

# Combining surveys and big data

Matthew J. Salganik  
Department of Sociology  
Princeton University

Summer Institute in Computational Social Science  
June 22, 2017



	Sampling	Interviews	Data environment
1st era	Area probability	Face-to-face	Stand-alone
2nd era	Random digital dial probability	Telephone	Stand-alone
3rd era	Non-probability	Computer-administered	<a href="#">Linked</a>

Will big data kill surveys?

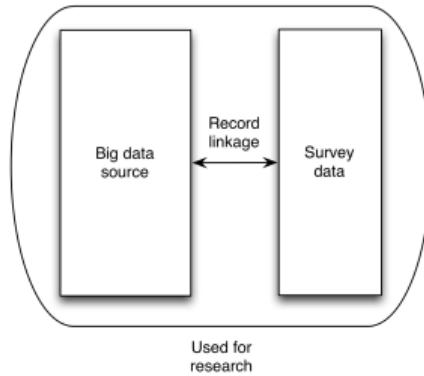


<http://schlitterblog.com/wp-content/uploads/2014/05/peanutbutterlover.jpg>

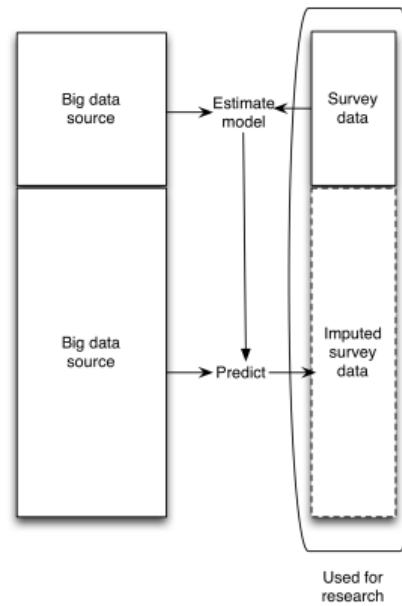


<http://schlitterblog.com/wp-content/uploads/2014/05/peanutbutterlover.jpg>

## Enriched asking



## Amplified asking

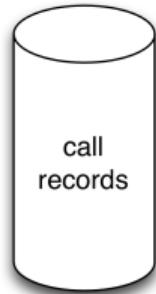


Note the different role of the big data in each case

# Predicting poverty and wealth from mobile phone metadata

Joshua Blumenstock,<sup>1\*</sup> Gabriel Cadamuro,<sup>2</sup> Robert On<sup>3</sup>

