

# Mapping subnational gender gaps in internet and mobile adoption using social media data

Johns Hopkins University

Casey F. Breen<sup>1</sup> Masoomali Fatehkia<sup>2</sup> Jiani Yan<sup>1</sup>  
Douglas R. Leasure<sup>1</sup> Ingmar Weber<sup>3</sup> Ridhi Kashyap<sup>1</sup>

December 13, 2023

---

<sup>1</sup>University of Oxford

<sup>2</sup>Qatar Computing Research Institute

<sup>3</sup>Saarland University

# Roadmap

1. Overview of **digital gender gaps** project
2. Our approach to using **social media data** to predict subnational digital gender gaps
3. Overview of subnational estimates

# Benefits of digital revolution

- ▶ The **digital revolution** has ushered in tremendous societal and economic benefits
  - ▶ Lower gender inequality, lower maternal/child mortality, higher contraception (Rotondi et al., 2020)
  - ▶ Boost social connectivity, social learning, access to vital services (Unwin, 2009; DiMaggio and Hargittai, 2001; Suri and Jack, 2016)
  - ▶ Increases levels of education, economic benefits (Hjort and Poulsen, 2019; Kho, Lakdawala and Nakasone, 2018; Kharisma, 2022)
- ▶ Benefits are often greatest in the most unequal, disadvantaged areas

# Tracking the digital divide

- ▶ Access to digital technologies such as mobile phones and internet remains **highly unequal**
  - ▶ Especially in low- and middle-income countries
  - ▶ Especially among women
- ▶ **UN Sustainable Development Goals (SDGs)**: Reducing inequalities in access to digital technologies by gender (SDG5) and reducing digital literacy gaps (SDG4)

# Digital gender gaps project overview

1. **Data infrastructure**: Map and understand gender gaps in digital connectivity and social media use
  - ▶ **Today - subnational estimates**
2. **Impacts research**: impacts of digital information and capabilities on women's economic and social empowerment outcomes
  - ▶ Cross-national, comparative perspective (low- and middle- income countries)

# Original “impacts” research

## Using Facebook ad data to track the global digital gender gap

Masoomali Fatehkia<sup>a</sup>, Ridhi Kashyap<sup>b</sup>  , Ingmar Weber<sup>c</sup>

Show more 

+ Add to Mendeley  Share  Cite

Regular article | [Open access](#) | [Published: 29 July 2021](#)

## Analysing global professional gender gaps using LinkedIn advertising data

Ridhi Kashyap  & Florianne C. J. Verkroost

[EPJ Data Science](#) **10**, Article number: 39 (2021) | [Cite this article](#)

**6421** Accesses | **13** Citations | **21** Altmetric | [Metrics](#)

RESEARCH ARTICLE | SOCIAL SCIENCES | 



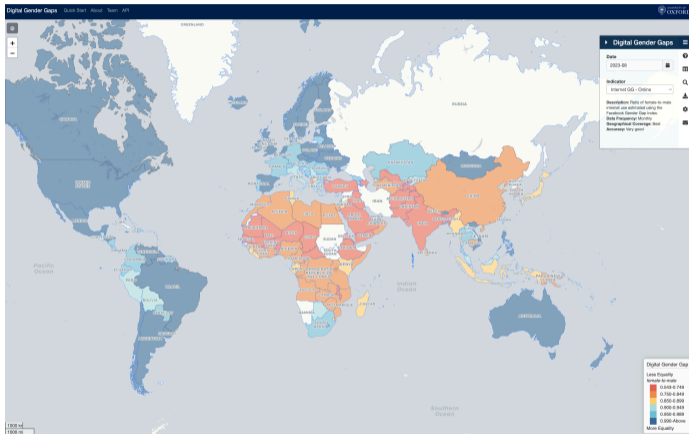
## Leveraging mobile phones to attain sustainable development

Valentina Rotondi  , Ridhi Kashyap , Luca Maria Pesando , , and Francesco C. Billari  [Authors info & Affiliations](#)

Edited by Barbara Entwisle, University of North Carolina at Chapel Hill, Chapel Hill, NC, and accepted by Editorial Board Member Mary C. Waters April 6, 2020 (received for review May 30, 2019)

