# A user-generated data based approach to enhancing location prediction of financial services in sub-Saharan Africa

Grant McKenzie<sup>a,b</sup>, R. Todd Slind<sup>b</sup>

<sup>a</sup>Department of Geography, McGill University, Montreal, Canada <sup>b</sup>Spatial Development International, Seattle, USA grant.mckenzie@mcgill.ca, tslind@spatialdev.com

# Abstract

The recent increase in user-generated content and social media adoption in developing countries offers an unprecedented opportunity to better understand the accessibility and spatial distribution of financial services in sub-Saharan Africa. Financial inclusion has been identified as a priority by multiple agencies in the region and on-the-ground efforts are currently underway to identify previously unknown financial access points in numerous developing African countries. Existing techniques for estimating the location of these access points rely on spatial analysis of often outdated or unsuitable publicly available datasets such as population density, road networks, etc., as well as expensive and time consuming surveys of locals in the region. In this work we propose an approach to augment existing spatial data analysis techniques through the inclusion of user-generated geo-content and geo-social media data. Through a comparison of standard regression models and machine learning techniques, this work proposes the use of alternative data sources to build prediction models for identifying financial access locations in countries where current estimation models are insufficient. With a better understanding of geospatial distribution patterns this work aims at reducing data acquisition costs and providing decision makers with critical data more quickly and efficiently. Finally, we present a mobile application built on the outcomes of this analysis that is currently being used to better inform on-the-ground data collection efforts.

*Keywords:* financial access, user-generated content, social media, Kenya, Uganda, random forest

Preprint submitted to Applied Geography

February 6, 2019

#### 1 1. Introduction

By current estimates, the number of individuals in sub-Saharan Africa 2 (SSA) with bank accounts at formal financial institutions is 25% [1], a num-3 ber that has remained relatively stagnant, growing by only a couple of per-4 centage points over the past four years [2]. By comparison, mobile money 5 accounts in East African countries, especially Kenya and Tanzania, have in-6 creased dramatically. The term *mobile money* here represents the use of 7 mobile devices to transfer money between users, pay bills, or purchase items. 8 Mobile money *providers* are those companies through which an individual 9 deposits or withdraws local currency to or from their mobile money account. 10 Mobile money providers are typically fixed-location, corner stores to which a 11 customer can go to exchange currency for *mobile money* (see Figure 1 for an 12 example). Safaricom, a leading Kenyan mobile network operator, launched 13 a mobile device-based payment system called M-Pesa in 2007 that revolu-14 tionized financial transactions across much of East Africa. In 2016, it was 15 estimated that mobile device penetration in Kenya surpassed 90%, an in-16 crease of over 6% in one year [3]. And while only a small portion of the 17 Kenyan population have traditional bank accounts, over 58% percent of in-18 dividuals in Kenya use mobile money [4] to transfer funds between people 19 and/or businesses or borrow money by way of a loan [5]. Mobile money has 20 such a dominant role in the Kenvan economy that in 2014 M-Pesa, by far 21 the leading mobile payment system, accounted for over 60% of the country's 22 gross domestic product [6]. 23

While the rise of mobile money has shown to reduce poverty rates [7] and 24 increased gender equality in many developing nations [8], there are concerns 25 over economic impact [9], taxation [10], and the influence of a single mo-26 bile network operator. The external focus on the striking growth in usage 27 of mobile money has also served to magnify the financial divide within the 28 country. During the FinAccess 2014 conference Njuguna Ndung'u, Governor 29 of the Central Bank of Kenya, gave a keynote address in which he encouraged 30 the expansion of financial inclusion in Kenya [11]. In this keynote, Professor 31 Ndung'u reiterated that while a considerable portion of the Kenyan popula-32 tion has access to mobile money infrastructure, a quarter of the population 33 remains entirely excluded. With the goal of increasing financial inclusion, the 34 Central Bank of Kenya, specified that a first step should include the iden-35



Figure 1: An example of a mobile money provider in Uganda. Source: Wikimedia Commons. License: CC 4.0

tification of all Financial Touch Points  $(FTP)^1$  within the country. While 36 there are on going efforts to collect location information on FTP providers 37 in Kenya [12, 13], the turn-over rate and movement of providers within the 38 country are high. In actuality, the locations of many FTP are still not known. 30 Efforts to better understand the distribution of financial services in Kenya 40 are on-going. These are focused on the spatial distribution of mobile money 41 infrastructure to identify opportunities for business expansion, agricultural 42 services, etc. [14, 15, 16]. On-the-ground data collection efforts continue 43 in SSA regions with the Humanitarian OpenStreetMap Team [17] following 44 other teams such as Brand Fusion [12] in their data collection efforts. Most of 45 these on-the-ground efforts involve canvassing entire countries on motorcycles 46 with GPS units in an attempt to identify new FTP locations or view the 47 identification of FTP as a secondary goal to mapping a country. Collectors 48 focus their efforts on highly populated regions, surveying locals and known 40 FTP providers [18]. In general though, there is a lack of informed strategy on 50

<sup>&</sup>lt;sup>1</sup>These include mobile money providers, brick and mortar banks, etc.

where to look for these financial touch points in the most efficient manner. 51 Population density maps and local knowledge are an important step and 52 our goal is that the methods proposed in this work can be used to augment 53 existing ones. To this end, this work aims to build a model for predicting the 54 location of financial touch points based not only on population densities, but 55 other publicly available datasets, both traditional authoritative (e.g., land 56 use, school locations) and user-contributed (e.g., volunteered information 57 and social media). 58

In the last year, the number of smartphone users in SSA has grown sub-59 stantially. The percentage of users in Kenya with smartphones was roughly 60 44% in 2016, a substantial shift from the previous year of 27% [19]. This 61 growth in smartphone access has also given rise to a substantial increase in 62 social media usage. Recent reports show social media usage at 58% of the 63 most popular activities conducted with a mobile device followed by search 64 engines at 39% and email at 30% [19]. Facebook, one of the most popular 65 social media platforms in the world has recently focused their attention on 66 SSA as a region for expansion [20]. These efforts are paying off with recent 67 statistics showing that 170 million Africans have joined Facebook, most of 68 which connect through their mobile device [21]. Of these, 6.1 million are from 69 Kenya. [22]. Twitter, has also seen an increase in adoption with monthly ac-70 tive users counted at roughly 2.2 million [23]. As users interact with these 71 platforms, they contribute significant amounts of digital content. This con-72 tent ranges from photographs and opinions to restaurant reviews and group 73 chats. The fact that much of this interaction happens via mobile device is 74 of importance as well. Many smart devices contain high resolution location 75 sensors such as GPS or Wi-Fi and social media applications make use of 76 this information which lead to social contributions that contain geographic 77 data such as places, local businesses and geotagged social posts. Through 78 the various application programming interfaces (APIs) offered by these plat-79 forms, researchers now have access to much of this published content. The 80 resolution of these data both spatially and temporally offer unique insight 81 into the behavior of individuals within the region. Not only can these data 82 be used to enhance low resolution (and often outdated) population density 83 maps but contributions such as those that mention local businesses can be 84 used to better predict the location of previously unmapped entities, such as 85 mobile money providers and other FTP. 86

Social media data are often defined as a subcategory of user-generated content (UGC), one that may contains geographic information, but is often

not contributed explicitly with the geographic content in mind [24]. Another 89 source of UGC common to the geography domain is volunteered geographic 90 information (VGI) [25]. One of the popular platforms for this type of in-91 formation is OpenStreetMap<sup>2</sup>, a rich set of geospatial data contributed to, 92 and curated by, thousands of citizens worldwide. In recent years there have 93 been substantial efforts to increase coverage and quality of geographic data 94 and maps in SSA.<sup>3</sup> These data in many cases are more up-to-date and have 95 greater coverage than many government or commercial geographic datasets 96 and knowing this, we propose their inclusion in our approach to predicting 97 financial access location in Kenya. 98

## 99 Research Contribution

The purpose of this work is to develop a method for predicting financial touch points in Kenya. Specifically, we are interested in determining if at least one FTP can be identified within a specific set of grid cells. Building on traditional authoritative datasets, we examine the fitness of emerging data sources for inclusion in an FTP prediction model and ultimately as a layer in a mobile application for data collection. To this end we address the following four research questions (RQ).

