g., micro- and nano-satellites) become available, we believe the RDPM will provide an integrated and cost effective method for simultaneously gathering spatial and survey data. Smallholder farms remain a critical factor for agricultural production across the world (Lowder *et al.*, 2016), and is closely tied to priority development goals, such as food security and poverty reduction. Poor quality agricultural data continues to be a challenging factor for allocating resources and monitoring progress on development (African Development Bank, 2004; FAO, 2008; Carletto *et al.* 2015b). While Carletto *et al.* (2015b) outline a number of areas for addressing this challenge, they note that methodological improvements are a critical pathway forward. We believe the RDPM provides a simple, time saving method that can reduce the burden of data collection on enumerators and project teams to collect data on plot boundaries, which are critical for understanding

agricultural practices, yield, land degradation, and other information necessary for informing development policies.

DRAFT

## References

- African Development Bank 2004. The Marrakech action plan for statistics: Better data for better results: An action plan for improving development statistics. 2<sup>nd</sup> International Roundtable on managing for Development Results, Marrakech.
- Ali, D.A., & Deininger, K. 2015. Is there a farm-size productivity relationship in African agriculture? Evidence from Rwanda. *Land Economics*, **91**(2): 317-343.
- André, C. & Platteau, J.P. 1998. Land relations under unbearable stress: Rwanda caught in the Malthusian trap. *Journal of Economic Behavior & Organization*, **34**(1): 1-47.
- Ayana, E.K., Fisher, J.R.B., Hamel, P., & Boucher, T.M. 2017. Identification of ridge-tillage and ditches using remote sensing for improved hydrological modeling. *International Journal of Remote Sensing*, **38**(16), 4611-4630.
- Barrett, C., Bellemare, M.F., & Hou, J.Y. 2010. Reconsidering conventional explanations of the inverse productivity–size relationship. *World Development*, **38**(1): 88–97.
- Baten, J., & Juif, D. 2013. A story of large landowners and math skills: Inequality and human capital formation in the long-run development. 1800-2000. *Journal of Comparative Economics*.
- Benjamin, D. 1995. Can unobserved land quality explain the inverse productivity relationship? *Journal of Development Economics*, **46**: 51-84.
- Bhalla, S.S., & Roy, P., 1988. Misspecification in farm productivity analysis: the role of land quality. *Oxford Economic Papers*, **40**: 55–73.
- Bilsborow, R.E., & Henry, S.J. 2012. The use of survey data to study migration-environment relationships in developing countries: Alternative approaches to data collection. *Population and Environment*, **31**(1): 13-141.
- Brown, M.E., Grace, K., Shively, G., Johnson, K.B., & Carroll, M. 2014. Using satellite remote sensing and household survey data to assess human health and nutrition response to environmental change. *Population and Environment*, **36**: 48–72.
- Carletto, C., Savastano, S., Zezza, A. 2013. Fact or artifact: The impact of measurement errors on the farm size–productivity relationship. *Journal of Development Economics*, **103**: 254–261.
- Carletto, C., Gourlay, S., & Winters, P. 2015a. From guesstimates to GPStimates: Land area measurement and implications for agricultural analysis. *Journal of African Economies*, 1-35.
- Carletto, C., Jolliffe, D., & Banerjee, R. 2015b. From tragedy to renaissance: Improving agricultural data for better policies. *Journal of Development Studies*, **51**(2): 133–148.
- Caletto, C., Gourlay, S., Murray, S., & Zezza, A. 2015c. Welcome to fantasyland: Comparing approaches to land area measurement in household surveys. *International Conference of Agricultural Statistics*. Milan.

- Census and Survey Processing System. 2015. U.S. Census Bureau, Washington, DC.  
<http://www.census.gov/population/international/software/cspro/>
- David, P. 1978. Non-sampling errors in agricultural surveys. Review, current findings, and suggestions for future research. *Philippine Statistical Association Annual Conference*. Manila.
- De Groote, H., & Traore, O. 2005. The cost of accuracy in crop area estimation. *Agricultural Systems*, **84**: 21-38.
- Deininger, K., & Squire, L. 1998. New ways of looking at old issues: inequality and growth. *Journal of Development Economics*, **57**(2): 259-287.
- Donaldson, D., & Storeygard, A. 2016. The view from above: Applications of satellite data in economics. *The Journal of Economic Perspectives*, **30**(4): 171–198.
- Food and Agricultural Organization. 2008. *The agricultural bulletin board on data collection, dissemination and quality of statistics*. Rome: FAO.
- Food and Agricultural Organization. 2014. *The State of Food and Agriculture 2014. Innovations in Family Farming*. Rome.
- Fermont, A., & Benson, T. 2011. Estimating yield of food crops grown by small-holder farmers. A review in the Uganda Context. International Food Policy Research Institute, Washington, D.C.
- Goldstein, M., & Udry, C. 1999. Agricultural innovation and risk management in Ghana. International Food Policy Research Institute, Washington, D.C.
- Hansen, M.C., Potapov, P.V., Moore, R., Hancher, M., Turubanova, S.A., Tyukavina, A., ... Townshend, J.R.G. 2013. High resolution global maps of 21<sup>st</sup>-century forest cover change. *Science*, 342(850), 850-853.
- Hofmann-Wellenhof, B., Lichtenegger, H., & Wasle, E. 2008. *GNSS – Global Navigation Satellite Systems*. Springer-Verlag, New York, NY.
- Holden, S.T., & Fisher, M. 2013. Can area measurement error explain the inverse farm size productivity relationship? Centre for Land Tenure Studies Working Paper 12/13, Norwegian University of Life Sciences.
- International Fund for Agricultural Development. 2010. *Rural poverty report 2011: New realities, New challenges, New opportunities for tomorrow's generation*. Rome, Italy.
- Johnson, K.B., Jacob, A., & Brown, M.E. 2013. Forest cover associated with improved child health and nutrition: Evidence from the Malawi Demographic and Health Survey and satellite data. *Global Health: Science and Practice*, **1**: 237–248.
- Keita, N., Carfagna, E., & Mu'Ammar, G. 2010. Issues and guidelines for emerging use of GPS and PDAs in agricultural statistics in developing countries. *The Fifth International Conference on Agricultural Statistics (ICAS V)*, Kampala.
- Kilic, T., Zezza, A., Carletto, C., & Savastano, S. 2017. Missing(ness) in action: selectivity bias in GPS-based land area measurements. *World Development*, **92**: 143-157.