June 1, 2020 | 117 (24) 13413-13420 | <https://doi.org/10.1073/pnas.1909326117>

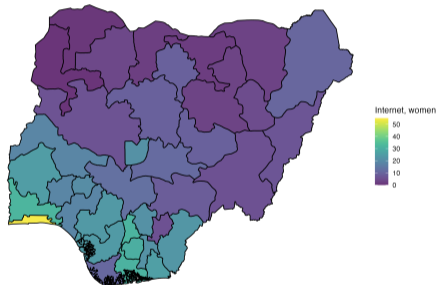
# Data infrastructure – digitalgendergaps.org



(Kashyap et al., 2020)

# Adoption of digital technology varies geographically

Observed

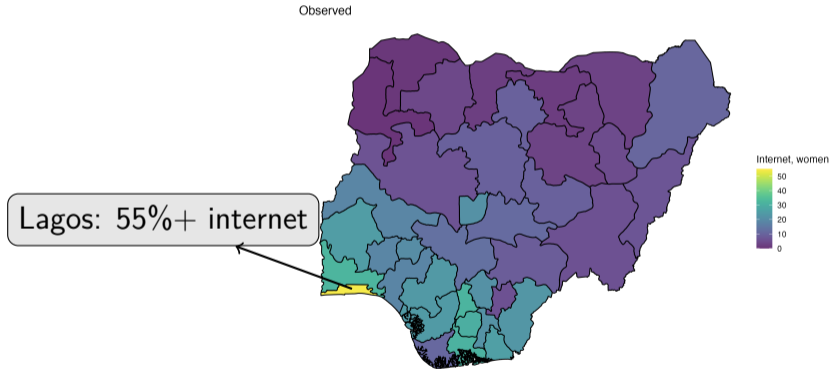


Source: Nigeria, Demographic and Health Survey



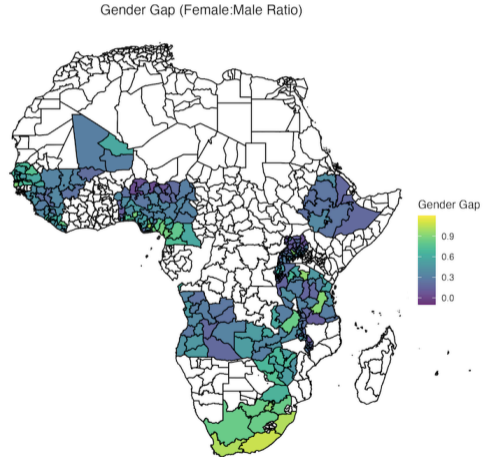


# Women using internet, past 12 months



# Develop subnational estimates of adoption

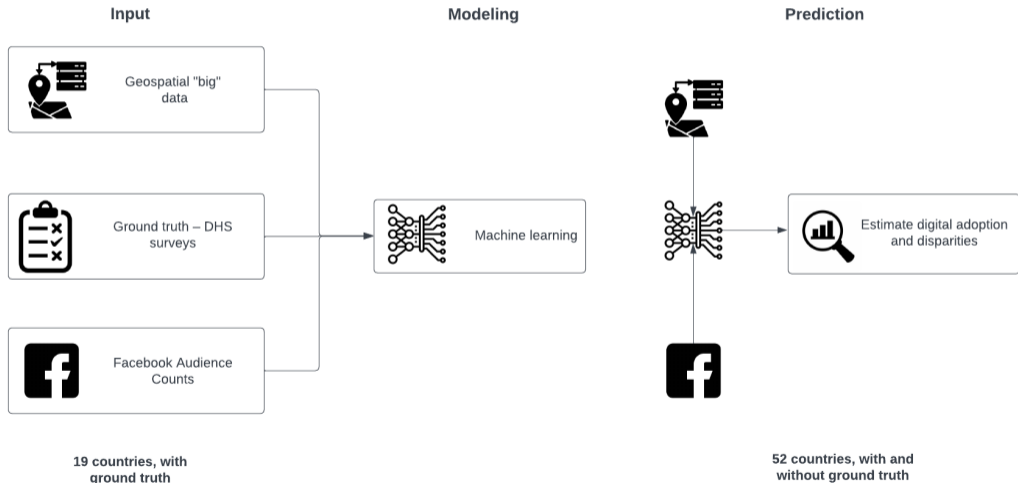
- ▶ **Goal:** Develop estimates of internet and mobile adoption by gender and digital gender gaps
- ▶ First GADM1 subnational level
  - ▶  $N = 874$