*RQ1* With the goal of identifying financial touch points in Kenya, how do
geo-tagged social media and volunteered geographic information fare
in comparison to authoritative datasets? To address this question,
we explore the distribution and correlation of various datasets with
known FTP in Kenya. We report on the accuracy of using these data
independently for estimating FTP counts and locations.

RQ2 Can social media data and volunteered geographic information be used 113 in combination with existing authoritative datasets to produce bet-114 ter FTP prediction models than those generated from the datasets in-115 dependently? Here we examine two traditional regression methods, 116 namely ordinary least squares and spatial lag as well as two machine 117 learning regression approaches, namely support vector regression and 118 random decision forest (RDF). The accuracy of these models are re-119 ported via three measures. 120

<sup>&</sup>lt;sup>2</sup>http://openstreetmap.org

<sup>&</sup>lt;sup>3</sup>https://hotosm.org/projects

*RQ3* Provided a best fit model, can we validate this approach through onthe-ground identification of previously unknown FTP? Secondly, how
accurate is the best fit model in identifying FTP in Kenya's neighboring
country of Uganda? We assess and report on the accuracy of the model
and identify important differences between the two countries that likely
impact the accuracy of the model.

RQ4 Can the FTP prediction model provide the foundation of a mobile application for FTP data capture and validation? We present a prototype
mobile application currently employed by users on-the-ground to add,
edit and delete FTP locations, driven by an FTP prediction layer generated from our best fit model.

The remainder of this article is organized as follows. In Section 2 we discuss existing research related to the topic and methods, and in Section 3 we present the various datasets used in this work. The methods used in predicting financial touch points are given in Section 4, with the results of the analysis shown in Section 5. Two different approaches for validating the data set are presented in Section 6 with an overview of the mobile application in Section 7. Finally, conclusions and next steps are stated in Section 8.

# 139 2. Related Work

Existing work in this area has highlighted the importance of understand-140 ing mobile financial services in sub-Saharan Africa specifically as it relates 141 to poor populations [26, 27]. Some of this research has used data collected 142 directly from mobile devices [28] while others have focused on the broader 143 impact of the technology [29]. Mobile money usage is not unique to sub-144 Saharan Africa. Many other countries have adopted mobile money systems, 145 China being one of the leading proponents of the technology [30]. Recent 146 reports have shown that payment systems suck as Alipay and WeChat pay 147 are having significant impacts in shaping the country's economy [31]. In re-148 cent years, the focus has shifted from the availability of mobile devices to 149 the actual usage patterns and applications. Short messaging service (SMS) 150 and social media usage have grown substantially and are having a sizable 151 impact on the developing world for everything from political movements [32] 152 to monitoring and tracking health epidemics (e.g., Ebola) [33]. 153

As social media usage and user-generated content grows in developing 154 countries, so does that availability of geotagged content [34]. The devel-155 opment of crowd-sourcing crisis tools such as Ushahidi [35] and Missing 156 Maps [36] have successfully demonstrated that geotagged social content can 157 have a substantial impact during crisis relief efforts. Recent work by Adams 158 et al. [37] has also shown that user-generated geo-tagged content from travel 159 blogs and Wikipedia articles can be used to identify thematic regions around 160 the world further emphasizing the power of crowd contributions. Exist-161 ing work by Linard et al. [38] has examined the inclusion of volunteered 162 geographic information in enhancing the WorldPop dataset. Their efforts 163 demonstrated that OpenStreetMap vector data can be used to combination 164 with satellite imagery to further refine global population estimates. Further 165 work has used a combination of VGI-based gazetteer data and social me-166 dia 'check-ins' to determine citizen locations [39] and prioritize evacuation 167 zones [40]. 168

From a methodological perspective, machine learning regression models 169 have been quite successful in a variety of scenarios. The range of literature in 170 this area speaks to the complexity and variety of models. Previous work on 171 the role of spatial autocorrelation in standard regression [41] is making it's 172 way into machine learning (e.g., SVM, RDF, etc.) discussions [42]. Existing 173 work from Song et al. [43] compared spatial econometric models to a random 174 decision forest approach in modeling fire occurrence and demonstrated the 175 benefits and disadvantages of the different approaches. Stevens et al. [44] 176 employed a RDF model in disaggregating census data for population map-177 ping with the goal of enhancing the WorldPop dataset and recent work on 178 identifying landscape preferences determined that an RDF approach applied 179 to Flickr photos produced the best results [45]. 180

# 181 **3. Data**

In this section, we provide an overview the datasets used in constructing the FTP identification models. The financial touch points are introduced as well as the predictors classified as VGI, Social Media, and Authoritative datasets.

# 186 3.1. Financial Touch Points

On-the-ground data collection efforts by  $Brand Fusion^4$  resulted in a 187 dataset of verified FTP in Kenya [12]. Brand Fusion estimates that these 188 data, collected in 2015, represent a high portion of all FTP within Kenya 189 but the data are non-exhaustive as FTP may have been missed by data 190 collectors, locations may have been established since the last round of data 191 collection, or FTP may have moved. The purpose of this paper in this case is 192 to use geospatial indicators near to these known FTP to predict and identify 193 previously unidentified FTP in Kenya. This 2015 Brand Fusion dataset iden-194 tified 83,273 FTP in Kenya and these form the basis on which our prediction 195 model is trained and tested. Figure 2 shows the distribution of these FTP in 196 Kenya as green markers. The Humanitarian OpenStreetMap Team (HOT) 197 collected FTP for neighboring Uganda [17]. In total, 45,417 verified FTP 198 were identified in Uganda and these points will form the basis of our follow-199 on analysis. Visually, the highest density of FTP appear to occur in densely 200 populated regions around Nairobi, Nyanza (Kenya), Kampala and Mbarara 201 (Uganda). Spatial analysis of these FTP locations through Moran's I [46] 202 and Ripley's K [47] functions confirm this, indicating clear spatial clustering 203 within these datasets. While the high population areas show the highest 204 numbers of FTP, it is the rural regions that are of particular interest to 205 government and non-government agencies. 206

# 207 3.2. Predictors

We compare and contrast a number of different datasets from a wide variety of sources with the purpose of determining how the inclusion of these data aid in predicting FTP locations. Table 1 lists these datasets along with their sources and our assigned category tag. These categories consist of two types of user-generated content, namely *volunteered geographic information* (*VGI*) and *social media* (*SM*) as well as more traditional datasets which by comparison we label *authoritative* (*AUTH*).

# 215 3.2.1. Authoritative Datasets

We define the *authoritative datasets* in this work as those not created through direct citizen contributions or social media data extraction. These datasets were generated using more authoritative and controlled mechanisms

<sup>&</sup>lt;sup>4</sup>http://www.brandfusion-africa.com/services/mobile-money



Figure 2: Financial Touch Points (FTP) in Kenya (83,273) and Uganda (45,417). Base map by ESRI.

and are therefore, allegedly, less prone to user bias or classification error. These data have been used in numerous other studies in estimating everything from population density and land use to human mobility and predicting disease outbreak [48, 49, 50, 33].