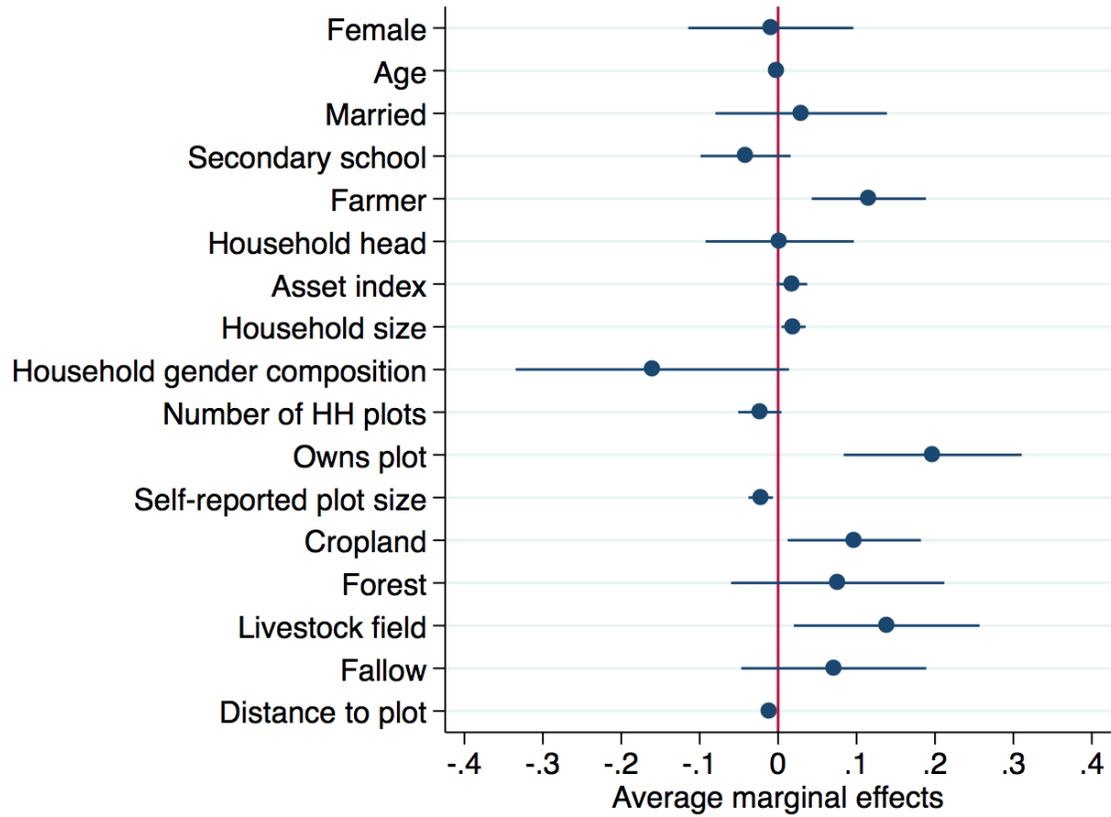
- Lamb, R.L. 2003. Inverse productivity: land quality, labor markets, and measurement error. *Journal of Development Economics*, **71**(1): 71–95.
- Lowder, S.K., Skoet, J. & Raney, T. 2016. The number, size, and distribution of farms, smallholder farms, and family farms worldwide. *World Development*, **87**: 16-29.
- Macours, K. 2011. Increasing inequality and civil conflict in Nepal. *Oxford Economic Papers*, **63**(1), 1-26.
- National Research Council. 1998. *People and pixels: Linking remote sensing and social science*. Liverman, D., Moran, E.F., Rindfuss, R.R., and Stern, P.C. (ed.). National Academy Press, Washington, DC.
- Nelson, S., & Swindale, A. 2013. *Feed the Future agricultural indicators guide: Guidance on the collection and use of data for selected Feed the Future agricultural indicators*. Westat, Rockville, MD.
- Savastano, S., Carletto, G., & Zezza, A. 2010. Using global position system for land measurement: Testing the farm size-productivity relationship. *International Conference on Agricultural Statistics*.
- Schøning, P., Apuuli, J.B.M., Menyha, E., & Muwanga-Zake, E.S.K. 2005. Handheld GPS Equipment for Agricultural Statistic Surveys. Experiments on Area Measurements Done During Fieldwork for the Uganda Pilot Census of Agriculture, 2003. Statistics Norway, Oslo.
- Thomlinson, J. R., P. V. Bolstad, & Cohen, W.B. 1999. Coordinating methodologies for scaling landcover classifications from site-specific to global: Steps toward validating global map products. *Remote Sensing of Environment*, **70**(1): 16–28.
- Vale, P.M., & Stabile, M.C.C. 2016. GIS without GPS: new opportunities in technology and survey research to link people and place. *Population and Environment*, **37**: 391-410.