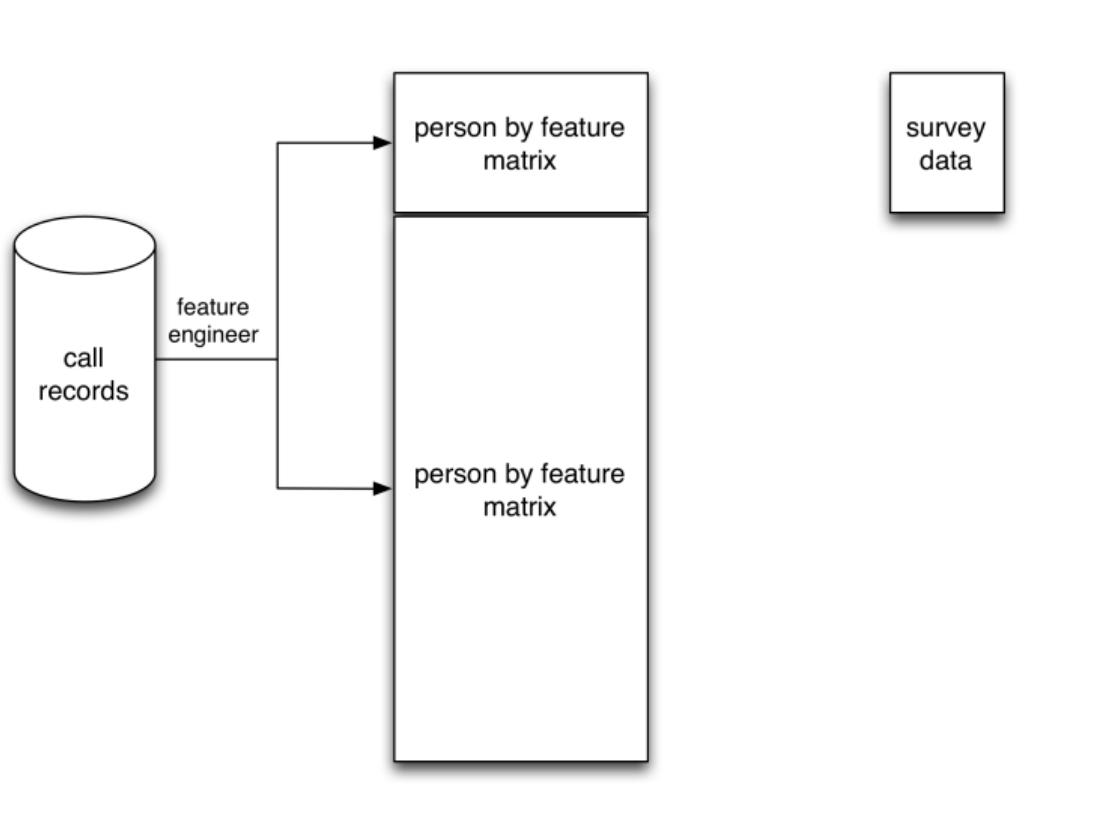


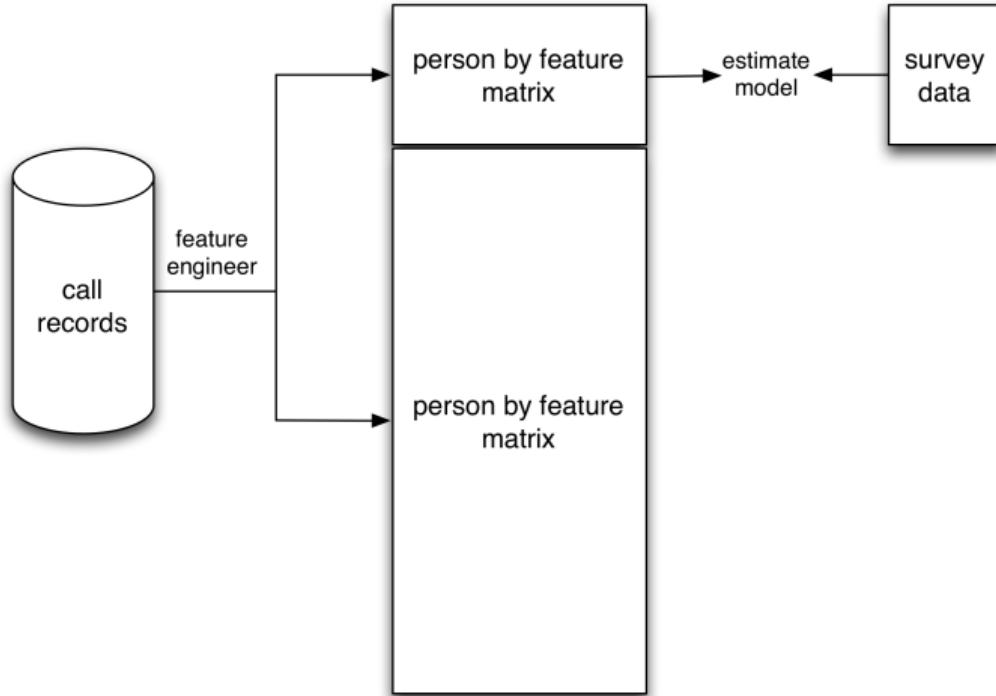