# Prediction framework - theoretical background

- ▶ Digital gender gaps will be shaped by overall levels of economic development and digital infrastructure
- ▶ **Patriarchal** norms and beliefs will moderate this relationship

# Overview of approach

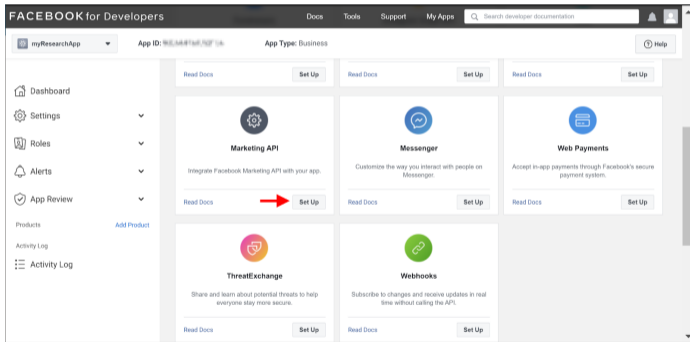


# Ground truth – Demographic and Health Surveys (DHS)

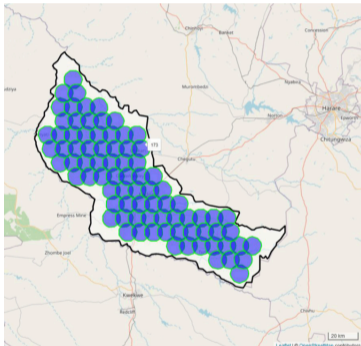
- ▶ Household surveys representative at the first subnational level
  - ▶ Standardized sample design, questionnaire, implementation, etc.
  - ▶ Questions on individual-level internet use and mobile phone use (wave 7 onwards)
- ▶ Focus on 19 different DHS surveys, 2016-2020

# Facebook audience counts

- ▶ Collected through public marketing API
- ▶ Specify geographic region (FB template or custom region)
- ▶ Disaggregated counts by gender, age, device type, etc.



# Facebook audience counts 'online predictors'



## ► Collected in 2021:

1. Facebook penetration 13+ female
2. Facebook penetration 13+ male
3. Facebook audience 13+ gender gap
4. iOS 13+ female fraction
5. iOS 13+ male fraction
6. WiFi age 13+ female fraction
7. WiFi age 13+ male fraction
8. 4G network age 13+ female fraction
9. 4G network age 13+ male fraction
10. FB rural WiFi mean (pop weighted)



# Geospatial and population data

- ▶ Include 'offline' predictors that are uniformly available and consistent across subnational units
  - ▶ Satellite-derived nightlights data
  - ▶ Population density
  - ▶ Subnational education index, income index, human development index (HDI), gender development index (GDI)

## Full set of offline predictors

| Variable Name                   | Source                    | Country (N) |
|---------------------------------|---------------------------|-------------|
| Educational Index Females       | Subnational Dev. Database | 50          |
| Educational Index Males         | Subnational Dev. Database | 50          |
| Income Index Females            | Subnational Dev. Database | 50          |
| Income Index Males              | Subnational Dev. Database | 50          |
| Subnational GDI                 | Subnational Dev. Database | 50          |
| Subnational HDI Females         | Subnational Dev. Database | 50          |
| Subnational HDI Males           | Subnational Dev. Database | 50          |
| WPop 2020 Age 15-49 Female Frac | WorldPop                  | 58          |
| WPop 2020 Age 15-49 Male Frac   | WorldPop                  | 58          |
| WPop 2020 Pop Density           | WorldPop                  | 59          |
| Nightlights Mean Pop Weighted   | Earth Observation Group   | 58          |

## Outcomes of interest (from DHS)

| Indicators               | Women | Men | Gender Gap |
|--------------------------|-------|-----|------------|
| Mobile Phone Ownership   | ✓     | ✓   | ✓          |
| Internet Use, Past 12 Mo | ✓     | ✓   | ✓          |

