

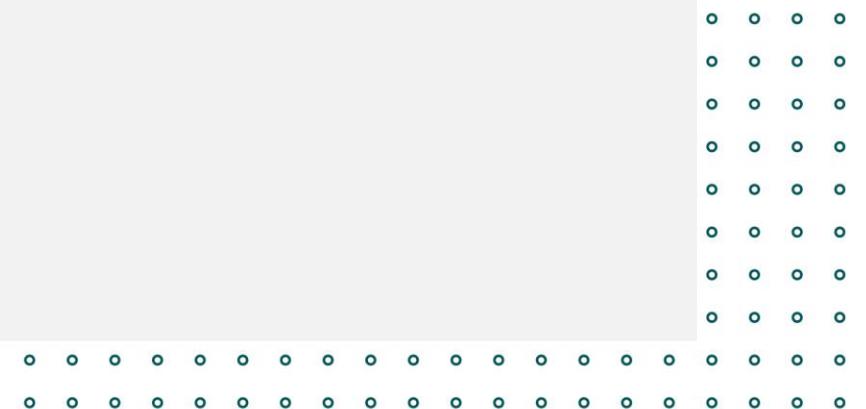


MEASURING AND UNDERSTANDING BIAS IN DIGITAL TRACE DATA

Emilio Zagheni & Tom Theile

Max Planck Institute for Demographic Research

SICSS Saarbruecken – Sept 10, 2025







ADVERTISING PLATFORMS LIKE THOSE OF FACEBOOK, LINKEDIN, ETC.

Locations

People living in this location

United States

New York, New York + 25mi

Include ▾ Search Locations Browse

Drop Pin

Add Locations in Bulk

Age
18 - 65+

Gender
All genders

Detailed Targeting
Include people who match ⓘ

Behaviors > Expats
Lived in UK (Formerly Expats - UK)

Add demographics, interests or behaviors Suggestions Browse

Audience Definition

Your audience is defined.

Potential Reach: 23,000 people ⓘ

Estimated Daily Results

Reach ⓘ
1.0K - 2.9K

Link Clicks ⓘ
8 - 23

The accuracy of estimates is based on factors like past campaign data, the budget you entered, market data, targeting criteria and ad placements. Numbers are provided to give you an idea of performance for your budget, but are only estimates and don't guarantee results.

Were these estimates helpful?

How many users who used to live in the UK, now live in New York City?

www.facebook.com/business

FACEBOOK ADS AS A RECRUITMENT TOOL FOR SURVEY PARTICIPANTS



 **Health Behavior Survey** Sponsored ...

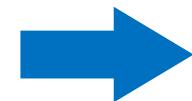
Do you live in the U.S.? We would like to learn about your health behavior!



© iStockphoto.com / Михаил Руденко

SURVEY3.GWDG.DE
We invite you to participate in our survey! [LEARN MORE](#)

2 Like Comment Share



MAX PLANCK INSTITUTE FOR DEMOGRAPHIC RESEARCH Load unfinished survey

0%

Health Behavior Survey United States

Thank you for your interest in our survey. This is an academic study led by researchers from the [Max Planck Institute for Demographic Research](#). Our goal is to understand health behaviors and help improve the health and well-being of people across countries at a time of increasing uncertainty. Your participation is crucial for our study.

The survey is directed at all people who are 18 years old or older, and it will take approximately 10-12 minutes to complete. Your participation is completely voluntary, you can stop participating at any time. In case you are not comfortable answering a particular question, you have the option to select "Prefer not to answer". Participants' data will be treated anonymously, and we will not ask for identifying information. If you wish to provide your email address, it will not be linked to your survey data. You can download our data protection policy [here](#).

If you have any questions about this research study, please contact Dr. André Grow and Dr. Daniela Perrotta at the Max Planck Institute for Demographic Research via healthsurvey@demogr.mpg.de.

Tick the box below and press 'Next' to continue.

I am willing to participate in this survey, I am at least 18 years old, and I have read the data protection policy.

[Next](#)

WARM UP

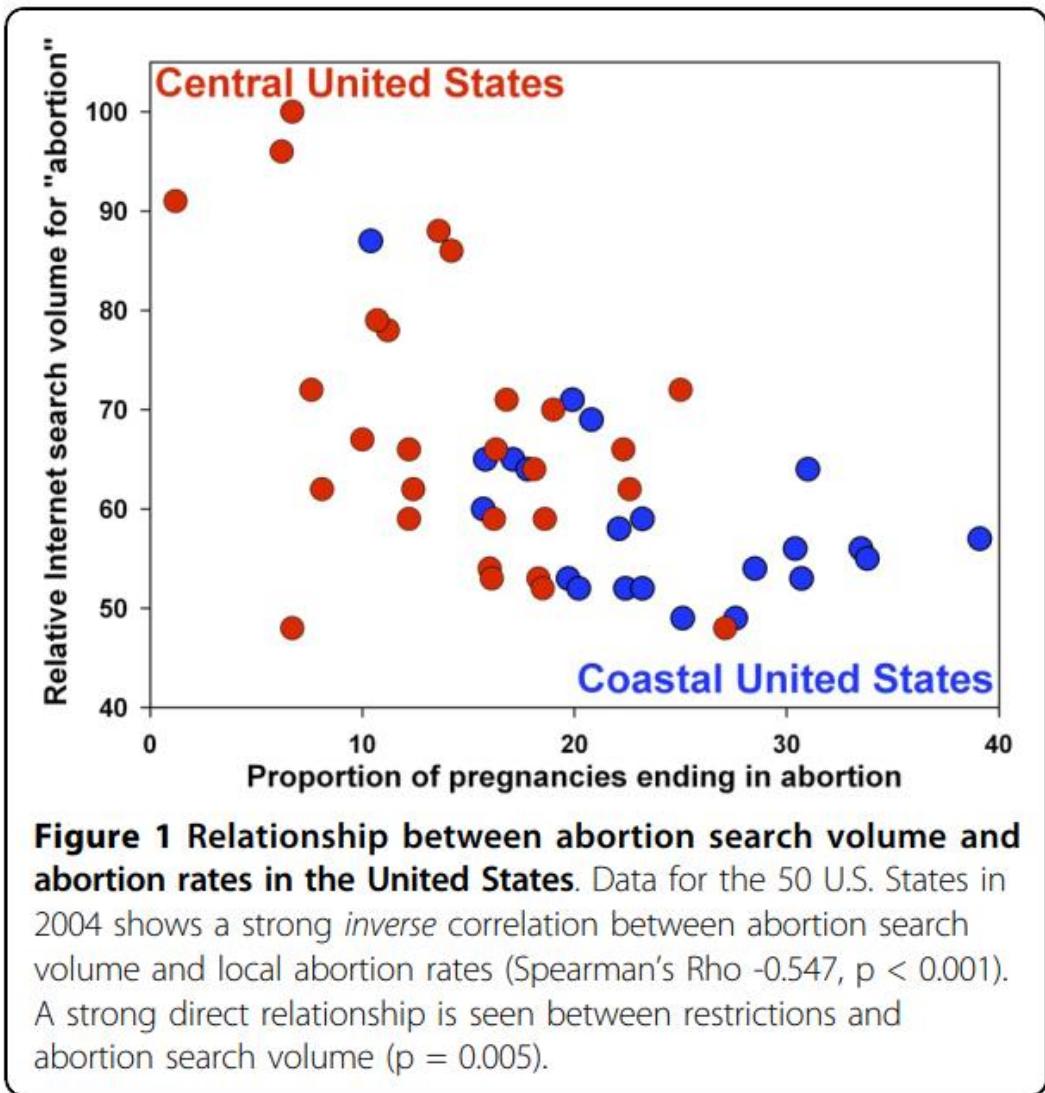


Figure 1 Relationship between abortion search volume and abortion rates in the United States. Data for the 50 U.S. States in 2004 shows a strong *inverse* correlation between abortion search volume and local abortion rates (Spearman's Rho -0.547, $p < 0.001$). A strong direct relationship is seen between restrictions and abortion search volume ($p = 0.005$).

Reis and Brownstein (2010)



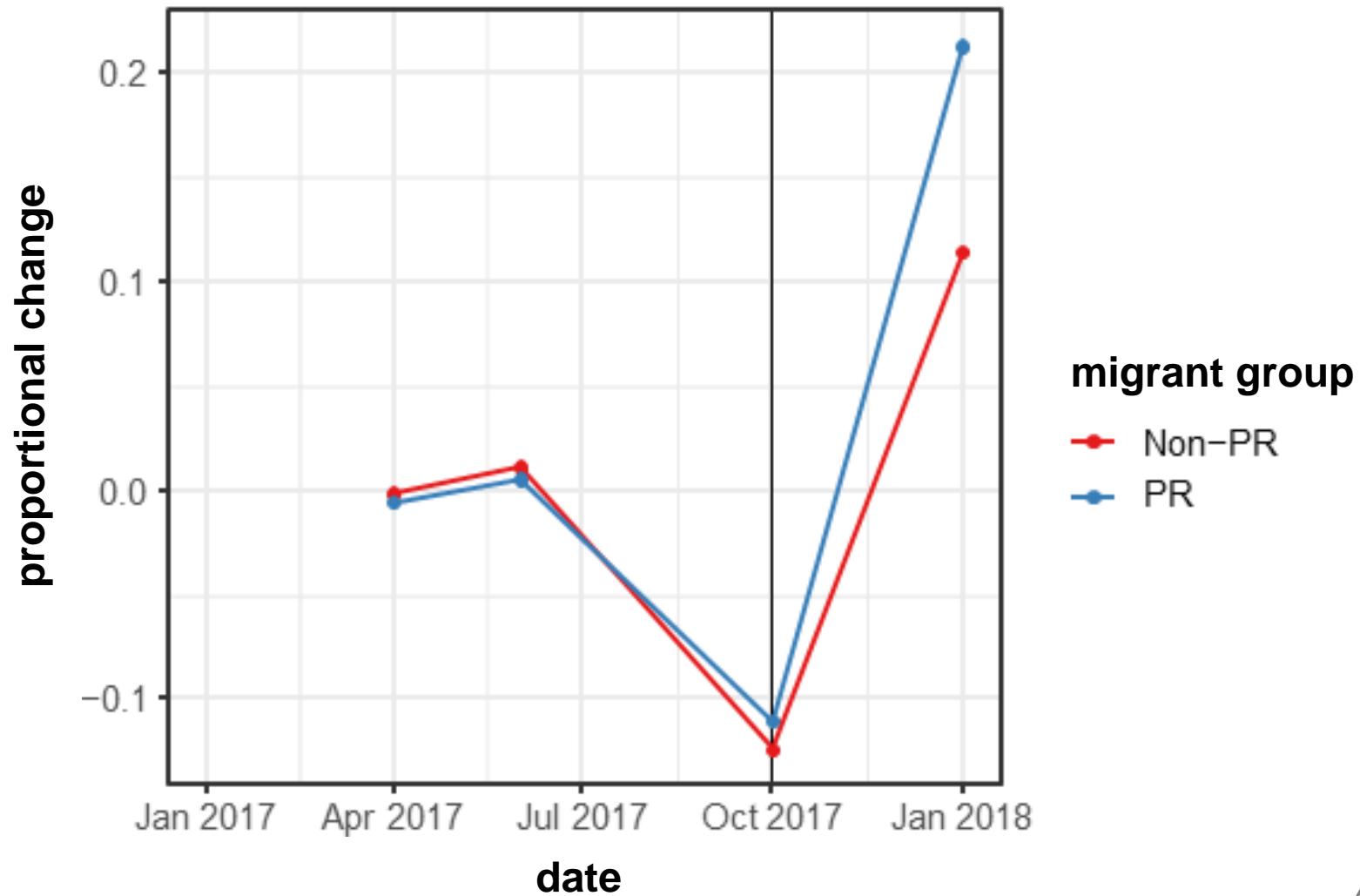
TODAY'S OUTLINE

- TRENDS OVER TIME AND FILTERING OUT BIAS
- SURVEYS AND NON-REPRESENTATIVE SAMPLES
- CALIBRATION AND WEIGHTING

MIGRATION AND MOBILITY AFTER HURRICANE MARIA MADE LANDFALL IN PUERTO RICO (SEPT. 2017)