Table 1: Descriptive statistics

	Full sample		Drew plot		Diff
	(1)	(2)	(3)	(1)-(2)	(3)
	Mean	SD	Mean	SD	(1)-(2)
<b>Individual characteristics</b>					
Female (%)	32	47	31	46	0.12
Household head (%)	81	39	80	40	1.9
Age	50	14	49	14	0.44
Completed secondary school (%)	30	46	30	46	0.30
Married (%)	78	42	80	40	-1.9
Farmer (%)	80	40	84	36	-4.6*
Years farming	24	15	25	15	-1.6
<b>Household characteristics</b>					
Household size	5.0	1.8	5.2	1.8	-0.14
Gender composition (% of HH that is female)	39	16	38	16	0.59
Wealth index	0.0	0.47	0.12	1.9	-0.39
<b>Plot characteristics (reported)</b>					
Number of plots per household	2.2	1.1	2.2	1.1	0.011
Own plot (%)	94	22	98	0.46	-3.5***
Self-reported plot size (Ha)	0.68	1.1	0.64	0.87	0.039
Plot used for cropland (%)	87	33	87	34	0.058
Plot kept fallow (%)	2.6	16	2.4	15	0.17
Plot used for grazing livestock (%)	5.2	21	5.6	23	-0.47
Plot used for woodland (%)	2.0	14	1.7	13	0.32
Plot used for other (%)	3.0	17	3.0	17	-0.074
Distance to plot (minutes)	5.7	15	2.9	6.0	2.7***
<i>n</i> (Households)	540		376		
<i>n</i> (Plots)	1,160		744		

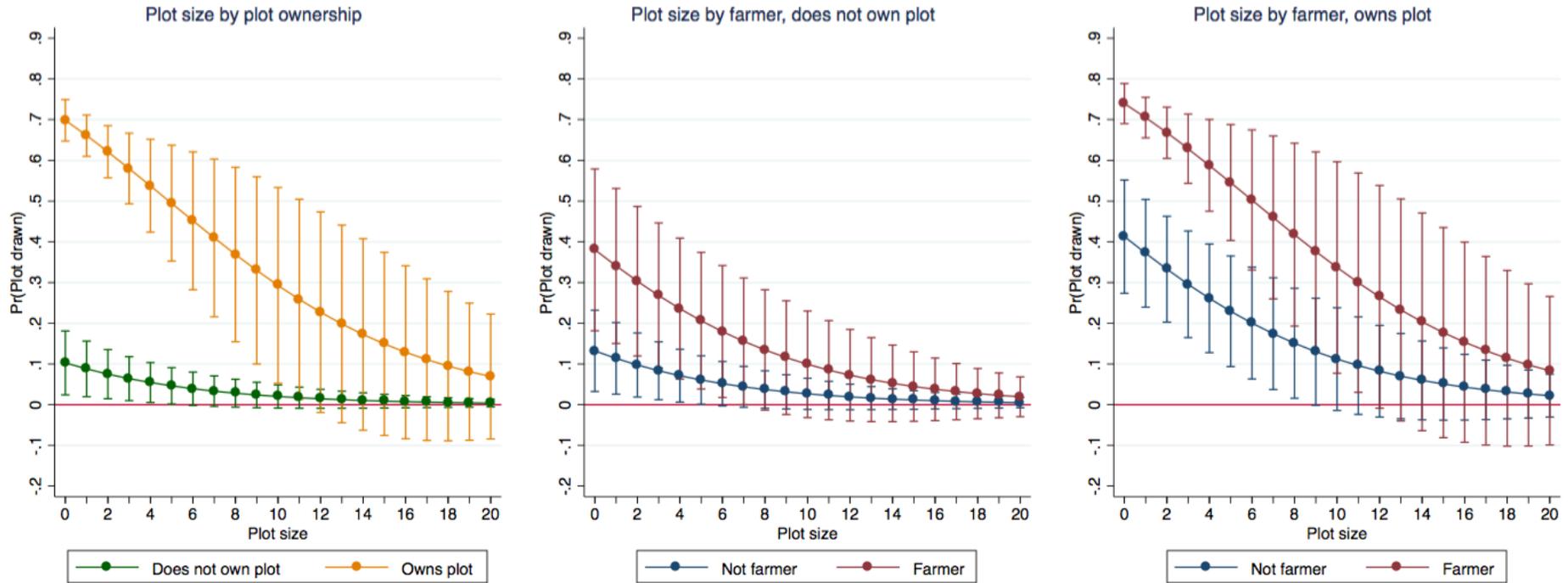
\*\*\* p&lt;0.01, \*\* p&lt;0.05, \* p&lt;0.10