# Defining a Digital Gender Gap

$$\text{Gender Gap} = \frac{\text{Indicator}_f / \text{Indicator}_m}{\text{Pop}_f / \text{Pop}_m} \quad (1)$$

where

- ▶  $\text{Indicator}_f$  is the number of female (male) users aged 15–49 (e.g., internet, past 12 months)
- ▶  $\text{Pop}_f$  is the total population of women (men) aged 15–49

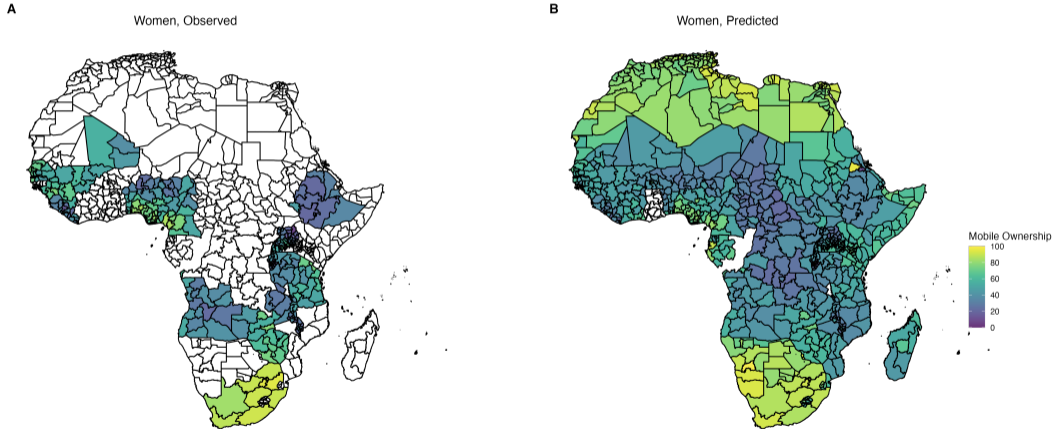
# Machine Learning Strategy

- ▶ How do you pick the **best** machine learning algorithm?
- ▶ Fit lots of algorithms, see which have the best performance
- ▶ Ensemble learning to combine algorithms and tests performance using cross-validation to estimate mean squared error for each algorithm (Van der Laan, Polley and Hubbard, 2007)

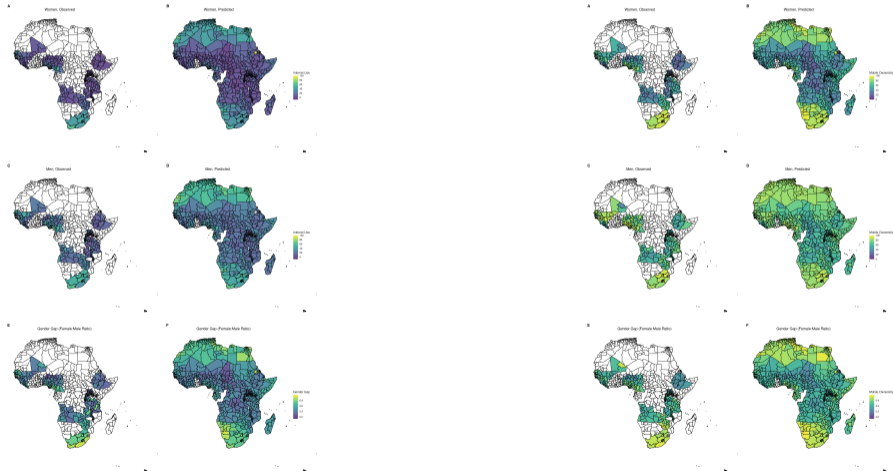
# Machine Learning Algorithms Considered

| Algorithm            | Description  |
|----------------------|--|
| glmnet (Lasso)       | Lasso Regression                                       |
| glmnet (Ridge)       | Ridge Regression                                       |
| glmnet (Elastic Net) | Elastic Net with 50% L1 Ratio                          |
| polspline            | Polynomial Spline                                      |
| ranger               | Random Forest with 100 Trees                           |
| gbm                  | Gradient Boosted Machine                               |
| glm                  | Generalized Linear Model                               |
| xgboost              | Extreme Gradient Boosting                              |
| SuperLearner         | Ensemble method combining multiple learning algorithms |