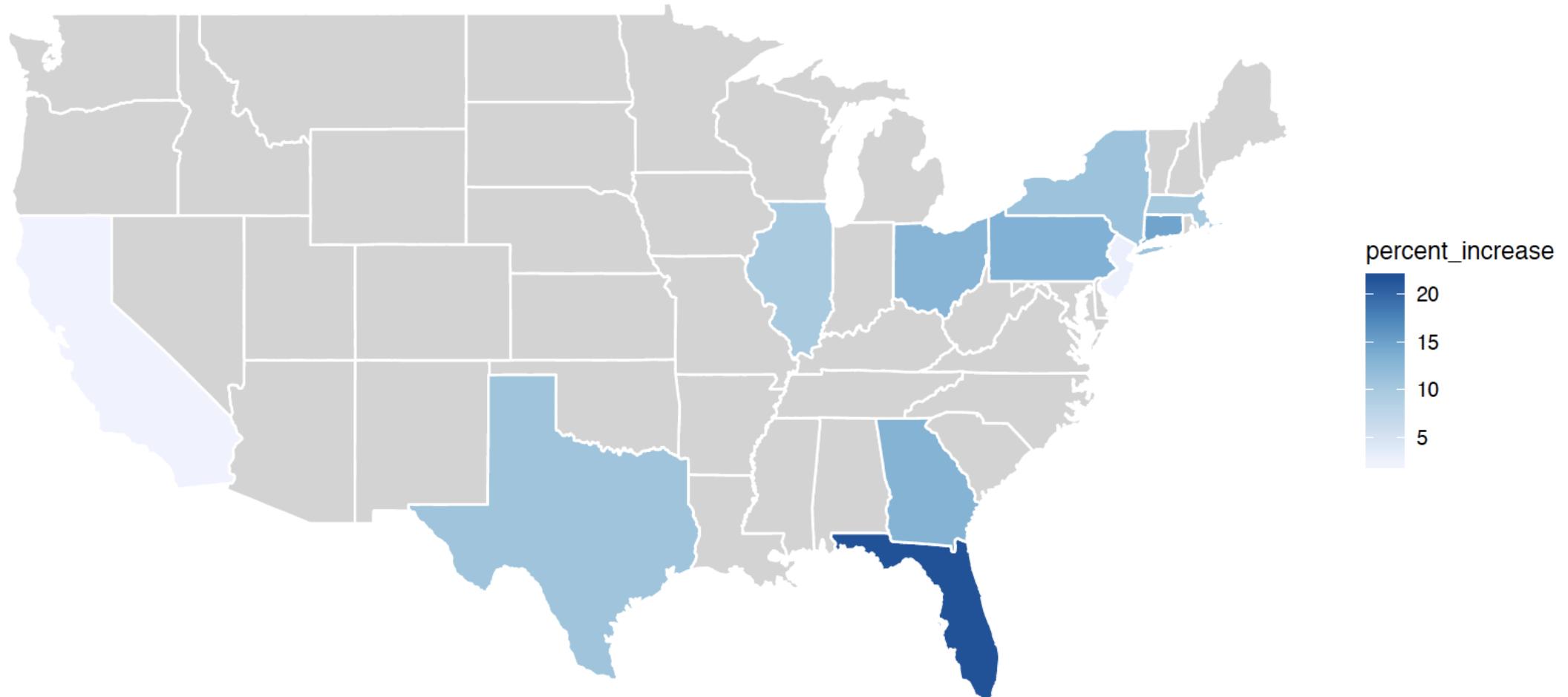


FACEBOOK DATA AND A DIFF-IN-DIFF APPROACH



Alexander, Polimis and Zagheni (2019)
Population and Development Review

PERCENT INCREASE IN PUERTO RICANS FROM OCT. 2017 TO JAN. 2018



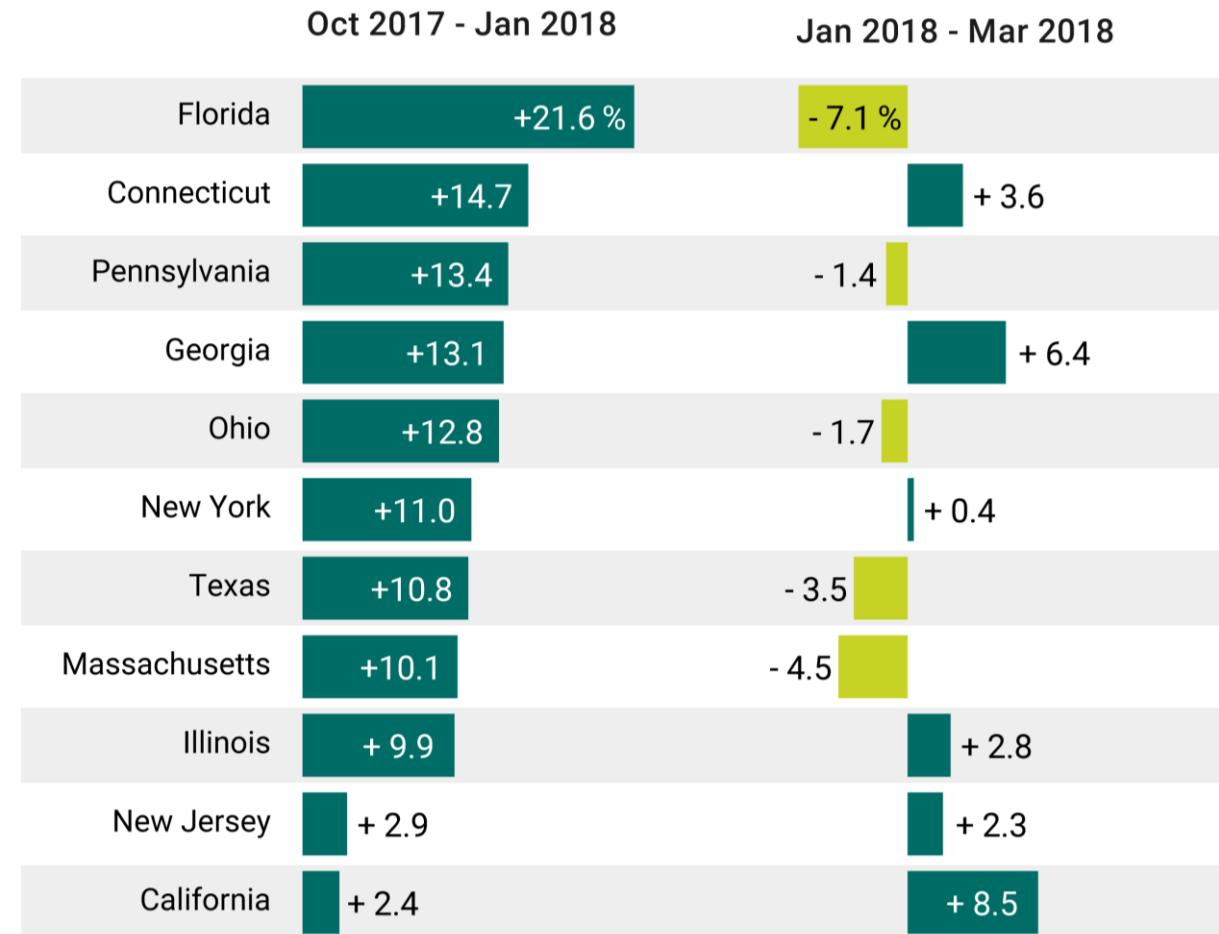
*Alexander, Polimis and Zagheni (2019)
Population and Development Review*



RETURN MIGRATION OR ONWARD RELOCATION

Migration after Hurricane Maria

Change of Puerto Rican migrant population in the US* (as percentages)



* For the nine US states with the largest populations of Puerto Rican migrants

Alexander, Polimis and Zagheni (2019)
Population and Development Review

CONSIDER THIS MODEL FOR DATA GENERATION



$$\underbrace{y_i^t}_{\text{Observation from social media for location } i} = \underbrace{n}_{\text{bias for location } i} + \underbrace{x_i^t}_{\text{"true" rate for location } i}$$

and

$$\underbrace{y_z^t}_{\text{Observation from social media for location } z} = \underbrace{m}_{\text{bias for location } z} + \underbrace{x_z^t}_{\text{"true" rate for location } z}$$



ASSUME THAT WE KNEW THE TRUE RATES (X) FOR TWO COUNTRIES: FRANCE (FR) AND SPAIN (SP)

$x_{FR}^{t+1} = 0.7$	$x_{SP}^{t+1} = 0.5$
$x_{FR}^t = 0.5$	$x_{SP}^t = 0.4$

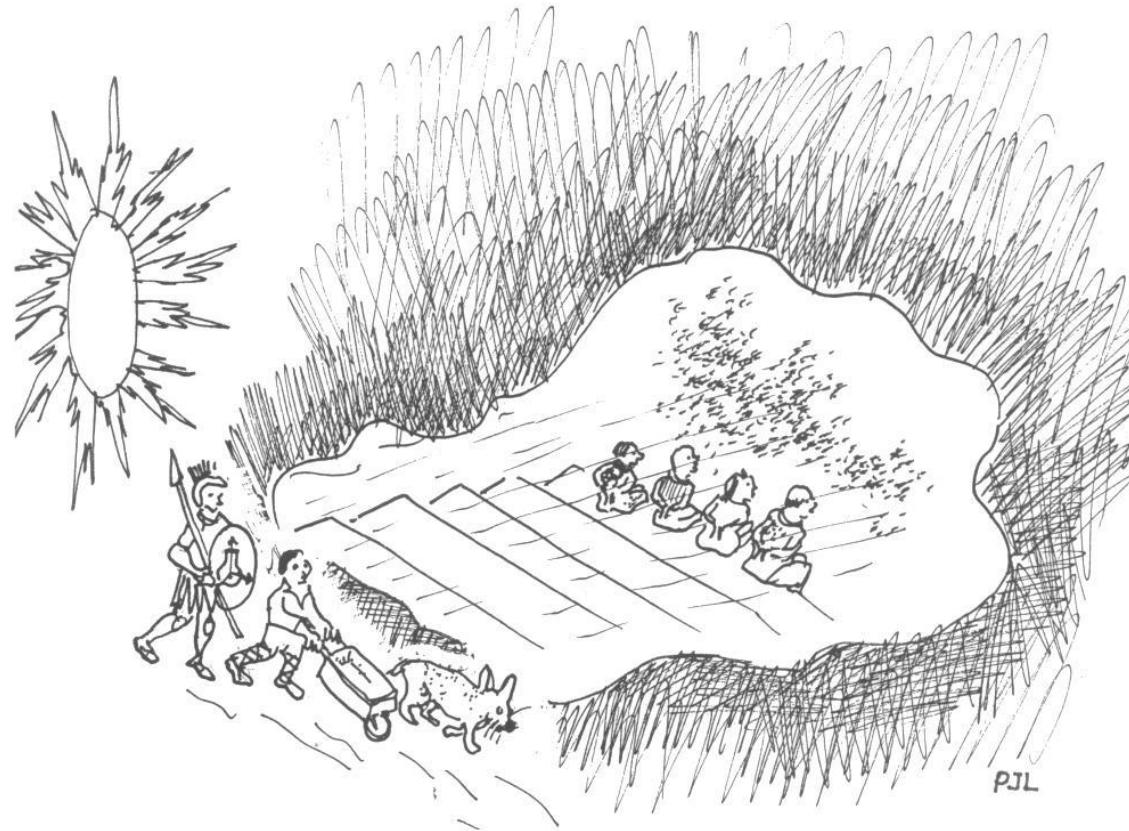
Let's define δ^{t+1} as the differential in the variation of these quantities of interest between time t and $(t + 1)$