Figure 1: Average marginal effects for logistic regression for drawing plot <sup>a</sup>



<sup>a</sup> Coefficients are displayed in Table B2. Bars shows 90% confidence intervals. All models estimated using robust standard errors clustered at the household-level and enumerator fixed effects.

Figure 2: Predicted probability of drawing the map by plot size, farmer occupation, and plot ownership<sup>a</sup>



<sup>a</sup> Predicted probabilities are estimated with variables are held at their means. Models are estimated using robust standard errors clustered at the household with 90% confidence intervals.

Figure 3: Spatial plots

Panel A: Sample of plot drawings

Panel B: Small agricultural plots



DRAFT



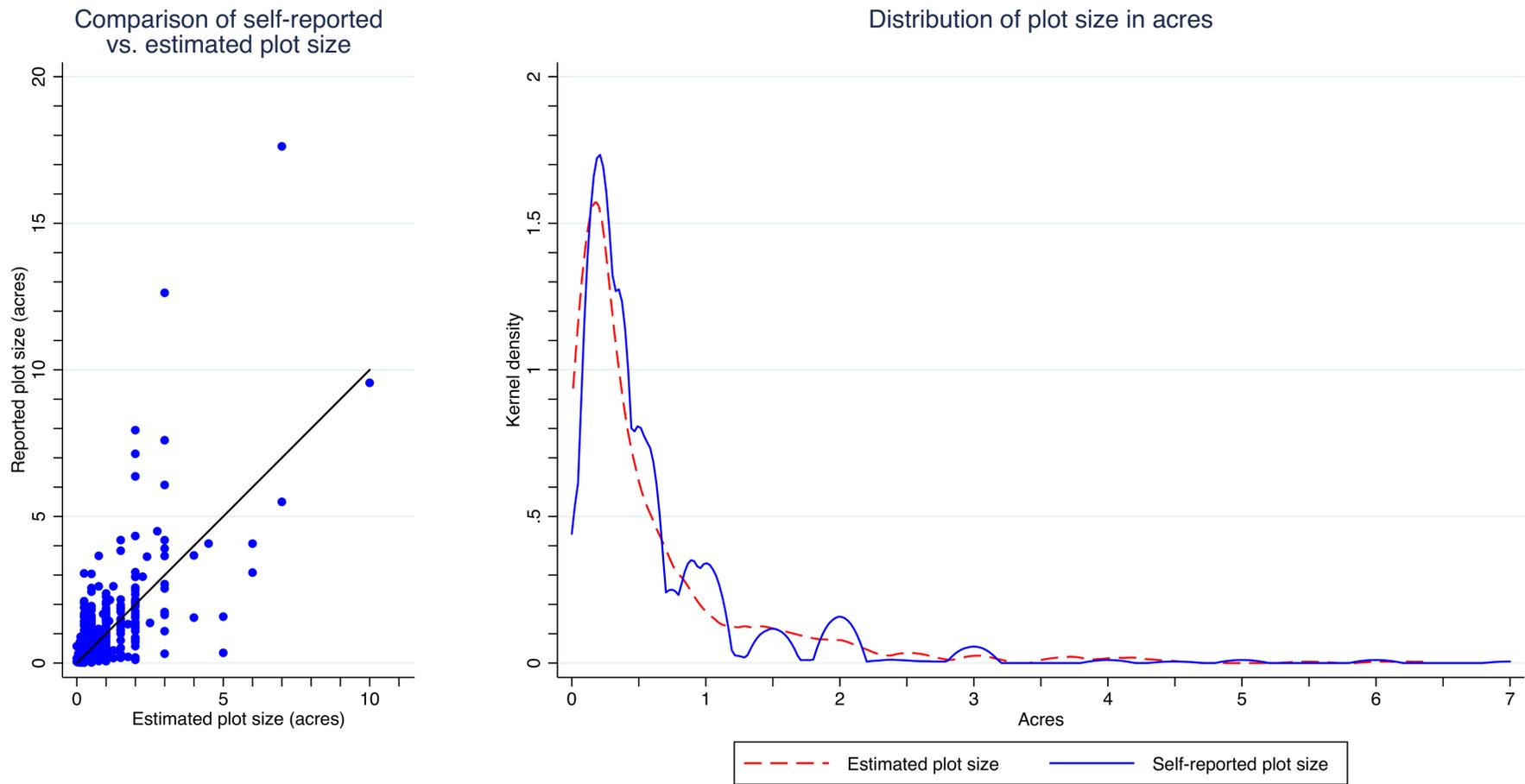
Table 2: Regression results for plot boundary quality and differences between survey and spatial data<sup>a</sup>