# Greatly expanded coverage of digital technology adoption



# Similar overall patterns for internet and mobile

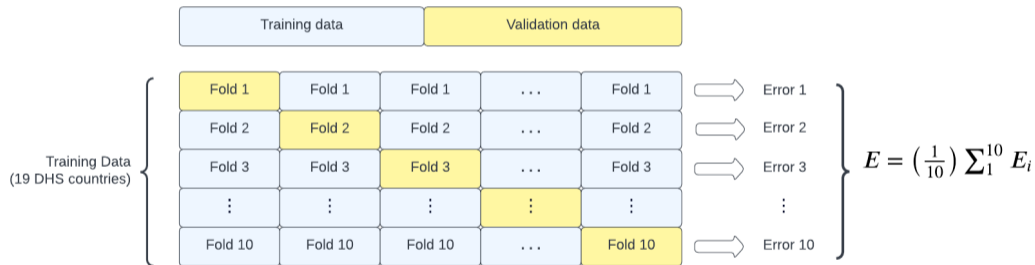




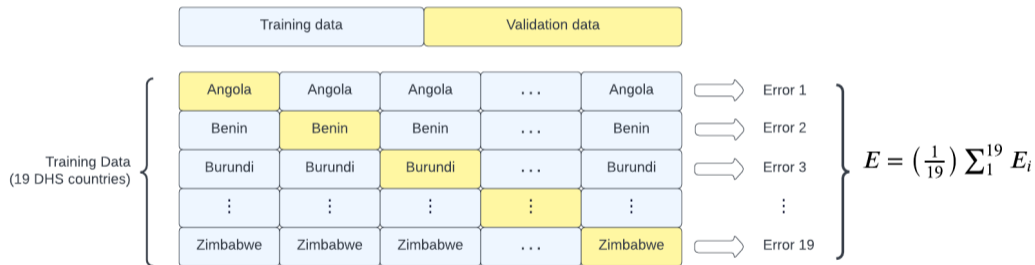
# Testing model performance

- ▶ How do we assess model performance?
- ▶ **Cross-validation** using 19 countries with ground truth data