$$\delta^{t+1} = \underbrace{(x_{FR}^{t+1} - x_{FR}^t) - (x_{SP}^{t+1} - x_{SP}^t)}_{\text{difference in the increments}} = ?$$

$$\delta^{t+1} = (0.7 - 0.5) - (0.5 - 0.4) =$$

$$= 0.2 - 0.1 = 0.1$$

PLATO'S ALLEGORY OF THE CAVE



Plato's Allegory of the Cave
(*The Republic*)

ALL WE SEE IS A DISTORTED IMAGE (Y) OF THE ‘TRUE’ UNDERLYING RATES (X)



$$\begin{array}{|c|c|} \hline y_{FR}^{t+1} & = 0.2 + 0.7 \\ \hline y_{FR}^t & = 0.2 + 0.5 \\ \hline \end{array} \quad \begin{array}{|c|c|} \hline y_{SP}^{t+1} & = 0.1 + 0.5 \\ \hline y_{SP}^t & = 0.1 + 0.4 \\ \hline \end{array}$$

What is δ^{t+1} ?

$$\delta^{t+1} = \underbrace{(y_{FR}^{t+1} - y_{FR}^t) - (y_{SP}^{t+1} - y_{SP}^t)}_{\text{difference in the increments}} = ?$$

$$\delta^{t+1} = (0.9 - 0.7) - (0.6 - 0.5) =$$

$$= 0.2 - 0.1 = 0.1$$

Same as before...

DIFFERENCE IN DIFFERENCES ESTIMATOR



- ▶ To the extent that the bias is additive and, within each country, is constant over short periods of time, DiD estimates from social media data:

$$\delta^{t+1} = (y_i^{t+1} - y_z^{t+1}) - (y_i^t - y_z^t)$$

are good estimates of the underlying differential:

$$\delta^{t+1} = \underbrace{(x_i^{t+1} - x_i^t) - (x_z^{t+1} - x_z^t)}_{\text{difference in the increments}}$$

- ▶ Additive values of the bias (m and n) cancel out



If the bias is expected to be multiplicative:

$$\underbrace{y_i^t}_{\text{Observation from social media for location } i} = \underbrace{n}_{\text{bias for location } i} \times \underbrace{x_i^t}_{\text{"true" rate for location } i}$$

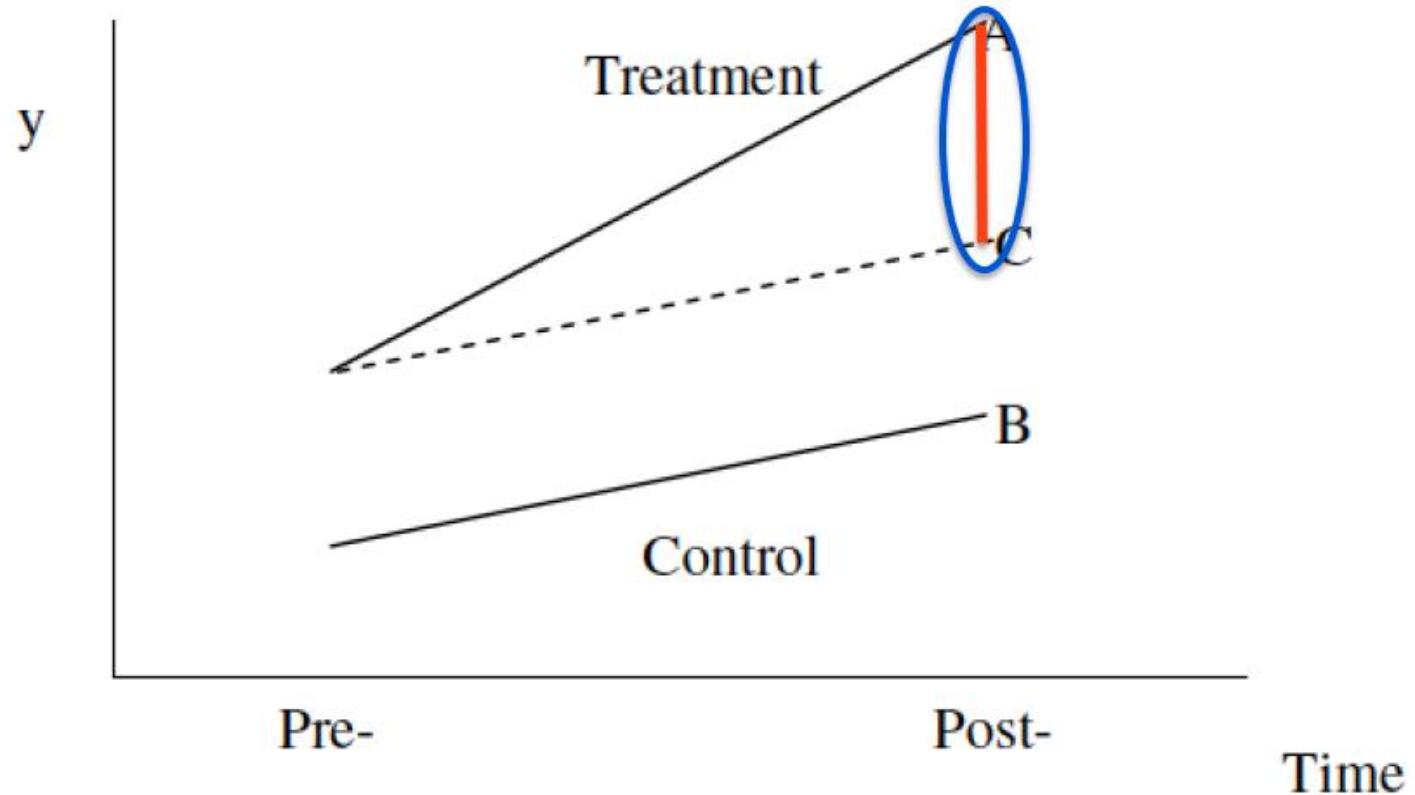
Use a logarithmic transformation

$$\log(y_i^t) = \log(n) + \log(x_i^t)$$

Then use the difference-in-differences estimator on the logs:

$$\delta^{t+1} = [\log(y_i^{t+1}) - \log(y_z^{t+1})] - [\log(y_i^t) - \log(y_z^t)]$$

DIFFERENCE IN DIFFERENCES – IN REGRESSION FORM



$$\hat{\delta} = (\bar{Y}_{treat,post} - \bar{Y}_{control,post}) - (\bar{Y}_{treat,pre} - \bar{Y}_{control,pre})$$

DIFFERENCE IN DIFFERENCES – IN REGRESSION FORM



$$y_{it} = \beta_0 + \beta_1 I(treat_{it}) + \beta_2 I(post_{it}) + \beta_3 I(treat_{it})I(post_{it}) + \epsilon_{it}$$

where:

$I(treat_{it}) = 1$ if treatment; 0 otherwise

$I(post_{it}) = 1$ if post; 0 otherwise

β_3 is the estimate of δ , the difference-in-difference estimator

Why?

DIFFERENCE IN DIFFERENCES – IN REGRESSION FORM



$$y_{it} = \beta_0 + \beta_1 I(treat_{it}) + \beta_2 I(post_{it}) + \beta_3 I(treat_{it})I(post_{it}) + \epsilon_{it}$$

- ▶ $E[y_{it}|I(treat_{it}) = 1, I(post_{it}) = 1] = \beta_0 + \beta_1 + \beta_2 + \beta_3$
- ▶ $E[y_{it}|I(treat_{it}) = 1, I(post_{it}) = 0] = \beta_0 + \beta_1$
- ▶ $E[y_{it}|I(treat_{it}) = 0, I(post_{it}) = 1] = \beta_0 + \beta_2$
- ▶ $E[y_{it}|I(treat_{it}) = 0, I(post_{it}) = 0] = \beta_0$

DIFFERENCE IN DIFFERENCES – IN REGRESSION FORM



$$y_{it} = \beta_0 + \beta_1 I(treat_{it}) + \beta_2 I(post_{it}) + \beta_3 I(treat_{it})I(post_{it}) + \epsilon_{it}$$

	Post	Pre	Difference
Treatment	$\beta_0 + \beta_1 + \beta_2 + \beta_3$	$\beta_0 + \beta_1$	$\beta_2 + \beta_3$
Control	$\beta_0 + \beta_2$	β_0	β_2
Difference	$\beta_1 + \beta_3$	β_1	β_3