	<b>Clear overlap</b>	<b>Subjective overlap</b>	<b>Clear and Subjective overlap</b>	<b>Absolute difference in plot size</b>	<b>Relative difference in plot size</b>
	<b>(1)</b>	<b>(2)</b>	<b>(3)</b>	<b>(4)</b>	<b>(5)</b>
Female	-0.42 (1.081)	0.85 (0.592)	0.46 (0.597)	0.03 (0.135)	0.43 (0.588)
Age	0.01 (0.018)	0.02 (0.011)	0.01 (0.011)	0.00 (0.002)	-0.00 (0.011)
Married	-0.83 (0.842)	1.28* (0.607)	0.83 (0.561)	-0.05 (0.180)	0.47 (0.578)
Secondary school	-0.49 (0.741)	-0.46 (0.326)	-0.46 (0.312)	-0.00 (0.066)	0.34 (0.263)
Is a farmer	1.70 (1.440)	-0.65 (0.438)	-0.35 (0.415)	0.04 (0.107)	-0.03 (0.264)
Is HH head	-0.63 (1.180)	-0.24 (0.587)	-0.39 (0.572)	0.01 (0.153)	-0.67 (0.444)
Asset index	0.09 (0.154)	0.13 (0.087)	0.12 (0.080)	-0.01 (0.022)	-0.08 (0.086)
Household size	0.04 (0.184)	-0.10 (0.119)	-0.08 (0.105)	0.01 (0.014)	0.01 (0.072)
HH gender composition	1.94 (1.912)	-0.95 (0.963)	-0.22 (0.922)	-0.16 (0.284)	0.49 (0.927)
Number of HH plots	-0.37 (0.312)	-0.07 (0.145)	-0.11 (0.137)	-0.01 (0.028)	0.09 (0.122)
Owns plot	- (-)	- (-)	- (-)	-0.70 (0.639)	-0.47 (0.368)
Self-reported plot size	-0.12 (0.315)	-0.18 (0.200)	-0.14 (0.165)	0.51*** (0.142)	-0.28* (0.149)
Cropland	-0.24 (0.708)	0.09 (1.019)	-0.08 (0.800)	0.13 (0.231)	-0.36 (0.434)
Forest	- (-)	- (-)	- (-)	0.61 (0.620)	-0.08 (0.784)
Livestock field	-1.20*** (0.409)	-1.01 (1.134)	-1.14 (0.854)	0.31 (0.365)	0.04 (0.579)
Fallow	- (-)	-0.12 (1.245)	-0.56 (1.106)	0.24 (0.251)	0.58 (0.504)
Distance to plot	0.02 (0.019)	0.04*** (0.014)	0.05* (0.025)	-0.00 (0.009)	0.02 (0.050)
Constant	-3.96 (3.497)	-2.32 (2.049)	-1.91 (1.812)	0.51 (0.686)	-0.67 (1.145)
Enumerator FE	YES	YES	YES	YES	YES
Pseudo R-squared	0.13	0.10	0.09	-	-
R-squared	-	-	-	0.37	0.10
Model	Logistic	Logistic	Logistic	OLS	OLS
<i>n</i> (Plots)	649	701	701	650	650

<sup>a</sup> Robust standard errors clustered at the household-level in parentheses. Models 4-5 estimated on a subset of plots that do not have clear or subjective overlaps. Some variables did not have sufficient variation, such as plot ownership and were dropped from the model. Full model results and models with household fixed effects in Tables B2-B5 in Appendix B.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.10

Figure 4: Reported versus estimated plot size<sup>a</sup>



<sup>a</sup> Solid line plots a 45 degree line. Kernel density plots present plots less than 7 acres. Kolmogorov-Smirnov test for equality of distributions is significant at  $p < 0.01$ .

## Appendix A: Measurement methods

Table A1 outlines the main methods for measuring agricultural plot sizes. The table is adapted from information from Keita *et al.* (2010), Fermont and Benson (2011), and USAID’s Feed the Future Agricultural Indicators Guide (Nelson and Swindale, 2013).