The 2015 WorldPop data contains high resolution (~100m cell size) hu-223 man population distribution estimates. The data was generated from a com-224 bination of remote sensed imagery, census and existing geospatial datasets 225 (e.g., road networks) [44, 51]. The Socioeconomic Data and Application 226 Center in NASA's Earth Observation System Data and Information System 227 group produces the Global Rural-Urban Mapping Project (GRUMP) data. 228 Similar to the WorldPop dataset, these data are produced through a combi-229 nation of census and satellite data (including night-time lights) at a resolution 230 of roughly 1km. Version 1 of this dataset was produced in 2011 and provides 231 rural and urban population density estimates for the year 2015 [52, 53]. Urban 232 land cover type regions were also extracted from the 0.5 km MODIS-based 233

Dataset Description	Source	Year	Category
Estimated persons per 3 arc-second (roughly 100m) cell	Worldpop	2015	AUTH
Primary & Secondary School Locations	OpenAfrica	2015	AUTH
0.5 km MODIS-based Global Land Cover Climatology	USGS	2014	AUTH
Global Rural-Urban Mapping Project (GRUMPv1)	NASA	2011	AUTH
GeoNames Places	GeoNames	2016	AUTH
LandScan-based Populated Places	Natural Earth	2016	AUTH
OSM Roads	OpenStreetMap	2016	VGI
OSM POI	OpenStreetMap	2016	VGI
Facebook Places	Instagram API	2016	SM
Tweets	Twitter API	2016	SM
Foursquare Venues	Foursquare API	2016	SM

Table 1: Datasets used in identifying financial touch points.

<sup>234</sup> Global Land Cover Climatology dataset [54] generated in 2014.

Dataset	Kenya	Uganda
Facebook Places	8107	4377
Twitter Tweets	204538	156426
Foursquare Venues	4016	2075
OpenStreetMap POI	16739	44203
OpenStreetMap Roads (km)	98381	48676
Schools (Primary & Secondary)	37317	29372
GeoNames Places	26038	25978
NE Populated Places	56	42

Table 2: Counts for the predictor datasets in Kenya and Uganda. Note that both the WorldPop and GRUMPv1 data are not count based datasets and so are not reported here.

Primary and Secondary school locations were accessed from OpenAfrica, 235 a web portal for open data in African countries. School locations for Kenya 236 were most recently updated in 2015 and contributed by the Kenya Open 237 Data Initiative [55]. Similarly, school locations for Uganda were collected by 238 the Uganda Bureau of Statistics and the Ministry of Education and Sports 239 from 2004–2010. Places were downloaded from the *GeoNames* placename 240 gazetteer which is made up of a number of sources, most notably the National 241 Geospatial-Intelligence Agency and the U.S. Board on Geographic Names for 242 regions outside of the United States. This point data represents everything 243 from mountain tops to water wells. Natural Earth Populated Places data 244 were used in this research which is based on LandScan-derived population 245 estimates [56]. Natural Earth devised the dataset based on regional signifi-246 cance of places over population census, differentiating it from the grid-based 247

systems previously mentioned [57]. Counts of these datasets are shown inTable 2.

#### 250 3.2.2. User-contributed Data

User-contributed data are those created either via volunteered geographic information (VGI) means or social media (SM) contribution. Typically contributions to these data are made from non-experts and do not rely on statistical models built from existing data sources. Anyone can add a place, venue, road, or post (tweet) to one of these datasets without requiring secondary approval.<sup>5</sup>

#### <sup>257</sup> Volunteered Geographic Information

*OpenStreetMap* Points of Interest were downloaded for Kenya using the 258 OsmPoisPbf extraction tool.<sup>6</sup> Table 2 lists the total number of POI with 250 roughly 2% (339) of these being tagged as MONEY BANK or MONEY EX-260 CHANGE. On examination of these tagged POI, the overwhelming major-261 ity of these were brick-and-mortar bank branches with few mobile money 262 providers or lenders. These mobile money providers and lenders are ei-263 ther corner stores / grocers or dedicated shops (e.g., M-Pesa). The Open-264 StreetMap Road data was also extracted in January 2016 and consists of 265 high resolution road network data contributed by volunteers. These data 266 are notably of a higher resolution and wider spatial coverage than the road 267 network datasets available from the Kenyan government GIS web portal. 268

# 269 Social Media Data

Social media data for this research involved three sources of geotagged content. Instagram and Foursquare both have digital gazetteers of place locations contributed by individuals while twitter allows contributors to geotag their posts with geospatial coordinates.

The Instagram locations API<sup>7</sup> was used to extract Points of Interest for Kenya. Instagram uses *Facebook Places* as it's gazetteer, with the purpose of allowing individuals to tag their photographs with a place name. Their API offers limited access to this gazetteer. In total, 8107 places were accessed in

<sup>&</sup>lt;sup>5</sup>Note that there is a community-based validation process in OpenStreetMap <sup>6</sup>https://github.com/MorbZ/OsmPoisPbf <sup>7</sup>https://www.instagram.com/developer

Kenva. The *Twitter* Streaming API<sup>8</sup> was used to access geotagged tweets 278 within Kenya over a 5 month time span from January through May 2016. 279 Only those tweets that included precise geographic coordinates and sourced 280 from the Android Twitter App or iPhone Twitter App were employed here. 281 In this work, only the geographic location of the tweets was relevant for this 282 research though future work may explore the content and language variation 283 within the text of the tweets. The Foursquare Venues Search  $API^9$  was 284 employed to access Points of Interest in the Foursquare gazetteer. Foursquare 285 began curating POI in March of 2009 and has been more transparent in how 286 they collect places [58] than Facebook. Notably Facebook has a much larger 287 user-base (2 billion vs. 45 million) however. 288

# 289 4. Methods

To start, a spatial grid was generated over the entire country of Kenya 290 at a resolution of 0.02 degrees, or approximately 2.2 km at the equator. 291 Selection of this resolution was based on trade off between reasonable travel 292 time within each grid (for on-the-ground collection efforts and actual FTP 293 users) and reduced computational complexity. This resulted in 120,111 grid 294 cells across Kenya. The grid was intersected with the FTP data producing an 295 FTP grid layer with aggregated count cells ranging in value from 0 to 2.402 (in 296 Nairobi). Similar layers were constructed for each of the predictor variables 297 using the same grid bounds and resolution. Finally, each gridded layer was 298 normalized to between 0 and 1. This was to ensure that each variable could 299 be compared to one another without one predictor overpowering the others. 300 While not essential in a linear or spatial regression model, it is particularly 301 important for a random decision forest approach [59]. 302

## 303 4.1. Individual predictors

The goal in the initial analysis for *RQ*1 is to determine how accurate each individual dataset is in identifying FTP. We first examine the correlation between each gridded dataset and the gridded FTP layer. Table 3 shows the Spearman's correlation matrix of all predictors. Notably, all datasets show positive correlation with the number of FTP per cell. The Worldpop, Grump and School datasets show the highest correlation with Facebook Places also

<sup>&</sup>lt;sup>8</sup>https://dev.twitter.com/streaming

<sup>&</sup>lt;sup>9</sup>https://developer.foursquare.com

showing a reasonably high value. Interestingly tweets have a relatively low correlation with FTP (0.11) and an even lower correlation with the other social media / user-generated content datasets (e.g., 0.05, 0.02) indicating that there is little similarity between our social media places and the geotagged tweets. On the other hand, GRUMP data are highly correlated with the WorldPop dataset.

We then calculate the *F*-score for each predictor against the FTP. *F*-score measures the relationship between the precision and recall of these datasets (Equation 1). *Precision*, in this case, is the number of FTP locations correctly identified divided by the total number of locations identified whereas *recall* is the number of FTP locations correctly identified divided by the total number of actual FTP locations.