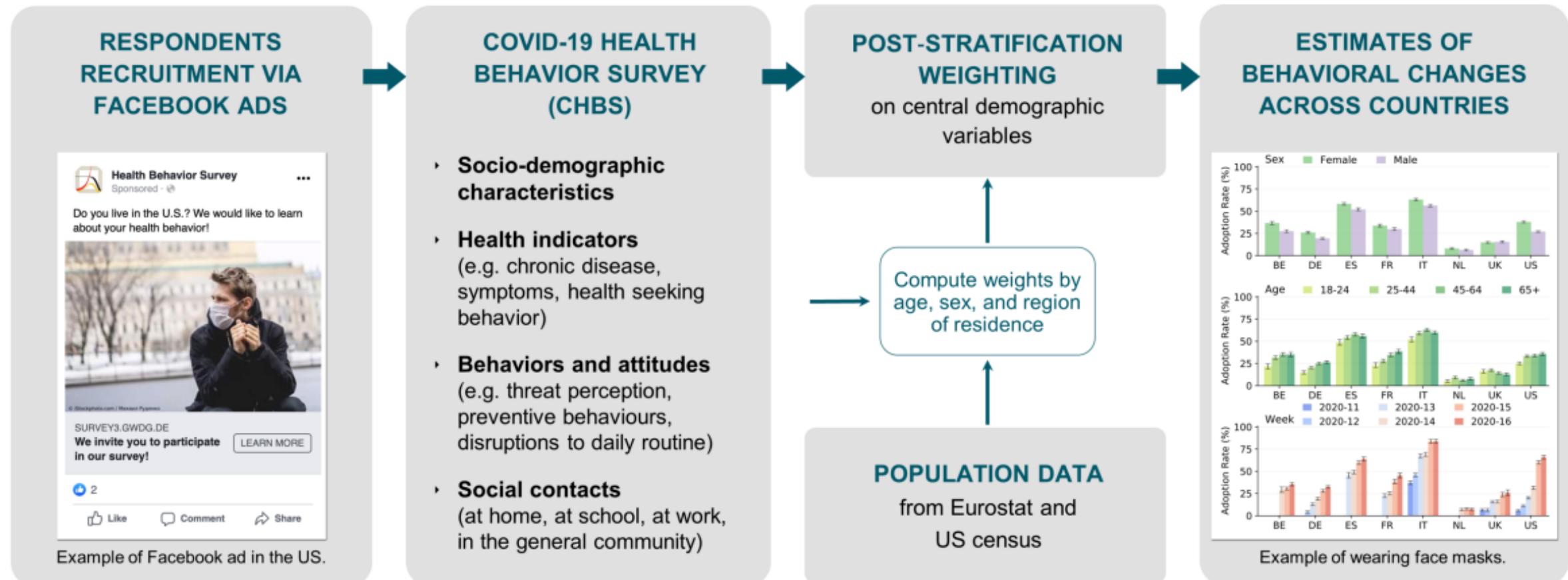


TODAY'S OUTLINE

- TRENDS OVER TIME AND FILTERING OUT BIAS
- SURVEYS AND NON-REPRESENTATIVE SAMPLES
- CALIBRATION AND WEIGHTING



THE COVID-19 HEALTH BEHAVIOR SURVEY (CHBS)



Grow, Perrotta, Del Fava, Cimentada, Rampazzo, Gil-Clavel and Zagheni (2020)
Journal of Medical Internet Research





Health Behavior Survey
Sponsored • ⓘ

Do you live in the U.S.? We would like to learn about your health behavior!



© iStockphoto.com / Михаил Руденко

SURVEY3.GWDG.DE
We invite you to participate in our survey!

[LEARN MORE](#)

2

Like Comment Share



MAX PLANCK INSTITUTE FOR DEMOGRAPHIC RESEARCH

Load unfinished survey

0%

Health Behavior Survey United States

Thank you for your interest in our survey. This is an academic study led by researchers from the [Max Planck Institute for Demographic Research](#). Our goal is to understand health behaviors and help improve the health and well-being of people across countries at a time of increasing uncertainty. Your participation is crucial for our study.

The survey is directed at all people who are 18 years old or older, and it will take approximately 10-12 minutes to complete. Your participation is completely voluntary, you can stop participating at any time. In case you are not comfortable answering a particular question, you have the option to select "Prefer not to answer". Participants' data will be treated anonymously, and we will not ask for identifying information. If you wish to provide your email address, it will not be linked to your survey data. You can download our data protection policy [here](#).

If you have any questions about this research study, please contact Dr. André Grow and Dr. Daniela Perrotta at the Max Planck Institute for Demographic Research via healthsurvey@demogr.mpg.de.

Tick the box below and press 'Next' to continue.

I am willing to participate in this survey. I am at least 18 years old, and I have read the data protection policy.

[Next](#)

Coverage: 8 Countries (BE, DE, ES, FR, IT, NL, UK, US)

Field period: Continuously running between March 13 and August 12, 2020

Topics: Socio-demographics, attitudes, behavior, health, social contacts

Total N: 144,034



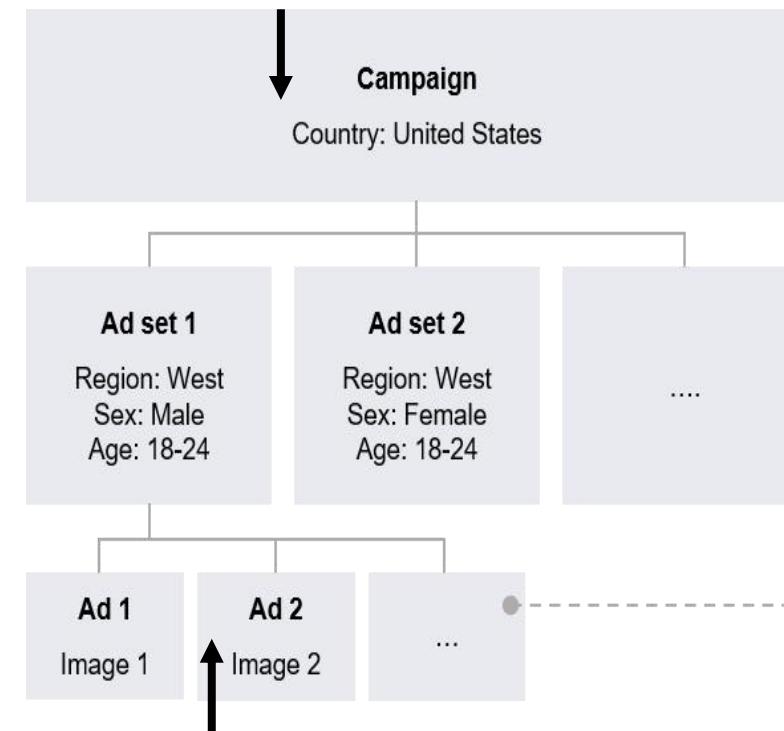


THE ADVERTISING CAMPAIGNS

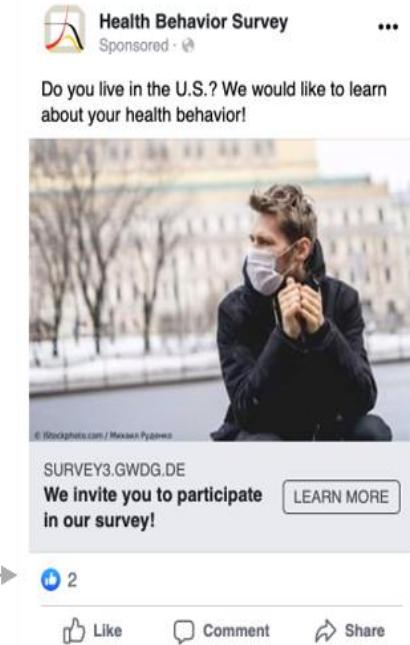
- One ad set/stratum per demographic group:
 - sex (m/f)
 - age (18-24, 25-44, 44-64, 65+)
 - Region of residence
(NUTS1/US Census regions)



One campaign per country



Six different ad images within each ad set

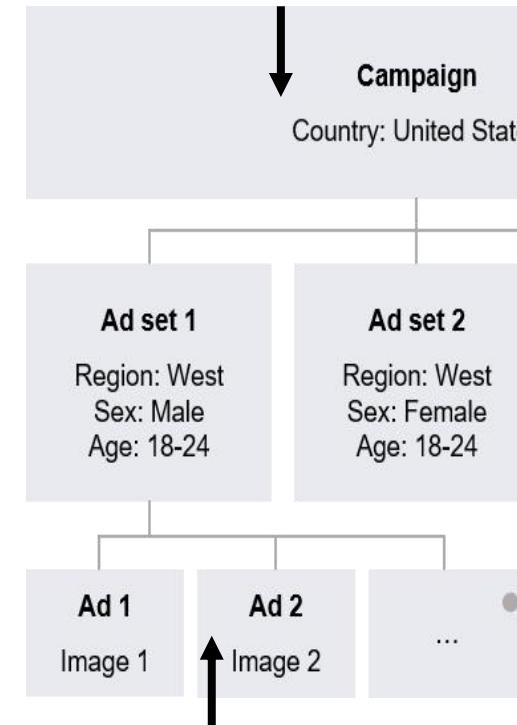


THE ADVERTISING CAMPAIGNS

- One ad set/stratum per demographic group:
 - sex (m/f)
 - age (18-24, 25-44, 44-64, 65+)
 - Region of residence
(NUTS1/US Census regions)



One campaign per country



1 – Male athlete
©Adobe Stock/grki



2 – Group of athletes
©Adobe Stock/nd3000



3 – Woman blowing nose
©iStockphoto/Goodboy Picture Company



4 – Couple blowing noses
©iStockphoto/Goodboy Picture Company



5 – Woman wearing mask
©Adobe Stock/shintartanya

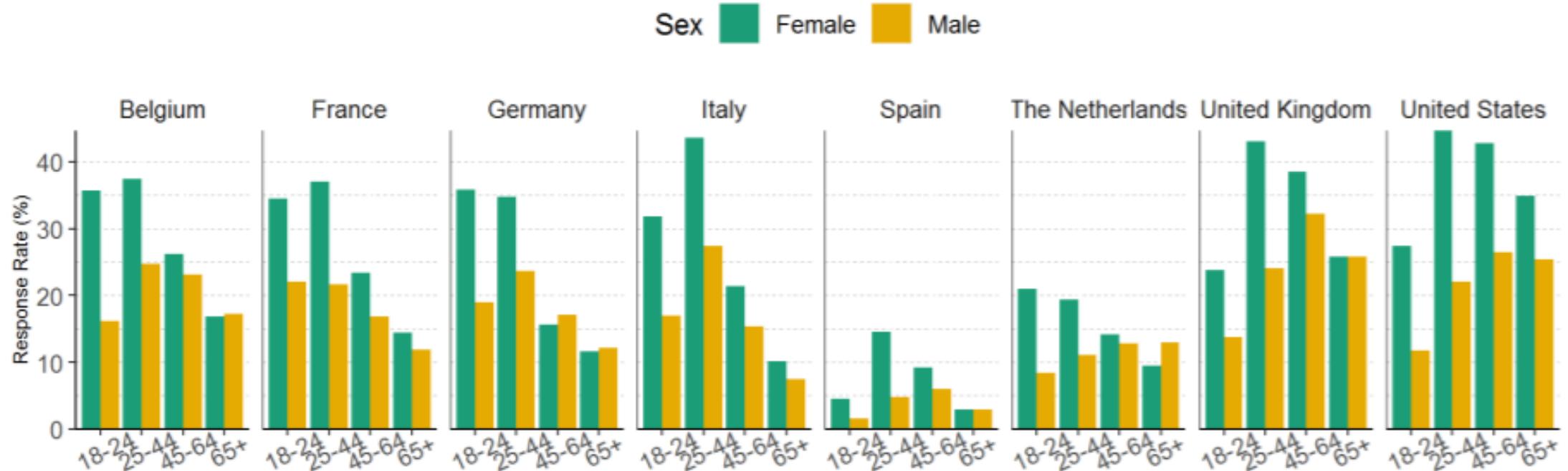


6 – Man wearing mask
©iStockphoto/Михаил Руденко

Six different ad images within each ad set

RESPONSE RATE

- **Response rate** (% Facebook users who completed the questionnaire after clicking on the ad) generally reasonable with some variability
- **Observation 1:** Spain
- **Observation 2:** Women > Men



D. Perrotta, A. Grow, F. Rampazzo, J. Cimentada, E. Del Fava, S. Gil-Clavel & E. Zagheni. 2020. "Behaviors and Attitudes in Response to the COVID-19 Pandemic: Insights from a Cross-National Facebook Survey". *medRxiv*, doi: <https://doi.org/10.1101/2020.05.09.20096388>.