Table A1: Methods for measuring plot size

Method	Description	Advantages	Disadvantages
Tape and compass	The “gold standard” for plot measurement. Method employs tape and compass to measure the sides of a plot and the angles of the corners to calculate total plot area. Can be used for irregularly shaped plots, although to do so measuring multiple polygons may necessary.	The gold standard for plot area measurement. Equipment is inexpensive (i.e., measuring tape and compass).	Can be extremely time consuming, especially when there are multiple plots or if plots are irregularly shaped. May require substantial training to minimize error.
GPS measurement	This method uses Global Positioning Systems devices to walk around plots for area measurement. Can be used for irregular plots. Plot measurements are taken when the GPS device connects with at least three satellites to measure latitude, longitude, and elevation. The enumerator will start at one corner of the field plot and walk fully around the perimeter. The average unit is accurate within 10–12 meters.	Time savings can be substantial, with some estimates being 300 percent compared to tape and compass methods (Fermont and Benson 2011).	Measurements are largely accurate, but for smaller plots this method may be limited. Weather, plot slope, and tree canopy cover may also disturb measurement.
Remote sensing	Using remotely sensed images to identify plots.	Can save time and money by avoiding direct measurement of plots.	If images are not high resolution it may be difficult to identify plots, especially small plots.
Direct farmer estimation	Relies on interviewing farmers and asks for estimates of plot size.	Can save time and money by avoiding direct measurement of plots.	Relies on farmers trusting enumerators, and farmers that lack formal education and quantitative skills may not be able to accurately report plot size (De Groot and Traoré 2005). Data quality also affected by size of the plot. Further, reported plot size may be “lumpy” – in other words, respondents may

			choose to answer in fixed unit increments for simplicity (David 1978).
Pacing	Method is often employed when respondents themselves may have low skills or knowledge to answer accurately about plot size. Measurements are taken by using the individual's pace, where one step is one unit. These paces are then converted to standard units.	Cheapest method that requires little to no skill.	Prone to error as an individual's pace can vary significantly by slope, season, or stability of the ground (levelness).

DRAFT

## Appendix B: Regression tables

Table B1: Logistic regression for successfully drawing polygon<sup>a</sup>

	(1)	(2)	(3)	(4)
Female	-0.28 (0.372)	-0.15 (0.415)	-0.07 (0.497)	-
Age	-0.01 (0.008)	-0.01 (0.009)	-0.02 (0.010)	-
Married	0.35 (0.382)	0.13 (0.436)	0.23 (0.515)	-
Secondary school	-0.23 (0.232)	-0.34 (0.245)	-0.32 (0.271)	-
Is a farmer	0.99*** (0.298)	1.02*** (0.300)	0.90** (0.353)	-
Is HH head	-0.17 (0.376)	-0.01 (0.381)	0.02 (0.446)	-
Asset index	-	0.11 (0.075)	0.14 (0.093)	-
Household size	-	0.16** (0.066)	0.15** (0.073)	-
HH gender composition	-	-0.70 (0.735)	-1.24 (0.829)	-
Number of HH plots	-	-0.06 (0.110)	-0.18 (0.129)	-
Owns plot	-	-	1.53*** (0.523)	1.21*** (0.369)
Self-reported plot size	-	-	-0.17** (0.075)	-0.14** (0.059)
Cropland	-	-	0.75* (0.406)	0.16 (0.368)
Forest	-	-	0.59 (0.648)	0.03 (0.559)
Livestock field	-	-	1.07* (0.565)	0.61 (0.480)
Fallow	-	-	0.55 (0.559)	-0.09 (0.530)
Distance to plot	-	-	-0.09*** (0.025)	-0.06*** (0.011)
Constant	1.10 (0.831)	0.76 (1.034)	0.12 (1.237)	-
Enumerator FE	YES	YES	YES	NO
Household FE	NO	NO	NO	YES
Observations	1,160	1,160	1,160	1,159
Pseudo R-squared	0.25	0.26	0.37	0.07

<sup>a</sup> Robust standard errors in parentheses

\*\*\* p<0.01, \*\* p<0.05, \* p<0.10

Table B2: Logistic regression for plot boundary quality<sup>a</sup>

	(1)	(2)	(3)
	Clear overlap	Subjective overlap	Clear and Subjective overlap
Female	-0.42 (1.081)	0.85 (0.592)	0.46 (0.597)
Age	0.01 (0.018)	0.02 (0.011)	0.01 (0.011)
Married	-0.83 (0.842)	1.28* (0.607)	0.83 (0.561)
Secondary school	-0.49 (0.741)	-0.46 (0.326)	-0.46 (0.312)
Is a farmer	1.70 (1.440)	-0.65 (0.438)	-0.35 (0.415)
Is HH head	-0.63 (1.180)	-0.24 (0.587)	-0.39 (0.572)
Asset index	0.09 (0.154)	0.13 (0.087)	0.12 (0.080)
Household size	0.04 (0.184)	-0.10 (0.119)	-0.08 (0.105)
HH gender composition	1.94 (1.912)	-0.95 (0.963)	-0.22 (0.922)
Number of HH plots	-0.37 (0.312)	-0.07 (0.145)	-0.11 (0.137)
Owns plot	-	-	-
Self-reported plot size	-0.12 (0.315)	-0.18 (0.200)	-0.14 (0.165)
Cropland	-0.24 (0.708)	0.09 (1.019)	-0.08 (0.800)
Forest	-	-	-
Livestock field	-1.20*** (0.409)	-1.01 (1.134)	-1.14 (0.854)
Fallow	-	-0.12 (1.245)	-0.56 (1.106)
Distance to plot	0.02 (0.019)	0.04*** (0.014)	0.05* (0.025)
Constant	-3.96 (3.497)	-2.32 (2.049)	-1.91 (1.812)
Enumerator FE	YES	YES	YES
Observations	649	701	701
Pseudo R-squared	0.13	0.10	0.09