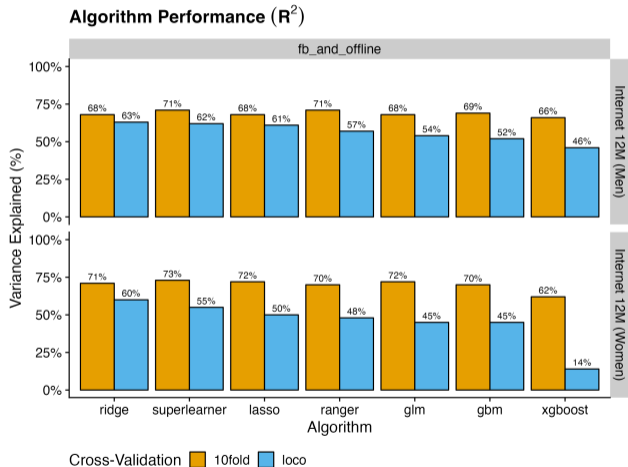
# 10-fold cross validation



# Leave-one-country-out cross validation

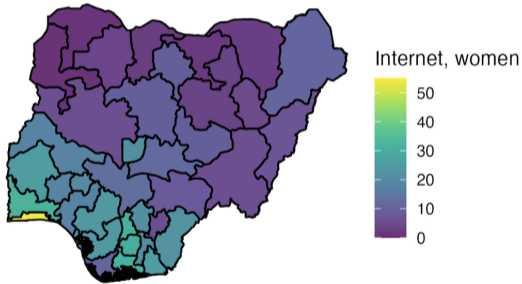


# Model performance

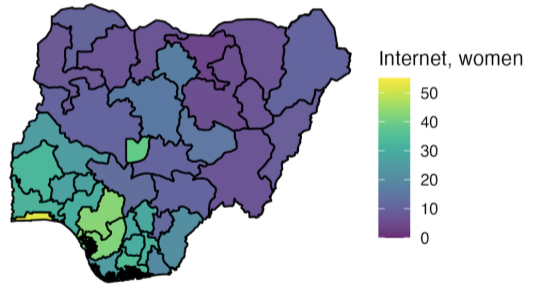


# Results for Nigeria (Leave-one-country-out)

Observed



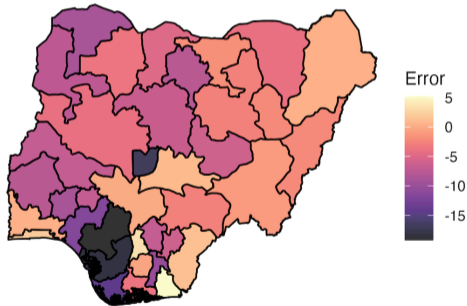
Predicted



# Assessing predictive accuracy

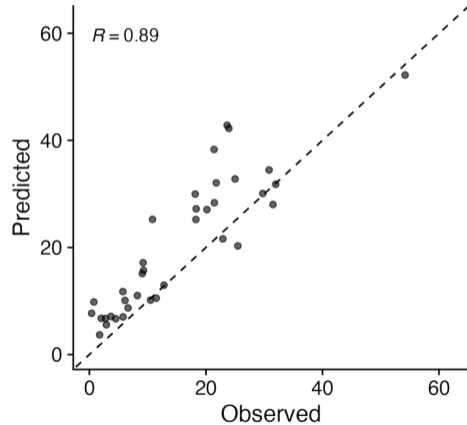
**C**

Error

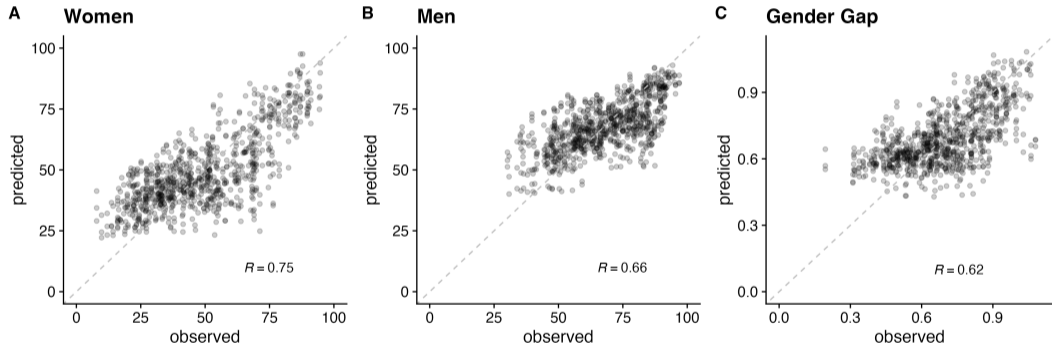


**D**

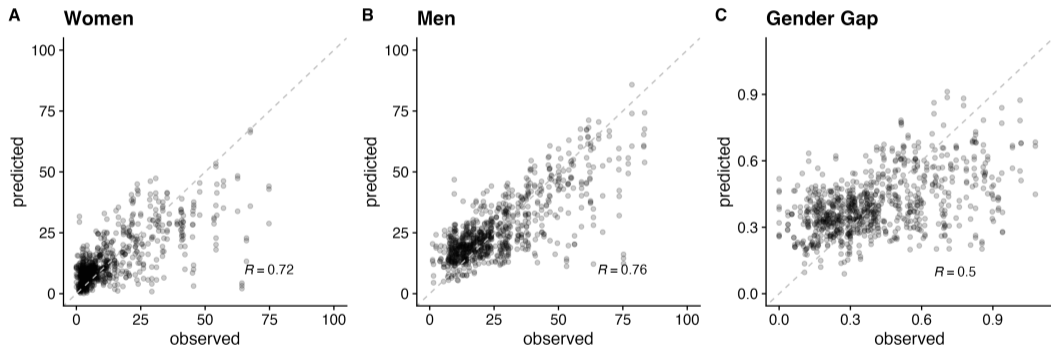
**Internet, women (error)**



# Overall predictiveness – mobile

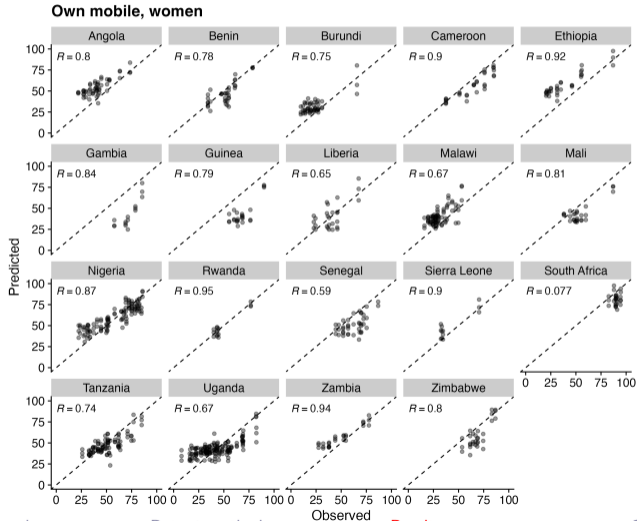


# Overall predictiveness – internet

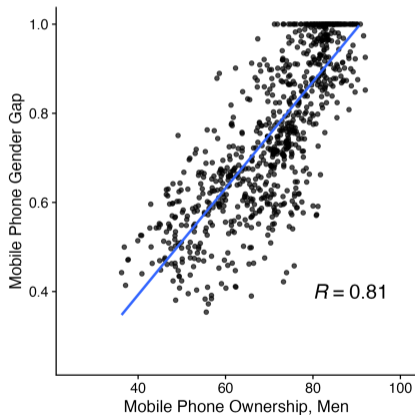




# Large variation in predictive accuracy across countries



# Relationship: levels of mobile phone penetration and gender gaps



## Next steps and future opportunities

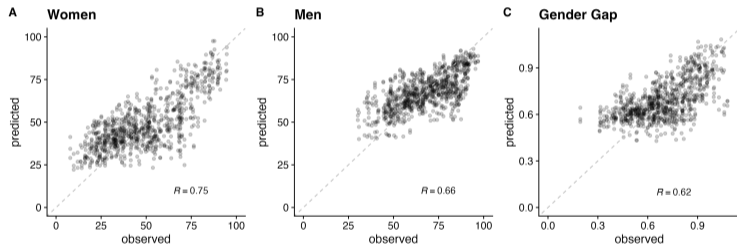
- ▶ Regular Facebook collections and pipeline to monitor trends over time
- ▶ Residual analysis + quantifying uncertainty: what factors explain where model does worse?

# Summary

- ▶ Using Facebook audience counts **greatly expands** our ability to accurately predict digital gender gaps in countries with no ground truth
- ▶ Huge **disparities** in access to mobile and internet technologies between and within countries
- ▶ New opportunities to study **population-level impacts** of digital technology using these subnational estimates

# Thank You

## ► Questions?



 caseyfbreen