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} \tag{1}$$

Assessing the accuracy of a predictor via the F-score involves a trade-322 off. Figure 3 shows precision versus recall for each of the predictor variables. 323 Notably, the authoritative datasets show a steeper decrease in recall as pre-324 cision drops below 0.4 whereas the user-contributed datasets tend to show 325 fairly low trade-offs between the precision and recall. The highest F-score 326 of 0.49 is found with the WorldPop data and a low of 0.07 with the Nat-327 ural Earth Populated Places location data (Table 4). While these F-scores 328 in combination with the correlation matrix show that the predictor datasets 329 are of value in estimating FTP locations, on their own they only correctly 330 identify a limited number of FTP in Kenya. 331

#### 332 4.2. Weighted combination of variables

Provided the accuracy of the predictors independently, we next explore a 333 number of methods for combining the predictors in order to better identify 334 the location of financial touch points in Kenya. Specifically, to address RQ2335 we test four approaches to FTP identification, namely ordinary least squares 336 regression, spatial lag regression, support vector regression, and random de-337 cision forest. The purpose of examining all of these methods is to determine 338 which approach most accurately predicts the location of known FTP and 339 produces a model on which to base further investigation into unknown FTP 340 locations. To be clear, the regression approaches produces probability values 341 that are used to in a prediction task of FTP in a grid cell or not. These 342

	FTP	Facebook	Foursquare	Twitter	IO4 MSO	OSM Roads	Schools	GRUMP	WorldPop	Land. Urban	GeoNames	ΣE
FTP	1.00	0.50	0.29	0.11	0.33	0.08	0.50	0.55	0.57	0.38	0.27	0.29
Facebook Places	0.50	1.00	0.38	0.05	0.13	0.14	0.39	0.27	0.29	0.39	0.15	0.37
Foursquare Venues	0.29	0.38	1.00	0.02	0.05	0.08	0.22	0.12	0.14	0.24	0.05	0.21
Twitter Tweets	0.11	0.05	0.02	1.00	0.02	0.01	0.08	0.19	0.18	0.07	0.03	0.02
IO4 MSO	0.33	0.13	0.05	0.02	1.00	0.03	0.32	0.54	0.52	0.20	0.63	0.04
OSM Roads	0.08	0.14	0.08	0.01	0.03	1.00	0.26	0.17	0.15	0.09	0.14	0.03
Schools	0.50	0.39	0.22	0.08	0.32	0.26	1.00	0.71	0.72	0.38	0.43	0.12
GRUMP	0.55	0.27	0.12	0.19	0.54	0.17	0.71	1.00	0.92	0.45	0.52	0.07
WorldPop	0.57	0.29	0.14	0.18	0.52	0.15	0.72	0.92	1.00	0.44	0.52	0.10
Landuse Urban	0.38	0.39	0.24	0.07	0.20	0.09	0.38	0.45	0.44	1.00	0.16	0.19
GeoNames Places	0.27	0.15	0.05	0.03	0.63	0.14	0.43	0.52	0.52	0.16	1.00	0.03
NE Pop. Places	0.29	0.37	0.21	0.02	0.04	0.03	0.12	0.07	0.10	0.19	0.03	1.00

Table 3: Spearman's correlation between Kenya dataset cell counts. All p<0.01.



Figure 3: Precision vs. recall graphs for all independent variables.

343 probability values are later used in the generation of a *prediction layer* for 344 inclusion in a mobile data collection application.

## 345 4.2.1. Ordinary Least Squares Model

A standard linear regression model was executed as a first step to deter-346 mine the impact of each independent variable (predictor dataset) on identify-347 ing FTP. The data were separated by category as shown in Table 1, namely 348 VGI, SM, or AUTH. Linear regression models were constructed for each 349 category independently as well as combined. The independent variables, co-350 efficients,  $R^2$ , and residual standard error (RSE) for each model are shown in 351 Table 5. Regarding multicollinearity between the independent variables, we 352 note some small changes in the regression coefficients as predictors are added 353 to the model. The most notable change here is in the OpenStreetMap POI 354

Dataset	Max F-score
Schools	0.38
GRUMP	0.44
WorldPop	0.49
Landuse Urban	0.12
GeoNames Places	0.31
NE Populated Places	0.07
Facebook Places	0.40
Foursquare Venues	0.23
Twitter Tweets	0.33
OSM POI	0.29
OSM Roads	0.21

Table 4: Maximum F-score values for each of the predictor variables independently.

dataset changing to having a negative influence on FTP identification when 355 combined with all other datasets. Similarly, we see the *tweets* dataset change 356 from having a significant impact on the model to not longer being significant. 357 We calculated the condition indices (condition number test), measures of ill-358 conditioning in the predictor matrices and found that the regression models 359 did not have significant multicollinearity. The conditional index values for 360 the respective regression models are 5.87 (AUTH), 1.53 (SM), 1.77 (VGI), 361 7.43 (Combined). 362

The AUTH-based regression produced an  $\mathbb{R}^2$  value of 0.412 with all co-363 efficients being significant (P < 0.001). Based on the coefficients, the World-364 Pop density values had the highest positive influence on the dependent FTP 365 variable with GRUMP data also showing a high value of influence. The 366 GeoNames places dataset had a small, but negative influence on the model. 367 The SM-based regression model produced a lower  $\mathbb{R}^2$  value meaning that 368 less of the known FTP locations could be explained by our place-based 369 and geotagged social media data. All coefficients were deemed significant 370 with Facebook places and Tweets producing larger positive coefficients than 371 Foursquare venues. The VGI-based linear regression models produced the 372 lowest R<sup>2</sup> value with OpenStreetMap POI having a much larger influence on 373 the model than OpenStreetMap Roads. Combining all independent variables 374 in one OLS linear regression model produced the highest  $\mathbb{R}^2$  value with all 375 coefficients having a significant impact with the exception of tweets and the 376 lowest residual standard error of the OLS models. As a first, but important, 377

step, these results are encouraging and indicate that a combination of social media, VGI and authoritative data produce better results for predicting
financial touch points than each data type independently.

Dataset	OLS Model Coefficient	Spatial Lag Model Coefficients
Authoritative Datasets (AUTH) Model		
Schools	3.80E-02	6.09E-02
GRUMP	9.56E-02	4.76E-02
WorldPop	2.23E-01	1.77E-01
Landuse Urban	1.06E-02	7.59E-03
GeoNames Places	-4.58E-02	-3.45E-02
NE Populated Places	5.54E-02	5.70E-02
Spatial Lag (Rho)	NA	2.33E-01
	$R^2$ 0.412, RSE 4.32E-03	$R^2$ 0.425, RSE 4.265E-03
Social Media Datas	ets (SM) Model	
Facebook Places	1.58E-01	1.33E-01
Foursquare Venues	5.55E-02	5.24E-02
Twitter Tweets	1.41E-01	3.39E-02
Spatial Lag (Rho)	NA	5.48E-01
	$R^2$ 0.267, RSE 4.82E-03	$R^2$ 0.423, RSE 4.27E-03
Volunteered Geogra	phic Information Datase	ts (VGI) Model
OSM POI	3.54E-01	2.34E-01
OSM Roads	8.57E-04	4.28E-04
Spatial Lag (Rho)	NA	5.52E-01
	$R^2$ 0.116, RSE 5.29E-03	$R^2$ 0.285, RSE 4.76E-03
Combined (All data	a) Model	
Schools	2.78E-02	3.05E-02
GRUMP	9.25E-02	5.22E-02
WorldPop	1.94E-01	1.60E-01
Landuse Urban	2.00E-03	-8.92E-04
GeoNames Places	-9.83E-02	-8.23E-02
NE Populated Places	3.04E-02	3.21E-02
Facebook Places	9.55E-02	9.59E-02
Foursquare Venues	4.30E-02	4.38E-02
Twitter Tweets	3.08E-02*	-8.72E-03*
OSM POI	-6.13E-04	1.02E-01
OSM Roads	1.05E-01	-6.23E-04
Spatial Lag (Rho)	NA	2.49E-01
	$R^2$ 0.489, RSE 4.02E-03	$R^2$ 0.502, RSE 3.97E-03

Table 5: Results of the OLS and Spatial Lag regression models with four combinations of predictor variables. All coefficients are significant (p < 0.001) except for Twitter OLS\* which is not significant and Twitter SLM\* with p < 0.05.