<sup>a</sup> Robust standard errors in parentheses. Some variables did not have sufficient variation, such as plot ownership and were dropped from the model.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.10

Table B3: Logistic regression for plot boundary quality with household fixed effects<sup>a</sup>

	(1)	(2)	(3)
	Clear overlap	Subjective overlap	Clear and Subjective overlap
Owens plot	0.03 (0.048)	0.11 (0.091)	0.14 (0.097)
Self-reported plot size	-0.00 (0.007)	-0.00 (0.014)	-0.00 (0.015)
Cropland	-0.02 (0.035)	0.02 (0.066)	-0.00 (0.071)
Forest	-0.05 (0.057)	-0.12 (0.108)	-0.17 (0.116)
Livestock field	-0.02 (0.043)	-0.03 (0.081)	-0.05 (0.087)
Fallow	-0.04 (0.051)	0.02 (0.097)	-0.03 (0.105)
Distance to plot	0.00 (0.001)	0.00** (0.002)	0.00** (0.002)
Constant	0.02 (0.059)	-0.03 (0.112)	-0.01 (0.121)
Household FE	YES	YES	YES
Observations	731	731	731
Pseudo R-squared	0.02	0.02	0.02

<sup>a</sup> Robust standard errors in parentheses. Some variables did not have sufficient variation, such as plot ownership and were dropped from the model.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.10

Table B4: Linear regression for the absolute difference between self-reported and estimated plot size

	(1)	(2)	(3)	(4)
Female	-0.04 (0.185)	-0.03 (0.170)	0.03 (0.135)	- -
Age	0.00 (0.003)	0.00 (0.003)	0.00 (0.002)	- -
Married	-0.10 (0.219)	-0.10 (0.225)	-0.05 (0.180)	- -
Secondary school	0.03 (0.075)	-0.02 (0.087)	-0.00 (0.066)	- -
Is a farmer	0.08 (0.144)	0.08 (0.161)	0.04 (0.107)	- -
Is HH head	0.05 (0.184)	0.10 (0.171)	0.01 (0.153)	- -
Asset index	-	0.05*** (0.020)	-0.01 (0.022)	- -
Household size	-	0.01 (0.022)	0.01 (0.014)	- -
HH gender composition	-	-0.34 (0.429)	-0.16 (0.284)	- -
Number of HH plots	-	-0.08*** (0.030)	-0.01 (0.028)	- -
Owns plot	-	-	-0.70 (0.639)	-0.63*** (0.189)
Self-reported plot size	-	-	0.51*** (0.142)	0.50*** (0.031)
Cropland	-	-	0.13 (0.231)	0.17 (0.148)
Forest	-	-	0.61 (0.620)	0.65*** (0.233)
Livestock field	-	-	0.31 (0.365)	0.24 (0.178)
Fallow	-	-	0.24 (0.251)	0.21 (0.216)
Distance to plot	-	-	-0.00 (0.009)	0.00 (0.006)
Constant	0.03 (0.309)	0.40 (0.365)	0.51 (0.686)	0.53** (0.242)
Household FE	NO	NO	NO	YES
Enumerator FE	YES	YES	YES	NO
Observations	650	650	650	650
R-squared	0.05	0.07	0.37	0.34

<sup>a</sup> Robust standard errors in parentheses. All models estimated on a subset of plots that do not have clear or subjective overlaps.

\*\*\* p<0.01, \*\* p<0.05, \* p<0.10

Table B5: Linear regression for relative difference in plot size<sup>a</sup>

	(1)	(2)	(3)	(4)
Female	0.50 (0.630)	0.48 (0.589)	0.43 (0.588)	- -
Age	-0.01 (0.011)	-0.00 (0.011)	-0.00 (0.011)	- -
Married	0.51 (0.579)	0.49 (0.574)	0.47 (0.578)	- -
Secondary school	0.28 (0.237)	0.37 (0.262)	0.34 (0.263)	- -
Is a farmer	-0.10 (0.273)	-0.09 (0.269)	-0.03 (0.264)	- -
Is HH head	-0.57 (0.429)	-0.65 (0.452)	-0.67 (0.444)	- -
Asset index	- -	-0.10 (0.080)	-0.08 (0.086)	- -
Household size	- -	-0.00 (0.073)	0.01 (0.072)	- -
HH gender composition	- -	0.55 (0.932)	0.49 (0.927)	- -
Number of HH plots	- -	0.15 (0.110)	0.09 (0.122)	- -
Owens plot	- -	- -	-0.47 (0.368)	-1.05 (0.810)
Self-reported plot size	- -	- -	-0.28* (0.149)	-0.15 (0.132)
Cropland	- -	- -	-0.36 (0.434)	-0.51 (0.632)
Forest	- -	- -	-0.08 (0.784)	-0.10 (0.994)
Livestock field	- -	- -	0.04 (0.579)	-0.01 (0.762)
Fallow	- -	- -	0.58 (0.504)	0.63 (0.924)
Distance to plot	- -	- -	0.02 (0.050)	0.01 (0.027)
Constant	-0.86 (0.762)	-1.64* (0.918)	-0.67 (1.145)	0.53 (1.033)
Household FE	NO	NO	NO	YES
Enumerator FE	YES	YES	YES	NO
Observations	650	650	650	650
R-squared	0.09	0.09	0.10	0.01