POST-STRATIFICATION WEIGHTING

After the data had been collected, adjustments were produced to obtain a closer approximation to a representative sample

$$w_i = \frac{p_i}{\hat{p}_i}$$

Weight assigned to respondents in stratum i in the survey

Proportion of individuals who belong to stratum i in the population

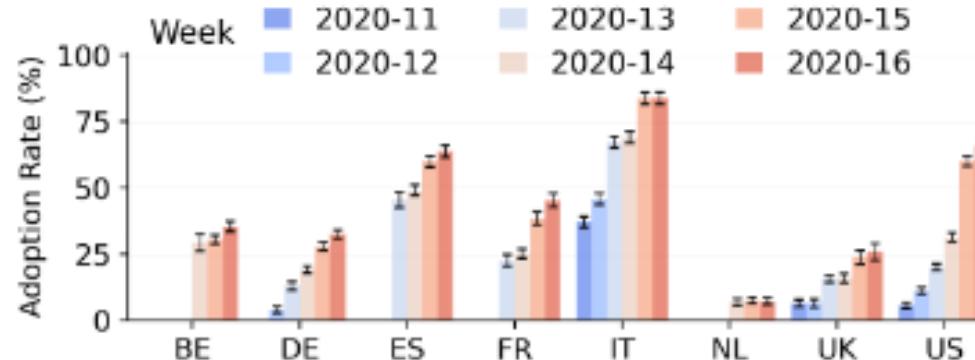
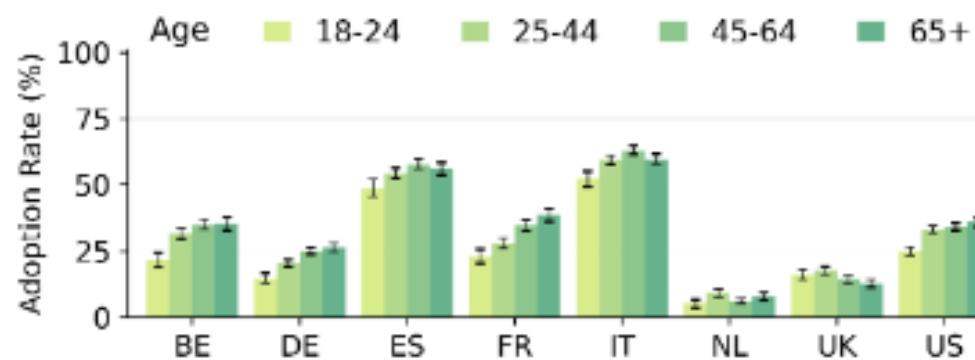
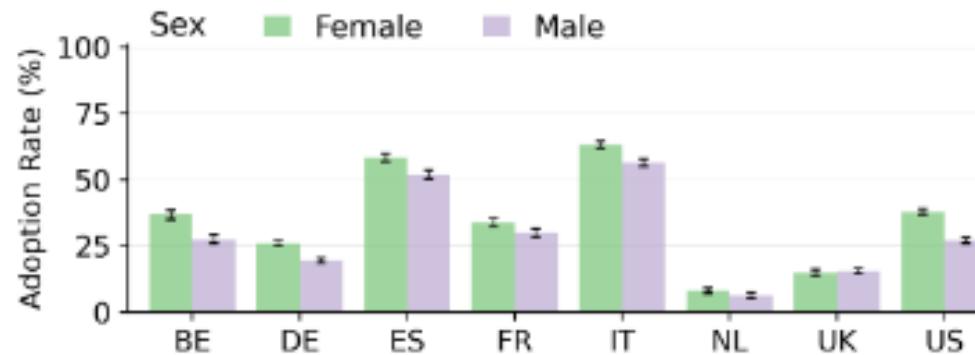
Proportion of respondents who belong to stratum i in the survey

Region	Sex	Age	Facebook		Eurostat		
			\hat{N}_i	\hat{p}_i	N_i	p_i	w_i
Central	Female	18-24	199	0.027	522,730	0.010	0.363
Central	Female	25-44	384	0.053	1,453,810	0.028	0.523
Central	Female	45-64	265	0.037	1,876,591	0.036	0.979
Central	Female	65+	104	0.014	1,605,018	0.031	2.133
Central	Male	18-24	119	0.016	572,911	0.011	0.665
Central	Male	25-44	214	0.030	1,454,223	0.028	0.939
Central	Male	45-64	176	0.024	1,758,096	0.034	1.380
Central	Male	65+	78	0.011	1,221,951	0.023	2.165

Table S1. Example of post-stratification reweigh using the central region of Italy.

D. Perrotta, A. Grow, F. Rampazzo, J. Cimentada, E. Del Fava, S. Gil-Clavel & E. Zagheni. 2020. "Behaviors and Attitudes in Response to the COVID-19 Pandemic: Insights from a Cross-National Facebook Survey". *medRxiv*, doi: <https://doi.org/10.1101/2020.05.09.20096388>.

ADOPTION RATE OF WEARING A FACE MASK



*Perrotta, Grow, Rampazzo,
Cimentada, Del Fava, Gil-Clavel and
Zagheni (2020)*

INITIAL VALIDATION OF FACEBOOK DATA

User's actual
characteristics



Facebook's
classification

FRACTION OF ALL CLASSIFICATIONS THAT ARE CORRECT

	sex	age	region
Belgium	0.98	0.96	0.91
France	0.98	0.93	0.95
Germany	0.98	0.95	0.97
Italy	0.99	0.95	0.97
Netherlands	0.98	0.96	0.98
Spain	0.99	0.94	0.97
United Kingdom	0.99	0.94	0.93
United States	0.99	0.94	0.98

Belgium and the
United Kingdom
somewhat stand out

Across countries, respondents were most likely to be correctly classified based on sex, and a little less likely based on age and region

FRACTION OF ALL CLASSIFICATIONS AS A SPECIFIC REGION THAT ARE CORRECT



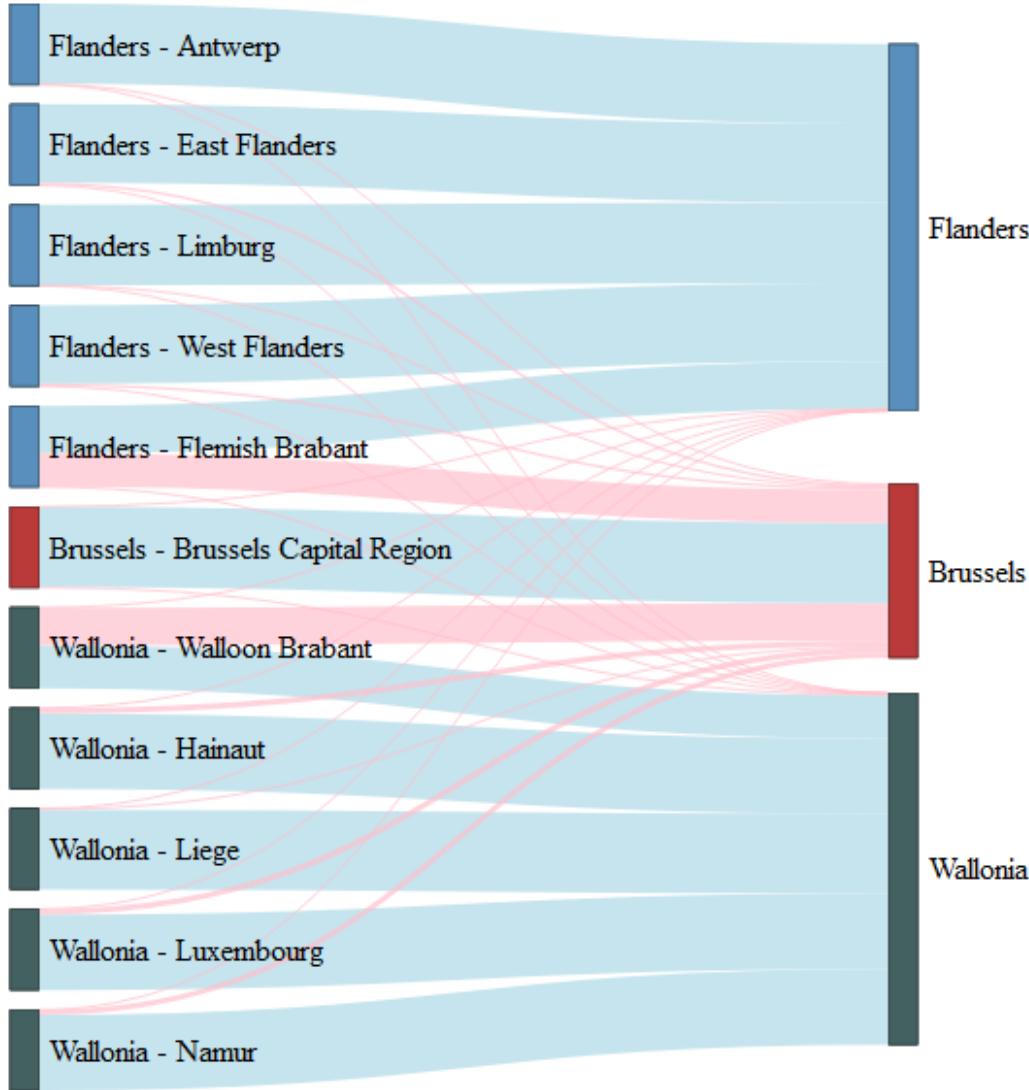
In most countries,
the fraction was
>.9 across regions

Country	Region	Precision
Belgium	Brussels	0.664
	Flanders	0.992
	Wallonia	0.992
	East	0.955
	North	0.976
	South	0.981
Germany	West	0.978
	England	0.992
	London	0.678
	Northern Ireland	0.993
	Scotland	0.991
	Wales	0.908
United Kingdom	Midwest	0.986
	Northeast	0.982
	South	0.979
	West	0.985
United States		

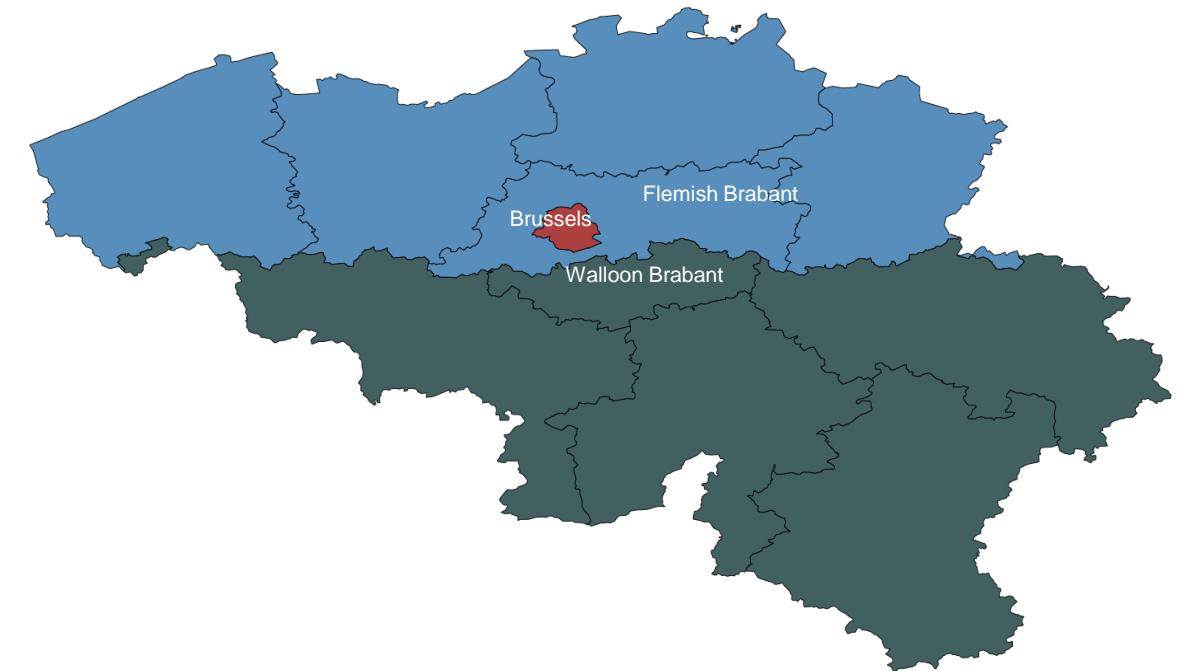
Belgium and the
United Kingdom
deviate from this