#### 381 4.2.2. Spatial Lag Model

Using the *Jarque-Bera* test [60], the variables in the OLS models were assessed for normality of the distribution of errors. All probability values for the tests were very low indicating non-normal distribution of the error terms. Our next step was to geospatially map the residuals of our best-fit

linear regression model in order to test for spatial autocorrelation in our pre-386 dictors. Visually, the residuals appeared to show a clear spatial pattern with 387 underestimation occurring near major cities such as Nairobi and overestimat-388 ing in more rural regions to the North. Moran's I analysis of the residuals 389 supported this assessment with significant global values of 0.305, 0.266, and390 0.100 for SM, VGI, and AUTH models respectively, with a distance threshold 391 of 0.02 degrees (distance to the nearest grid cell). Local Moran's I analysis 392 also found highly significant spatial clustering around the high density FTP 393 regions, predominantly major cities. These results, combined with low prob-394 ability values from Breusch-Pagan tests [61] for heteroskedasticity indicate a 395 need to account for spatial autocorrelation in our regression analysis. 396

A spatial lag [62] regression model (Equation 2) was constructed rely-397 ing on an Euclidean distance weighted matrix using Queen contiguity at a 398 threshold of 0.02 degrees. Y represents the vector of response variables,  $\rho$  the 390 coefficients of spatial regression terms, making WY the spatial lag weighted 400 response. X is the matrix of independent predictors,  $\beta$  the coefficient matrix 401 of X and  $\epsilon$  the random error vector. The results of the Spatial Lag regression 402 models for the 3 groups of predictor variables and the combined model are 403 shown in Table 5. 404

$$Y = \rho W Y + \beta X + \epsilon \tag{2}$$

In all cases, there was an increase in the amount of variance explained 405 (R-squared) over the OLS regression models, and a relative decrease in the 406 standard error of the residuals. The WorldPop population dataset still had 407 a large influence in the combined dataset model (based on the coefficient 408 value) while Tweets remained low in contribution and significance. The spa-409 tial lag (Rho) coefficients all had significant impacts on the respective models 410 demonstrating that accounting for spatial dependency in such a model pos-411 itively influenced the ability to predict FTP in Kenya. These results again 412 indicate that combining datasets from various user-generated and authori-413 tative sources positively influence the ability to predict FTP and that the 414 inclusion of a spatial lag term positively contributes to an explanation of the 415 variance in our model. 416

#### 417 4.2.3. Support Vector Regression

<sup>418</sup> Support vector machine (SVM) analysis takes a different approach to <sup>419</sup> prediction than the previous two analyses. SVM is nonparametric and ap-<sup>420</sup> proaches regression through a kernel function [63, 64]. To start, we used an

epsilon ( $\epsilon = 0.1$ ) type of regression with a linear kernel.<sup>10</sup> This approach 421 attempts to find a separating hyper-plane between the two classes, in our 422 cases occurrence of FTP in a grid cell or not, with a maximum gap between. 423 In general, SV regression perform better with a higher number of dimen-424 sions, or predictor variables in our case, and really only if the combination 425 of these variables almost certainly leads to a known FTP. In our cases, nei-426 ther of these conditions hold true as the number of datasets (dimensions) 427 is relatively small and based on our previous OLS and spatial lag analysis, 428 the variance explained is low. While this form of analysis was tested on our 420 dataset, it primarily acts as a first *comparison* step in a machine learning 430 approach to this problem. 431

# 432 4.2.4. Random Decision Forest

Random decision forests (RDF) [65] are an ensemble learning method for regression, in our case, that construct a set of decision trees for the purpose of prediction. An optimal threshold value for identifying the occurrence of an FTP or not in a grid cell is calculated. A random forest aims to correct for overfitting, known to happen in a standard decision tree approach [66].

Dataset	IncNodePurity
GRUMP	2.32E-02
WorldPop	2.29E-02
Schools	2.22E-02
Landuse Urban	5.23E-03
Twitter	1.06E-02
OSM POI	9.62E-04
Geonames	4.27E-03
Facebook	1.69E-02
OSM Roads	5.98E-05
NE Major Towns	1.00E-03
Foursquare	4.75E-03

Table 6: Incremental Node Purity of the variables in the random decision forest model.

The *R* RandomForest package<sup>11</sup> was used to fit a random decision forest regression model to the FTP data based on each of the category predictor variables independently as well as all together. This resulted in a  $1.39 \times 10^5$ mean of squared residuals explaining 55% of the variance. This approach used 500 trees with 4 variables tried at each split. The incremental node

 $<sup>^{10}\</sup>mathrm{R}$  package: https://cran.r-project.org/web/packages/e1071

<sup>&</sup>lt;sup>11</sup>https://cran.r-project.org/web/packages/randomForest/

<sup>443</sup> purity for the model is shown in Table 6 and reports on the average change <sup>444</sup> of impurities of a tree node (in which the variable was used) before and after <sup>445</sup> a split. Plotting the percentage increase in mean square error (MSE) for <sup>446</sup> the combined approach (Figure 4) we find that many of the authoritative <sup>447</sup> datasets are the most important to the regression fit. Tweets, OSM POI and <sup>448</sup> Facebook places all positively contribute to the model, with OSM Roads, NE <sup>449</sup> Populated Places and Foursquare venues having little impact on the RDF fit.



Figure 4: Percentage increase in mean square error of prediction as a result of variable shuffling. In essence, the higher the value, the more important that variable is to the RDF regression model.

Given the known spatial dependency of the predictor variables (based 450 on global and local Moran's I measures), we elected to construct a separate 451 RDF model which included latitude and longitude coordinates as covariables. 452 There is some evidence in the existing literature that the inclusion of geospa-453 tial variables in such a model can influence the accuracy of prediction [42]. 454 Given the non-parametric nature of RDF, these variables could be included 455 in the model and used in the prediction assessment. This led to a slightly 456 higher percentage explained variance (0.56 vs. 0.55) and latitude was found 457 to be the second most important contributing variable as determined by 458 the percentage increase in mean square error. Again, though the prediction 459

method has changed substantially, the findings again support the fact that
user-contributed data are important in location prediction.

# 462 5. Results

In this section we present the results of the analyses performed in the 463 previous sections. Running each of the regression models (OLS, Spatial Lag, 464 SVM, and RDF) with datasets from each of our categories (VGI, SM, AUTH) 465 as well as a combination of all datasets (COMBO) produced a set of FTP 466 prediction values for each cell in our Kenya grid, 16 different FTP predic-467 tion grids. These regression-based prediction grids were each then compared 468 to our known FTP grid and three measures of accuracy were calculated for 469 each prediction. Table 7 shows a comparison of the four regression techniques 470 used in this work along with values for assessing accuracy of prediction in-471 cluding maximum F-score, Spearman's Correlation and root mean square 472 error (RMSE). The SVM and RDF methods also show results for regression 473 models that included all predictor variables as well as latitude and longitude 474 centroids of the grid cells. 475

In general, the random decision forest regression model approach pro-476 duced the best results across most categories. The RDF model that included 477 variables of all data categories, including latitude and longitude (LL) coor-478 dinates, produced the most accurate predictions as reported across all three 479 measures. A maximum F-score of 0.74 is quite high considering the multitude 480 of factors that may contribute to establishing an FTP. Similarly, a Spear-481 man's correlation of 0.96 is extremely high but should by understood in the 482 context of the sparsity of the FTP locations and predictions (most grid cells 483 are 0). Lastly, the reported RMSE is low relative to the comparable RMSE 484 values from all other methods and data categories. 485

Figure 5 further explains the F-scores for highest performing RDF model by plotting precision versus recall for the random decision forest models split by data category. In comparison to Figure 3, the combined approach of all datasets produces a much better trade-off between precision and recall, specifically addressing *RQ*2 as stated in the introduction.

Next, the residuals of the best-fit RDF regression model are mapped back
to the location data. Visual inspection identifies very little clustering within
the residuals and a Moran's I analysis confirms this with a bootstrapped
observed value of less than 0.001 implying a high degree of spatial randomness
in these RDF-based residuals.