 casey.breen@demography.ox.ac.uk

# References

- DiMaggio, Paul and Eszter Hargittai. 2001. "From the 'Digital Divide' to 'Digital Inequality': Studying Internet Use as Penetration Increases." p. 25.
- Hjort, Jonas and Jonas Poulsen. 2019. "The Arrival of Fast Internet and Employment in Africa." *American Economic Review* 109(3):1032–1079.
- Kashyap, Ridhi, Masoomali Fatehikia, Reham Al Tamime and Ingmar Weber. 2020. "Monitoring Global Digital Gender Inequality Using the Online Populations of Facebook and Google." *Demographic Research* 43:779–816.
- Kharisma, Bayu. 2022. "Surfing Alone? The Internet and Social Capital: Evidence from Indonesia." *Journal of Economic Structures* 11(1):8.
- Kho, Kevin, Leah K Lakdawala and Eduardo Nakasone. 2018. "Impact of Internet Access on Student Learning in Peruvian Schools."
- Rotondi, Valentina, Ridhi Kashyap, Luca Maria Pesando, Simone Spinelli and Francesco C. Billari. 2020. "Leveraging Mobile Phones to Attain Sustainable Development." *Proceedings of the National Academy of Sciences* 117(24):13413–13420.
- Suri, Tavneet and William Jack. 2016. "The Long-Run Poverty and Gender Impacts of Mobile Money." *Science* 354(6317):1288–1292.
- Unwin, P. T. H. 2009. *ICT4D: Information and Communication Technology for Development*. Cambridge University Press.
- Van der Laan, Mark J., Eric C. Polley and Alan E. Hubbard. 2007. "Super Learner." *Statistical Applications in Genetics and Molecular Biology* 6(1).