A CLOSER LOOK AT REGIONS IN BELGIUM

REPORTED



FACEBOOK

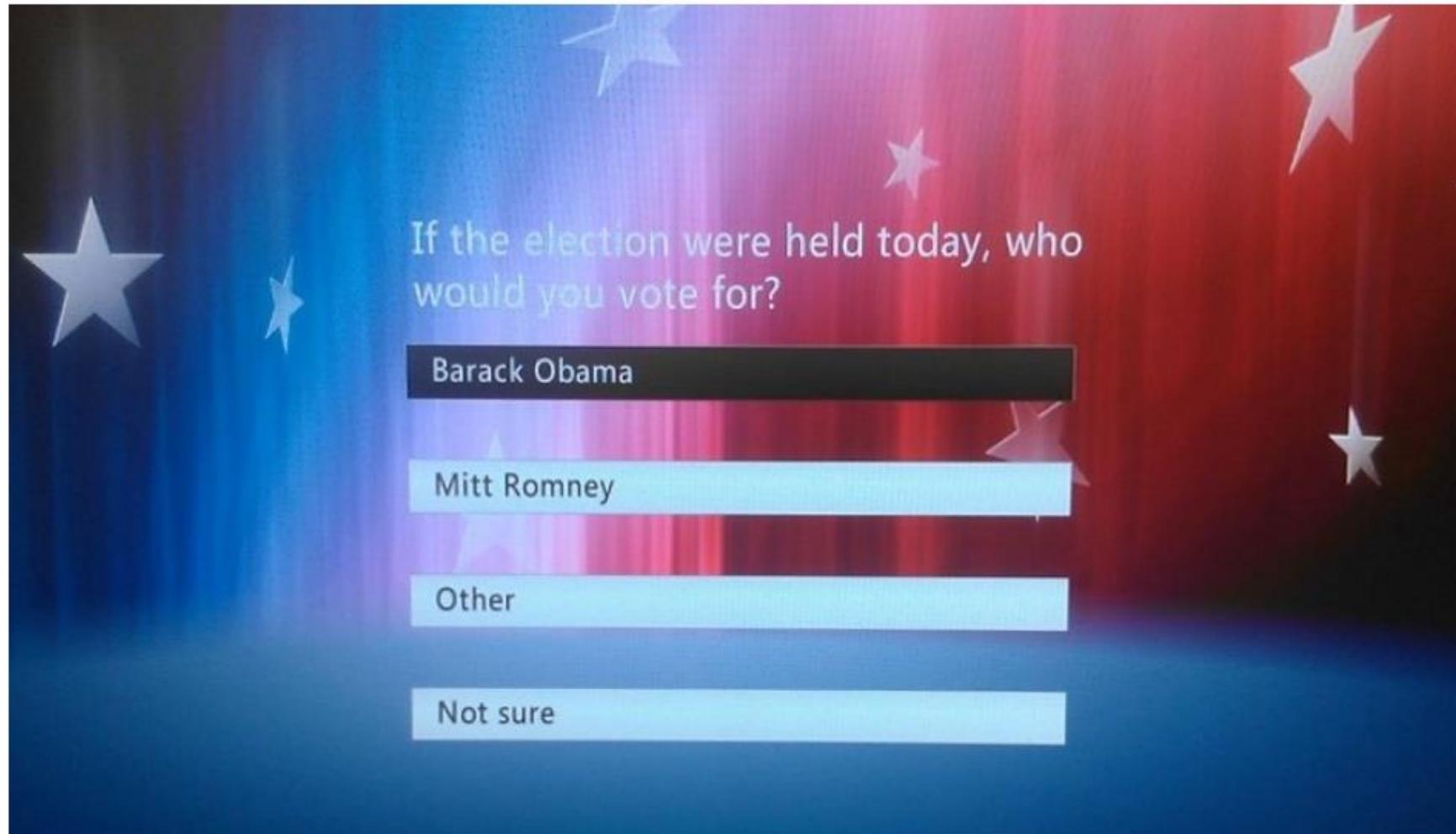


Misclassifications for Brussels mostly concern people who live in adjacent regions



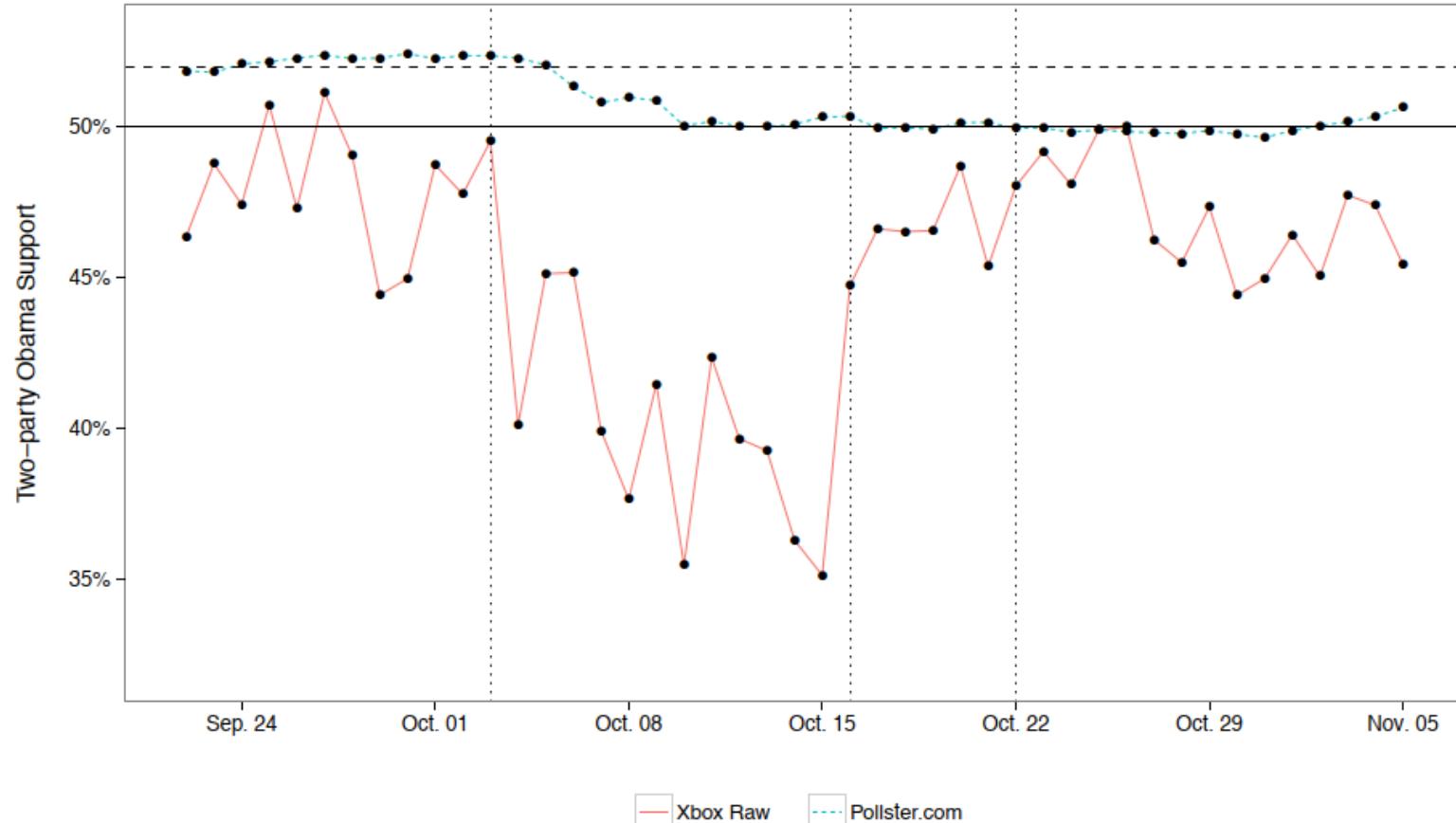
MORE NON-REPRESENTATIVE SAMPLES: A SURVEY OF XBOX USERS

MORE NON-REPRESENTATIVE SAMPLES: A SURVEY OF XBOX USERS





RAW DATA FROM THE XBOX SURVEY



Wang, Rothschild, Goel and Gelman (2014)

APPROACH

Post-stratification:

- ▶ Partition the population into cells, j (e.g., based on combinations of various demographic attributes)
- ▶ Use the sample to estimate the response variable, \hat{y}_j , within each cell
- ▶ Aggregate the cell-level estimates, by weighting them for their relative proportion in the population:

$$\hat{y}^{PS} = \frac{\sum_{j=1}^J N_j \hat{y}_j}{\sum_{j=1}^J N_j}$$



WHAT IS THE PROBLEM?



WHAT IS THE PROBLEM?