Method	Category	Max F-Score	Correlation	RMSE
	VGI	0.31	0.340	5.29E-03
OIS	$\mathbf{SM}$	0.43	0.516	4.82E-03
OLS	AUTH	0.49	0.678	4.32E-03
	COMBO	0.51	0.699	4.02 E- 03
	VGI	0.36	0.429	5.13E-03
Spotial Law	$\mathbf{SM}$	0.42	0.518	4.82E-03
Spatial Lag	AUTH	0.49	0.637	4.34E-03
	COMBO	0.51	0.694	4.05E-03
	VGI	0.28	0.301	5.35E-03
SVM	$\mathbf{SM}$	0.35	0.417	5.17E-03
5 V IVI	AUTH	0.47	0.582	5.21E-03
	COMBO	0.55	0.590	5.17E-03
	COMBO & LL	0.56	0.587	5.17E-03
	VGI	0.31	0.606	4.69E-03
DDE	$\mathbf{SM}$	0.43	0.849	3.20E-03
	UGC	0.46	0.855	3.32E-03
	AUTH	0.57	0.930	2.25E-03
	COMBO	0.62	0.955	1.85E-03
	COMBO & LL	0.74	0.960	1.79E-03

Table 7: Prediction results of the regression methods split by category of dataset. The maximum F-score, Spearman's Correlation and root mean square error are reported. Note that all Spearman correlation values are significant (p < 0.01).

# 496 6. Validation

#### 497 6.1. Ground-truthing in Kenya

One primary goal of this work was to build a prediction model that 498 would allow researchers in the field to identify previously unidentified FTP 499 in Kenya. With this goal in mind we used the best fit random decision forest 500 model (reported in the previous section) to predict FTP locations across 501 Kenya. The predicted number of FTP locations was subtracted from the pre-502 viously known number of FTP per cell to produce a residuals map showing 503 the difference between known and predicted FTP. Of these residual cells, we 504 further investigated 47 that contained no known FTP and showed large neg-505 ative values (indicating high probability of finding FTP). Identifying these 506 locations with high potential is important as a single, previously unknown, 507 FTP could potentially be servicing a number of inhabitants; Inhabitants that 508 were thought to be without access to financial services. 509



Figure 5: Precision vs. Recall for Kenya RDF predictions.

These 47 potential FTP cells were ranked based on the size of the residual 510 and the latitude and longitude coordinates of the centroids were shared with 511 researchers on the ground in Kenya (see Figure 6). The selection of these 512 specific locations was also based on availability of data collection personnel 513 in the region around Eldoret city in eastern Kenya. Data collectors traveled 514 to these high-FTP-potential locations and recorded the presence and location 515 of any FTP they found within 1km radius of the cell centroid (represented as 516 square markers in Figure 6). In essence, the data collectors used the ranking 517 of residuals for binary classification (decision to travel to location or not) and 518 then counted the total number of FTP found within the vicinity of the marked 519 location. In total, 203 previously unidentified FTP were recorded within the 520 vicinity of these locations. In total, 41 of the 47 locations reported at least 521 one previously unknown FTP location within a 1.1 km radius. Assigning 522



Figure 6: Previously unknown FTP location (47) identified by the prediction model. Blue color density indicates rank based on probability of finding at least one FTP within 1.1km of the marked location.

the count of identified FTP to their nearest marked location (again, see 523 Figure 6) allowed us to compute the correlation between estimated FTP 524 potential and count of actual FTP identified. The resulting Spearman's 525 correlation was 0.233 (p < 0.01), a small but positive correlation indicating 526 that the magnitude of the residuals, not just the binary threshold, have a 527 role to play in FTP identification. It should be noted that a 1.1 km cell 528 radius is quite a large distance to explore and while quite a few new FTP 529 were identified, it is likely that other FTP may existed in the area but were 530 not identified. 531

The identification of these previously unidentified FTP offers validation to the RDF machine learning approach suggested in this research, and addresses RQ3. This approach presents a data-driven based method for uncovering previously unidentified FTP locations and has the potential to significantly reduces the on-the-ground efforts of individuals that previously relied on 537 qualitative assessment and brute force search methods.

## 538 6.2. Applicability to neighboring countries

In order to test the limits of our RDF prediction approach, the best-fit regression model constructed from numerous datasets in Kenya was applied to datasets collected in the neighboring country of Uganda. The countries of Kenya and Uganda, while similar in many ways, also differ substantially. We are currently in the process of collecting further on-the-ground data to test the transferability of this model to the neighboring country of Uganda.

In the mean time, our naive approach was again to rely on the same 545 publicly available datasets and use the best-fit model from the Kenya data 546 to predict locations and number of FTP in Uganda. Figure 7 graphs the 547 precision versus recall for three data categories independently as well as the 548 combined RDF regression model. Not surprisingly, the RDF model trained 549 on Kenya data produces poorer results in Uganda than Kenya. The F-scores 550 for the three data categories of SM, VGI and AUTH are 0.43, 0.44 and 551 0.36 respectively with a combined F-score of 0.44. The best Spearman's 552 correlation value was 0.69 for the combined model with a RMSE of 6.08E-553 03. In fact, just using OpenStreetMap POI data produced accuracy values 554 (F-score, Correlation and RMSE) similar to the combined RDF model built 555 from Kenya data. 556

There are numerous reasons for the drop in accuracy scores compared 557 to Kenya. The most obvious answer is that these are different countries 558 with unique economic, information & communications technologies (ICT), 559 and socio-demographic properties. It is naive to assume that a model built 560 on data from one country could be applied to a completely different country 561 without a loss of accuracy. Second, the FTP location data were collected 562 and reported by a different provider in Uganda than in Kenya (Humanitarian 563 OpenStreetMap vs. Brand Fusion). There are likely differences in the data 564 collection techniques, number of people involved and technology employed. 565 Future work will explore these differences with the purpose of identifying key 566 ways in which a model can be altered to account for regional differences. 567



Figure 7: Precision vs. recall for Uganda RDF Predictions.

# 568 7. Mobile Application

One of the outcomes of this research, and the focus of RQ4, is an Androidbased mobile application for identifying, creating, editing and deleting financial touch points within sub-Saharan Africa. The current prototype application functions both with and without a stable Internet connection and currently focuses on Kenya.

#### 574 7.1. Prediction overlay

<sup>575</sup> Based on the best-fit RDF prediction model developed in Section 4.2.4, <sup>576</sup> a raster layer containing FTP location predictions was constructed at a res-<sup>577</sup> olution of 0.02 degrees. This raster layer was styled on a white to green <sup>578</sup> color ramp using natural break classification and tiled to allow efficient data <sup>579</sup> transfer and visualization on the mobile mapping application (Figure 8a).



Figure 8: The FTP mobile prediction and capture application.

#### 580 7.2. Financial touch point locations

<sup>581</sup> Upon loading, the mobile application prompts the user to download <sup>582</sup> known FTP locations for one or more of Kenya's 70 districts. The purpose <sup>583</sup> of this is to allow a user to download only the data required, thus reducing <sup>584</sup> data usage and device storage. Before leaving an area of stable connectivity, <sup>585</sup> the user will download the known FTP locations for the district(s) in which <sup>586</sup> they will be traveling.

Users are invited to zoom and pan the map as they would on any standard 587 mobile mapping application (Figure 8c). The FTP locations are shown as 588 point markers on the map and clustered depending on zoom scale. When 589 the user selects a marker on the map, they are presented with the *Details* 590 interface. This interface shows information collected about the FTP by the 591 original party. The user can choose to edit this information (Figure 8b) 592 or delete the FTP entirely. Finally, the user has the option of zooming 593 into their current location on the map, either through panning/zooming or 594 selecting the *locate me* button. Once the map is at a reasonable scale, the 595 user can tap the map to add a new FTP. In this case, the unpopulated *Edit* 596 interface is presented to the user. Once the user is finished editing, adding 597

and deleting FTP, they have the option (selection from the context menu) to upload the changes to the database. Again, this allows for offline editing and reduces overhead of constant communication with the server whenever a FTP is edited. The application is currently in use by data collection teams in Kenya.