<sup>a</sup> Robust standard errors in parentheses. All models estimated on a subset of plots that do not have clear or subjective overlaps. Relative plot size is calculated as  $(Self - reported\ plot\ size - Estimated\ plot\ size) / Estimated\ plot\ size$ .  
\*\*\* p<0.01, \*\* p<0.05, \* p<0.10

## Appendix C: Land cover comparison

We also compared land cover classifications identified by spatial and survey data. We unfortunately were not able to ground truth land cover on plots, so we rely on remotely sensed land cover data instead. While land cover data are not the same as ground truthed data, the land cover classification product has high accuracy with ground truthed data.

We identified land cover using a supervised classification process in ArcGIS Feature Analyst on Pleiades satellite imagery with 0.5 m pan-sharpened spatial resolution. The imagery was collected on June 22, 2015, right after the harvest season and shortly before the survey implementation. We utilized 326 ground control points with measured land cover and agricultural practices (collected during two weeks of field work in June 2015). Ground control points were used in training the land cover classification. We focused on four land cover categories: trees, vegetated cropland, grass fields, and non-vegetated cropland (temporarily bare cropland/recently planted or harvested cropland). Producer accuracy (or the probability that a randomly selected pixels is classified correctly) of land cover classifications varied, but the overall accuracy was 81 percent although this varied by land cover type. The producer accuracy is 91 percent for vegetated cropland, 79 percent for non-vegetated cropland, 84 percent for grass fields, and 74 percent for trees. Thomlinson *et al.* (1999) establish an overall accuracy target of 85 percent, with each category having an accuracy of 75 percent or greater. Given the spatial heterogeneity and small plot size of our study area, we believe 81 percent accuracy to be reasonably high.

We focus on four land cover types: trees, vegetated cropland, grass fields, and non-vegetated cropland (temporarily bare cropland/recently planted or harvested cropland). Plots were identified as being in a category if the survey data indicated the plot was used for crops, had forest, kept fallow, or was a field for livestock, while a separate variable identified land cover type from spatial data. We find that plots reported as vegetated cropland in the survey have the largest proportion of overlap (95 percent, see Table 2). Trees, non-vegetated cropland, and grass fields had low congruence with remotely sensed land cover data with overlap at 29, 19, and 12 percent, respectively. But even with vegetated cropland we find the reliability of survey responses to be low. For instance, approximately 27 percent of plots identified as vegetated cropland in the survey were not identified as crops in the spatial data ( $n=162$ ), and 75 percent of plots that were reportedly not cropland in the survey were spatially identified as being cropland. Correlations between land cover classifications from spatial and survey data are also extremely low, as no classification has correlation coefficients above 0.15.

These results have several implications. The results suggest ground truthing are needed to verify land cover data from surveys, as our results indicate survey data may grossly misclassify land cover for small agricultural plots. The spatial land classification's accuracy assessment suggests the spatial land cover classification is highly reliability, so we believe our findings raise considerable questions about land cover questions in survey data for small agricultural plots. The probability of misclassifying land cover in the spatial data and survey data is extremely low. For instance, a 1.2 acre plot (0.5 hectares) at 2.5 meter resolution consists of approximately 20,000 pixels; since only a single pixel is required to constitute a match with the survey data (indicating a given land cover type is "present"), the probability of misclassifying land cover in all pixels is close to zero. Spatial land cover classifications and household survey land cover data are likely to operate at different levels (pixel vs. plot). For instance, respondents may not identify a plot with a small patch of grass as having grass cover in the survey. This suggests future work should ensure survey questions have greater alignment with spatial data and take into account mixed land cover/use in a single plot (e.g., asking about the proportion of the plot with grass cover).

Table C1: Plot land cover congruence by land cover classification<sup>a</sup>

	<b>Vegetated cropland (1)</b>	<b>Trees (2)</b>	<b>Non-vegetated cropland<sup>b</sup> (3)</b>	<b>Grass fields (4)</b>
Self-reported plot land cover (# plots)	450	75	51	55
Spatially identified land cover (# plots)	474	258	267	457
% match	95	29	19	12
$\phi$ coefficient	0.01	0.11	0.14	0.04

<sup>a</sup> Data were limited to high quality plots (n=650). Spatially identified land cover indicates the presence of each land cover type on a given plot, with multiple land covers possible on a single plot.

<sup>b</sup> This classification includes temporarily bare cropland/recently planted or harvested cropland

DRAFT