- ▶ When considering all possible combination of sex, age, race, education, state, etc., the number of cells becomes very big
- ▶ Some cells would be zeroes, or would have high stochasticity

- ▶ The proposed solution: a multilevel regression model
- ▶ The model is used to estimate the expected value in each cell



LOGISTIC MODEL FOR BINARY RESPONSE (E.G., VOTER VS NON-VOTER)

$$\begin{aligned}\Pr(Y_i \in \{\text{Obama, Romney}\}) = \\ \text{logit}^{-1}(\alpha_0 + \alpha_1(\text{state last vote share}) \\ + a_{j[i]}^{\text{state}} + a_{j[i]}^{\text{edu}} + a_{j[i]}^{\text{sex}} + a_{j[i]}^{\text{age}} + a_{j[i]}^{\text{race}} + a_{j[i]}^{\text{party ID}} + b_{j[i]}^{\text{ideology}} + b_{j[i]}^{\text{last vote}})\end{aligned}$$



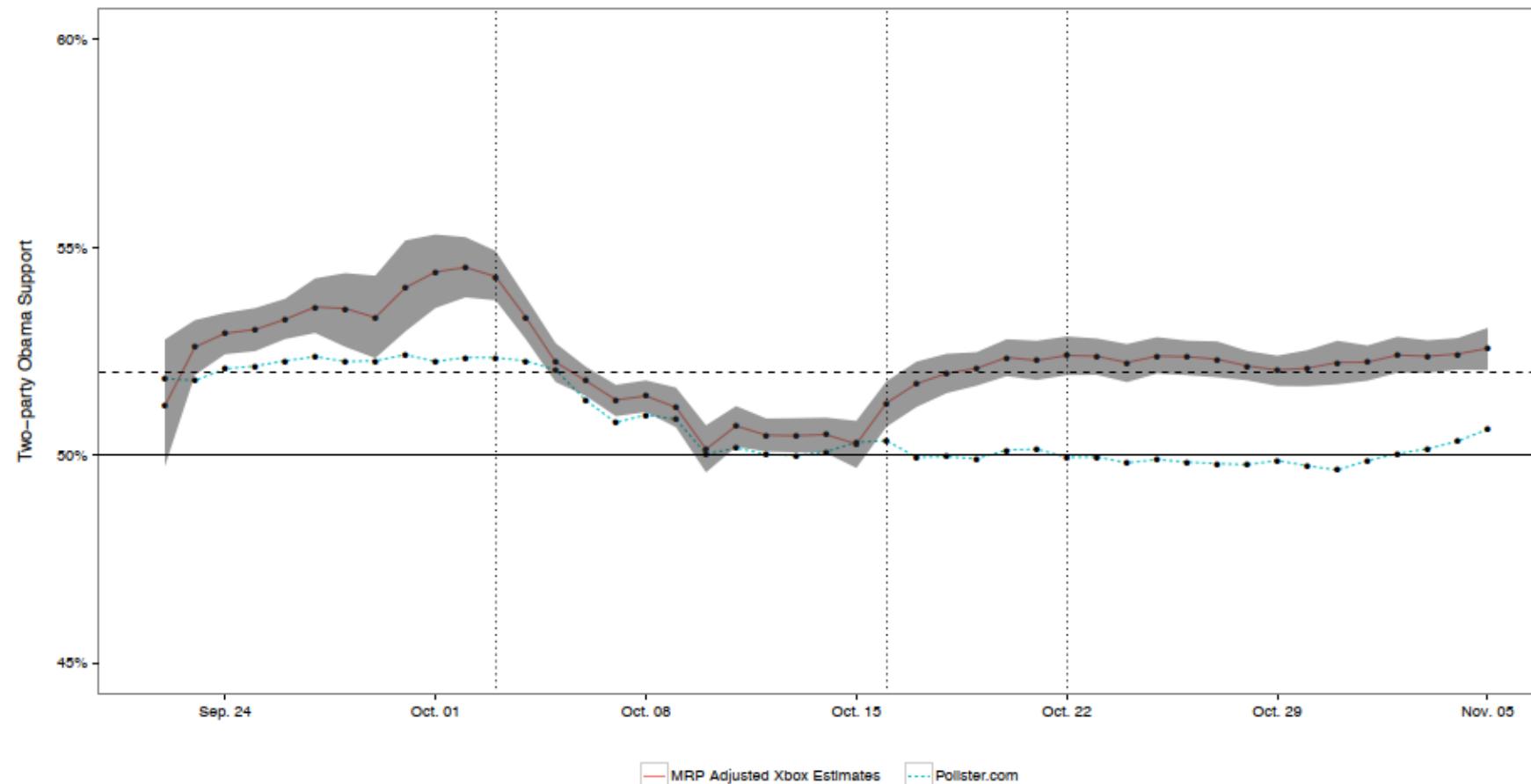
LOGISTIC MODEL FOR BINARY RESPONSE (E.G., OBAMA CONDITIONAL ON BEING A VOTER)

$$\begin{aligned}\Pr(Y_i = \text{Obama} \mid Y_i \in \{\text{Obama}, \text{Romney}\}) &= \\ \text{logit}^{-1}(\beta_0 + \beta_1(\text{state last vote share}) \\ &\quad + b_{j[i]}^{\text{state}} + b_{j[i]}^{\text{edu}} + b_{j[i]}^{\text{sex}} + b_{j[i]}^{\text{age}} + b_{j[i]}^{\text{race}} + b_{j[i]}^{\text{party ID}} + b_{j[i]}^{\text{ideology}} + b_{j[i]}^{\text{last vote}})\end{aligned}$$

$$\begin{aligned}b_{j[i]}^{\text{var}} &\sim N(0, \eta_{\text{var}}^2), \\ \eta_{\text{var}}^2 &\sim \text{inv-}\chi^2(\mu, \eta_0^2).\end{aligned}$$



RESULTS FROM THE XBOX SURVEY AFTER POST-STRATIFICATION ADJUSTMENTS





TODAY'S OUTLINE

- TRENDS OVER TIME AND FILTERING OUT BIAS
- SURVEYS AND NON-REPRESENTATIVE SAMPLES
- CALIBRATION AND WEIGHTING

THE PROBLEM



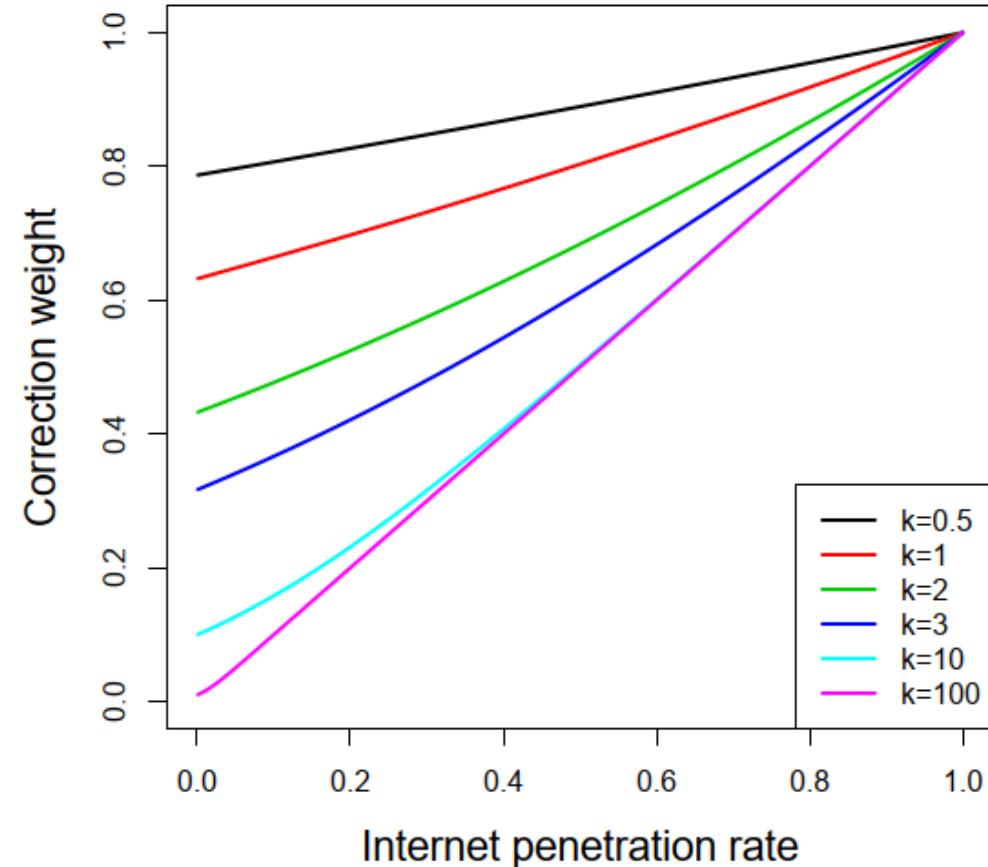
- We have some statistics considered as “silver standard” for Europe
- We have digital trace data on migration flows by age and sex for most countries, but with biases
- We want to be able to say something about migration flows for other regions of the world



THE APPROACH – WHAT ARE THE ASSUMPTIONS?

Define a function (of a parameter k) that corrects for over representation, given age- and sex-specific Internet penetration rates.

$$CF = \frac{pen_{gac}(e^{-k} - 1)}{(e^{-k \times pen_{gac}} - 1)}$$



Source: Zagheni and Weber (2012)

KEY STEPS

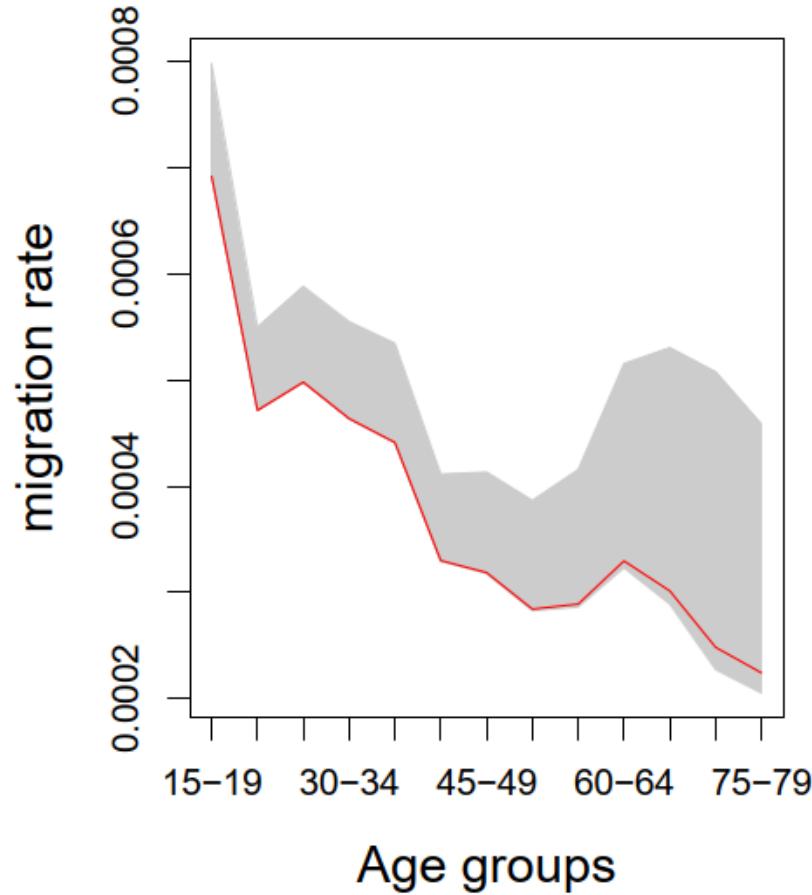


- Estimate the parameters that make the corrected estimates from digital trace data as close as possible to the “official” or “best quality” migration data available (e.g., for Europe)
- Use the estimated parameters and their level of uncertainty to extrapolate to other countries/regions