#### **8.** Conclusions & Future Work

In this work we present a novel approach to identifying financial touch 604 points in Kenya through combined use of geosocial media data, volunteered 605 geographic information, and authoritative geospatial datasets (RQ1 and RQ2). 606 We showed that we can significantly increase the ability to identify FTP lo-607 cations by including both spatial and platially tagged social media posts in 608 our analysis. Current state-of-the-art machine learning techniques were com-609 pared to existing ordinary least squares and spatial regression models and it 610 was shown that a random decision forest model using combined data from all 611 three sources best identified existing financial touch points and can be used 612 to identify the location of previously unknown FTP (RQ3). With this goal 613 in mind, we developed a mobile application for on-the-ground data collec-614 tion that uses the results of the RDF model as a geospatial estimation layer 615 through which users are be better informed on where to locate FTP (RQ4). 616 The application is currently in use in Kenya and has aided in the identi-617 fication of previously unknown financial touch points. Data collection done 618 using this application (with the inclusion of the prediction layer) has the po-619 tential to substantially impact financial services in countries such as Kenya 620 and Uganda. Provided detailed maps of access to financial services in sub-621 Saharan Africa, local government and international agencies are better in-622 formed when formulating policies and regulating financial services. The goal 623 of this work is to facilitate this discussion by providing access to the most 624 up-to-date geospatial data. 625

This analysis does come with some limitations. Given the country-level 626 analysis that was executed, a trade off was made when determining the cell 627 size for analysis. Increasing or decreasing this cell size would understandably 628 impact the accuracy of the identification model. Access to known FTP loca-629 tions is another limiting aspect of this type of analysis. Two different data 630 sets were collected from two different organizations in two different coun-631 tries. The methods of data collection varied and there is likely bias in how 632 the data was collected (e.g., accessibility of roads, daylight restrictions, etc.). 633

While these biases potentially impacted the final results of the analysis, they had little influence on the methods of analysis that were employed. A limitation of the validation approach lies in the lack of collected information related to true and false FTP negatives. Data collection teams in Kenya did not report on the lack of FTP in regions that were identified as not having FTP as it was not their primary mandate. Future data collection campaigns will aim to collect these data.

Future work in this area will continue to focus on refining the identifica-641 tion model through inclusion of additional datasets, updating known FTP 642 locations, and feedback from on-the-ground data collection efforts. Though 643 this work is primarily focused on leveraging the relationship between exter-644 nal datasets and FTP, the role of *nearby* FTP within a known touch point 645 dataset could potentially have an impact on the identification of new FTP as 646 well. Additionally, we are in the midst of assessing the accuracy of our exist-647 ing model and refining new models based on data from neighboring countries 648 in the region. Further examination of neighboring country-specific datasets 649 will lead to a better understanding of the impact that socio-economics, de-650 mographics, ICT adoption, etc. have on the ability to successfully identify 651 FTP locations at a broader scale. 652

## 653 Acknowledgements

This research was made possible in part through a grant from the Bill & Melinda Gates Foundation, Grand Challenges Explorations initiative (Grant number OPP1140466). We would also like to acknowledge the mobile application development work by the software and data science team at Spatial Development International and data collection efforts by our partners at GeoNAREM. Additional thanks go to Brand Fusion and the Humanitarian OpenStreetMap Team.

## 661 References

- [1] European Investment Bank, Banking in sub-Saharan Africa:Recent
   Trends and Digital Financial Inclusion, Technical Report, European In vestment Bank, 2016.
- [2] T. Triki, I. Faye, Financial Inclusion in Africa, Technical Report, African
   Development Bank, 2013.

- [3] C. A. of Kenya, Quarterly Sector Statistics Report, Technical Report,
   Communications Authority of Kenya, 2016.
- [4] World Bank, The Global Findex Database 2014: Measuring Financial
   Inclusion around the World. Policy Research Working Paper 7255, Technical Report, World Bank, 2015.
- 672 [5] M. Ochieng, The new money lenders of nairobi,
  673 https://www.fsdafrica.org/knowledge-hub/blog/
  674 the-new-money-lenders-of-nairobi/, 2016.
- [6] Economist Intelligence Unit, M-pesa: Out of africa, into romania, The Economist (2014).
- [7] T. Suri, W. Jack, The long-run poverty and gender impacts of mobile money, Science 354 (2016) 1288–1292.
- [8] H. Sekabira, M. Qaim, Can mobile phones improve gender equality and
  nutrition? panel data evidence from farm households in uganda, Food
  Policy 73 (2017) 95–103.
- [9] T. Suri, Mobile money, Annual Review of Economics 9 (2017) 497–520.
- [10] A. Olingo, Kenya, uganda in plans to pull informal sector into tax
   bracket, The East African (2016).
- [11] N. Ndung'u, Understanding and expanding financial inclusion in kenya,
   2014. Keynote Speech at FinAccess GIS Mapping of all Financial Access
   Touch Points.
- [12] Brand Fusion, Financial Inclusion Research Project Handbook, Tech nical Report, 2015.
- [13] FSD Kenya, FinAccess geospatial mapping 2013, Technical Report, FSD
   Kenya, 2015.
- [14] J. Kim, Reaching the rural regions in kenya through mobile money, http://finclusionlab.org/es/node/519/, 2016.
- [15] O. K. Kirui, J. J. Okello, R. A. Nyikal, G. W. Njiraini, et al., Impact
  of mobile phone-based money transfer services in agriculture: evidence
  from kenya, Quarterly Journal of International Agriculture 52 (2013)
  141–162.

- <sup>698</sup> [16] N. Hughes, S. Lonie, M-pesa: mobile money for the "unbanked" turning <sup>699</sup> cellphones into 24-hour tellers in kenya, Innovations 2 (2007) 63–81.
- <sup>700</sup> [17] P. Uithol, Mapping financial inclusion in uganda (2015).
- [18] Brand Fusion, Kenya Multi Sector GIS Mapping Project Final Report,
   Technical Report, Brand Fusion, 2015.
- [19] A. Xylouris, Connected Consumer Survey 2016: mobile services and
   devices in Sub-Saharan Africa, Technical Report, Analysys Mason, 2016.
- [20] Facebook People Insights, Journeys of connectivity: How people in sub saharan africa come online, Facebook IQ (2017).
- T. Shapshak, Facebook has 170 million african users, mostly on mobile,
   Forbes (2017).
- [22] B. A. of Kenya, State of Internet in Kenya 2016, Technical Report,
   Bloggers Association of Kenya, 2016.
- [23] M. Kaigwa, O. Madung, S. Costello, Nendo 2014/15 Social Media Trend
   Report, Technical Report, Nendo Consultancy, 2015.
- [24] G. McKenzie, K. Janowicz, Coerced geographic information: The notso-voluntary side of user-generated geo-content, in: Eighth international
  conference on geographic information science.
- <sup>716</sup> [25] M. F. Goodchild, Citizens as sensors: the world of volunteered geogra-<sup>717</sup> phy, GeoJournal 69 (2007) 211–221.
- [26] G. Porter, Mobile phones, livelihoods and the poor in sub-saharan africa:
   Review and prospect, Geography Compass 6 (2012) 241–259.
- [27] A. Tanle, A. M. Abane, Mobile phone use and livelihoods: qualitative
  evidence from some rural and urban areas in ghana, GeoJournal (2017)
  1-11.
- <sup>723</sup> [28] B. Dillon, Using mobile phones to collect panel data in developing <sup>724</sup> countries, Journal of international development 24 (2012) 518–527.
- [29] S. A. Asongu, J. C. Nwachukwu, The role of governance in mobile phones
  for inclusive human development in sub-saharan africa, Technovation 55
  (2016) 1–13.