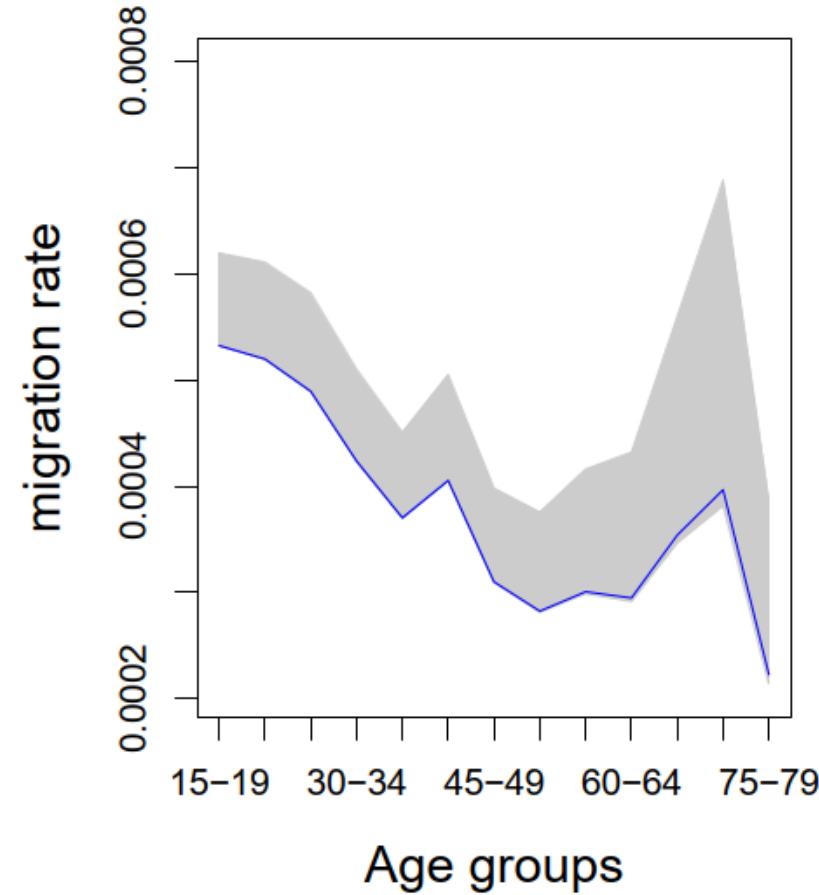
EXAMPLE FOR BRAZIL



BR – Female



BR – Male



Source: Zagheni and Weber (2012)



POST-STRATIFICATION WEIGHTS + ADJUSTMENTS

LET'S EXPLORE IT TOGETHER

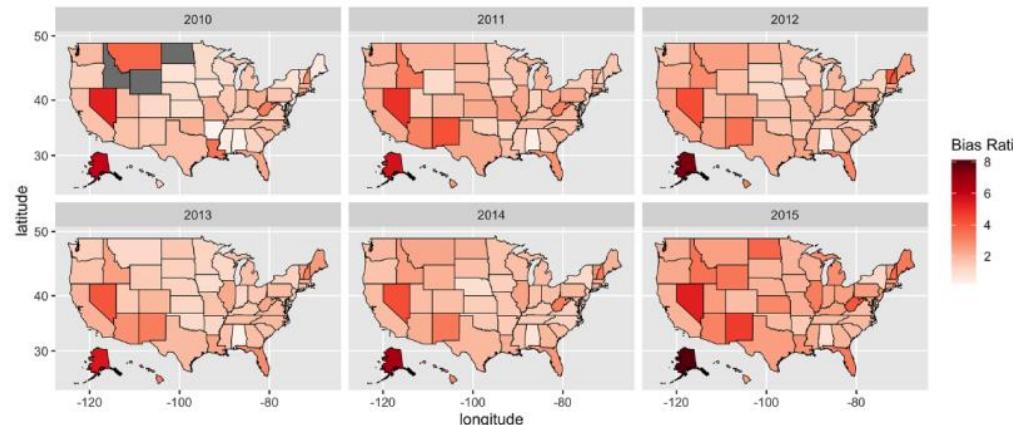
$$\text{Income}_o = \frac{GNIpc_o}{\max(GNIpc_o)}$$

$$W_{o,t} = \frac{1}{\text{Income}_o \times \frac{FBUsers_{o,t}}{Population_{o,t}} + (1 - \text{Income}_o) \times r_t}$$

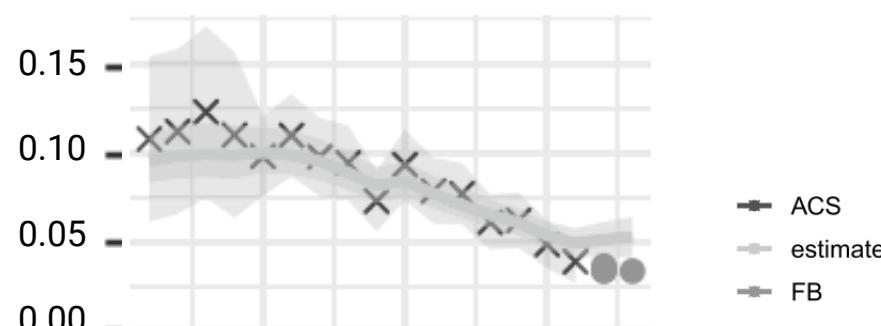
MORE ADVANCED STATISTICAL METHODS TO LEVERAGE TIMELY SOCIAL MEDIA DATA FOR PROVISIONAL ESTIMATES



Bias ratios by state



Fraction of Mexican men, 25-29, in Georgia US



Original Article

Modeling the Bias of Digital Data: An Approach to Combining Digital With Official Statistics to Estimate and Predict Migration Trends

Sociological Methods & Research
1-39
© The Author(s) 2023
Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/00491241221140144
journals.sagepub.com/home/smri



Yuan Hsiao¹ , Lee Fiorio², Jonathan Zagheni¹

Population Research and Policy Review (2022) 41:1–28
<https://doi.org/10.1007/s11113-020-09599-3>

ORIGINAL RESEARCH



Combining Social Media and Survey Data to Nowcast Migrant Stocks in the United States

Monica Alexander¹ , Kivan Polimis², Emilio Zagheni³

Received: 11 April 2020 / Accepted: 20 July 2020 / Published online: 18 August 2020
© The Author(s) 2020

Demography (2021) 58(6):2193–2218
DOI 10.1215/00703370-9578562 © 2021 The Authors
This is an open access article distributed under the terms of a Creative Commons license (CC BY-NC-ND 4.0).

Published online: 9 November 2021

A Framework for Estimating Migrant Stocks Using Digital Traces and Survey Data: An Application in the United Kingdom

Francesco Rampazzo, Jakub Bijak, Agnese Vitali, Ingmar Weber, and Emilio Zagheni

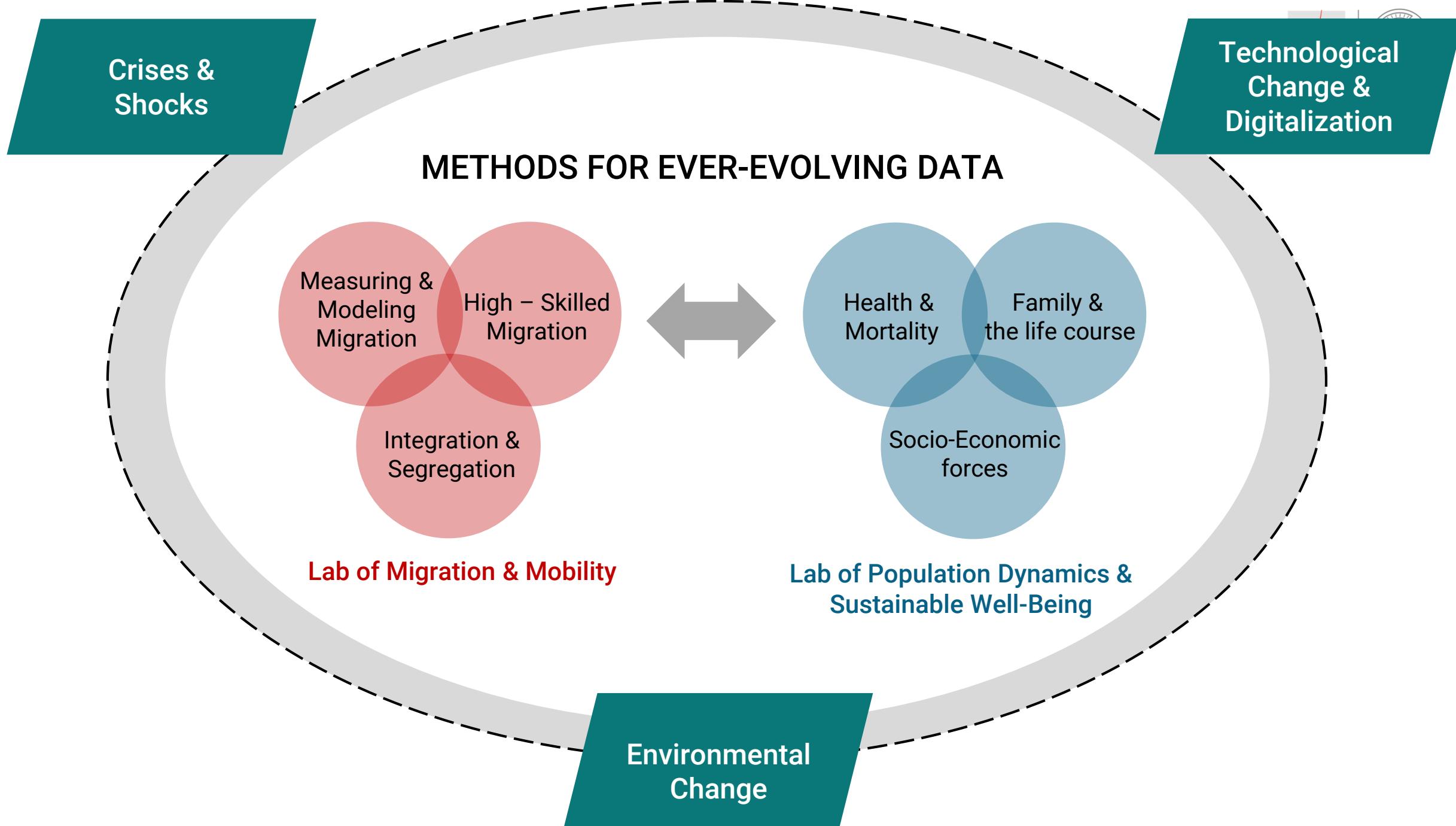
ABSTRACT An accurate estimation of international migration is hampered by a lack of timely and comprehensive data, and by the use of different definitions and measures of migration in different countries. In an effort to address this situation, we complement traditional data sources for the United Kingdom with social media data: our aim is to understand whether information from digital traces can help measure international migration. The Bayesian framework proposed is used to combine data from the Labour Force Survey (LFS) and the Facebook Advertising Platform to study the number of European migrants in the United Kingdom, with the aim of producing more accurate estimates of the numbers of European migrants. The overarching model is divided into a Theory-Based Model of migration and a Measurement Error Model. We review the

This has important implications for understanding policy effectively and for allocation and mobility are often lacking, and in a timely manner. Social media data offer new opportunities for generating more accurate demographic estimates and to complement traditional data sources. The Facebook Advertising Platform, for example, is a valuable source of information that is regularly updated. This is an open access article distributed under the terms of a Creative Commons license (CC BY-NC-ND 4.0).



TODAY'S OUTLINE

- TRENDS OVER TIME AND FILTERING OUT BIAS
- SURVEYS AND NON-REPRESENTATIVE SAMPLES
- CALIBRATION AND WEIGHTING





DIGITAL AND COMPUTATIONAL DEMOGRAPHY @MPIDR



ACKNOWLEDGMENTS & FURTHER DISCUSSION



THE COMMUNITY AT THE MPI FOR DEMOGRAPHIC RESEARCH