- <sup>728</sup> [30] J. Guo, H. Bouwman, An ecosystem view on third party mobile payment <sup>729</sup> providers: a case study of alipay wallet, info 18 (2016) 56–78.
- [31] P. Armstrong, Y. Wang, Is alibaba losing to tencent in china's trilliondollar payment war?, Forbes (2018).
- [32] P. N. Howard, M. R. Parks, Social media and political change: Capacity,
  constraint, and consequence, Journal of communication 62 (2012) 359–
  362.
- [33] A. Wesolowski, C. O. Buckee, L. Bengtsson, E. Wetter, X. Lu, A. J.
  Tatem, Commentary: containing the ebola outbreak-the potential and challenge of mobile network data, PLoS currents 6 (2014).
- [34] A. Stefanidis, A. Crooks, J. Radzikowski, Harvesting ambient geospatial
   information from social media feeds, GeoJournal 78 (2013) 319–338.
- [35] O. Okolloh, Ushahidi, or testimony: Web 2.0 tools for crowdsourcing
   crisis information, Participatory learning and action 59 (2009) 65–70.
- [36] L. Palen, R. Soden, T. J. Anderson, M. Barrenechea, Success & scale in a data-producing organization: The socio-technical evolution of openstreetmap in response to humanitarian events, in: Proceedings of the 33rd annual ACM conference on human factors in computing systems, ACM, pp. 4113–4122.
- [37] B. Adams, G. McKenzie, M. Gahegan, Frankenplace: interactive thematic mapping for ad hoc exploratory search, in: Proceedings of the
  24th International Conference on World Wide Web, International World
  Wide Web Conferences Steering Committee, pp. 12–22.
- [38] C. Linard, A. Tatem, F. R. Stevens, A. Gaughan, N. N. Patel, Z. Huang,
  Use of active and passive vgi data for population distribution modelling:
  experience from the worldpop project (2014) 1–16.
- [39] G. McKenzie, K. Janowicz, Where is also about time: A locationdistortion model to improve reverse geocoding using behavior-driven
  temporal semantic signatures, Computers, Environment and Urban Systems 54 (2015) 1–13.

- [40] Y. Hu, K. Janowicz, H. Couclelis, Prioritizing disaster mapping tasks
  for online volunteers based on information value theory, Geographical
  Analysis 49 (2017) 175–198.
- [41] L. Anselin, Spatial econometrics, A companion to theoretical econo metrics 310330 (2001).
- [42] M. J. Cracknell, A. M. Reading, Geological mapping using remote sensing data: A comparison of five machine learning algorithms, their response to variations in the spatial distribution of training data and the
  use of explicit spatial information, Computers & Geosciences 63 (2014)
  22–33.
- [43] C. Song, M.-P. Kwan, W. Song, J. Zhu, A comparison between spatial econometric models and random forest for modeling fire occurrence, Sustainability 9 (2017) 819.
- [44] F. R. Stevens, A. E. Gaughan, C. Linard, A. J. Tatem, Disaggregating
  census data for population mapping using random forests with remotelysensed and ancillary data, PloS one 10 (2015) e0107042.
- [45] O. Chesnokova, M. Nowak, R. S. Purves, A crowdsourced model of
  landscape preference, in: LIPIcs-Leibniz International Proceedings in
  Informatics, volume 86, Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.
- [46] P. A. Moran, Notes on continuous stochastic phenomena, Biometrika
  37 (1950) 17–23.
- [47] B. D. Ripley, The second-order analysis of stationary point processes,
   Journal of applied probability 13 (1976) 255–266.
- [48] C. Linard, M. Gilbert, R. W. Snow, A. M. Noor, A. J. Tatem, Population
  distribution, settlement patterns and accessibility across africa in 2010,
  PloS one 7 (2012) e31743.
- [49] M. A. Friedl, D. K. McIver, J. C. Hodges, X. Zhang, D. Muchoney,
  A. H. Strahler, C. E. Woodcock, S. Gopal, A. Schneider, A. Cooper,
  et al., Global land cover mapping from modis: algorithms and early
  results, Remote Sensing of Environment 83 (2002) 287–302.

- [50] N. W. Ruktanonchai, P. DeLeenheer, A. J. Tatem, V. A. Alegana, T. T.
  Caughlin, E. zu Erbach-Schoenberg, C. Lourenço, C. W. Ruktanonchai,
  D. L. Smith, Identifying malaria transmission foci for elimination using
  human mobility data, PLoS computational biology 12 (2016) e1004846.
- [51] P. Deville, C. Linard, S. Martin, M. Gilbert, F. R. Stevens, A. E.
  Gaughan, V. D. Blondel, A. J. Tatem, Dynamic population mapping using mobile phone data, Proceedings of the National Academy of Sciences 111 (2014) 15888–15893.
- [52] D. Balk, U. Deichmann, G. Yetman, F. Pozzi, S. Hay, A. Nelson, Deter mining global population distribution: methods, applications and data,
   Advances in parasitology 62 (2006) 119–156.
- [53] S. Freire, T. Kemper, M. Pesaresi, A. Florczyk, V. Syrris, Combining GHSL and GPW to improve global population mapping, in: Geoscience and Remote Sensing Symposium (IGARSS), 2015 IEEE International, IEEE, pp. 2541–2543.
- <sup>804</sup> [54] P. D. Broxton, X. Zeng, D. Sulla-Menashe, P. A. Troch, A global land
  <sup>805</sup> cover climatology using modis data, Journal of Applied Meteorology
  <sup>806</sup> and Climatology 53 (2014) 1593–1605.
- <sup>807</sup> [55] H. Rahemtulla, J. Kaplan, B.-S. Gigler, S. Cluster, J. Kiess, C. Brigham,
  <sup>808</sup> Open Data Kenya: Case Study of the Underlying Drivers, Principal
  <sup>809</sup> Objectives and Evolution of One of the First Open Data Initiatives in
  <sup>810</sup> Africa, Open Development Technology Alliance (ODTA), 2012.
- <sup>811</sup> [56] J. E. Dobson, E. A. Bright, P. R. Coleman, R. C. Durfee, B. A. Worley,
  <sup>812</sup> Landscan: a global population database for estimating populations at
  <sup>813</sup> risk, Photogrammetric engineering and remote sensing 66 (2000) 849–
  <sup>814</sup> 857.
- [57] Natural Earth, Populated places, http://www.naturalearthdata.
   com/downloads/10m-cultural-vectors/10m-populated-places/,
   2014.
- <sup>818</sup> [58] S. Perez, Foursquare begins crowdsourcing local business data collection
  <sup>819</sup> with questions that appear after check-ins (2013).

- [59] P. O. Gislason, J. A. Benediktsson, J. R. Sveinsson, Random forests for
  land cover classification, Pattern Recognition Letters 27 (2006) 294–300.
- [60] C. M. Jarque, A. K. Bera, Efficient tests for normality, homoscedasticity
  and serial independence of regression residuals, Economics letters 6
  (1980) 255–259.
- [61] T. S. Breusch, A. R. Pagan, A simple test for heteroscedasticity and
  random coefficient variation, Econometrica: Journal of the Econometric
  Society (1979) 1287–1294.
- [62] L. Anselin, Spatial econometrics: methods and models, volume 4,
   Springer Science & Business Media, 2013.
- [63] H. Drucker, C. J. Burges, L. Kaufman, A. J. Smola, V. Vapnik, Sup port vector regression machines, in: Advances in neural information
   processing systems, pp. 155–161.
- [64] C. Cortes, V. Vapnik, Support-vector networks, Machine learning 20 (1995) 273–297.
- [65] T. K. Ho, Random decision forests, in: Document Analysis and Recognition, 1995., Proceedings of the Third International Conference on,
  volume 1, IEEE, pp. 278–282.
- [66] J. Friedman, T. Hastie, R. Tibshirani, The elements of statistical learn ing, volume 1, Springer series in statistics New York, 2001.