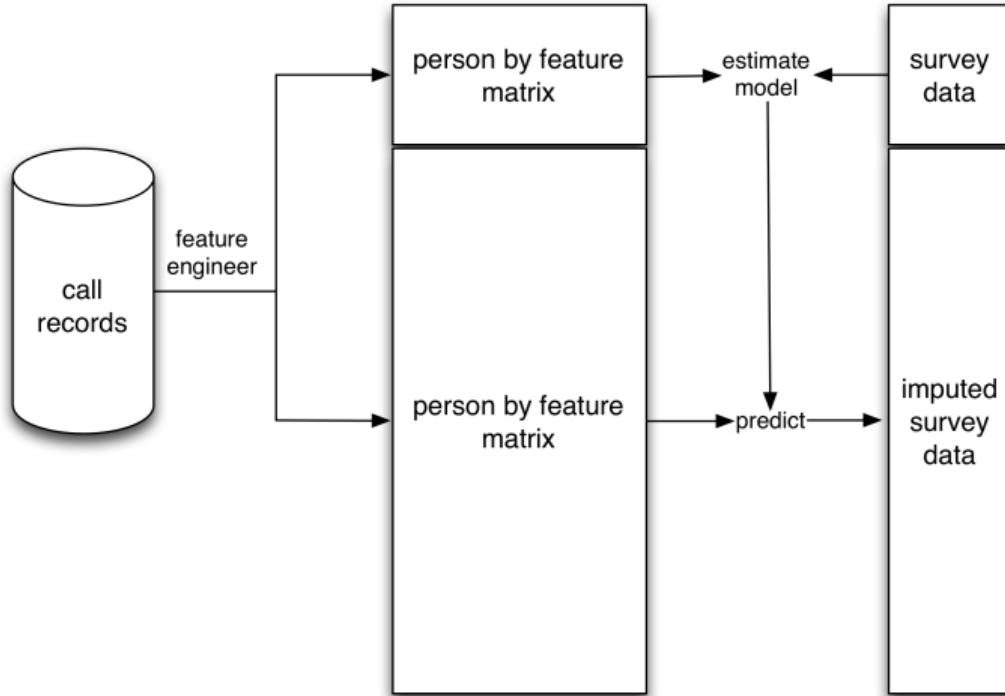
call  
records

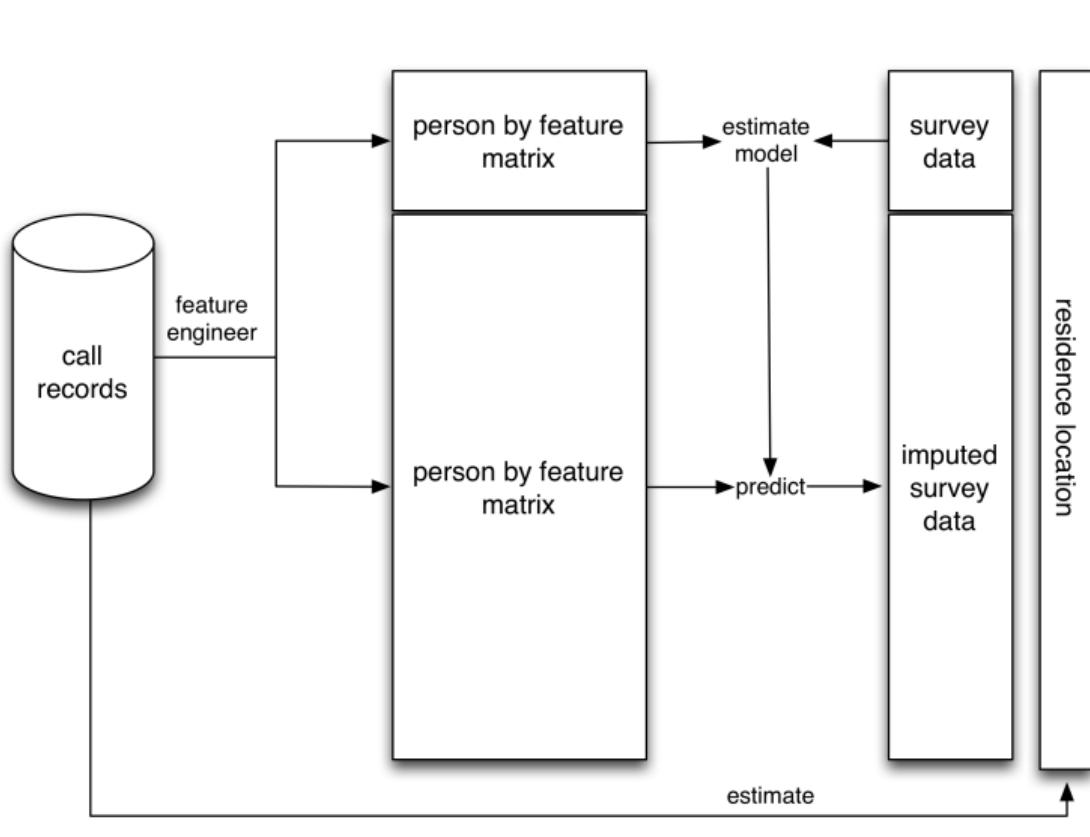


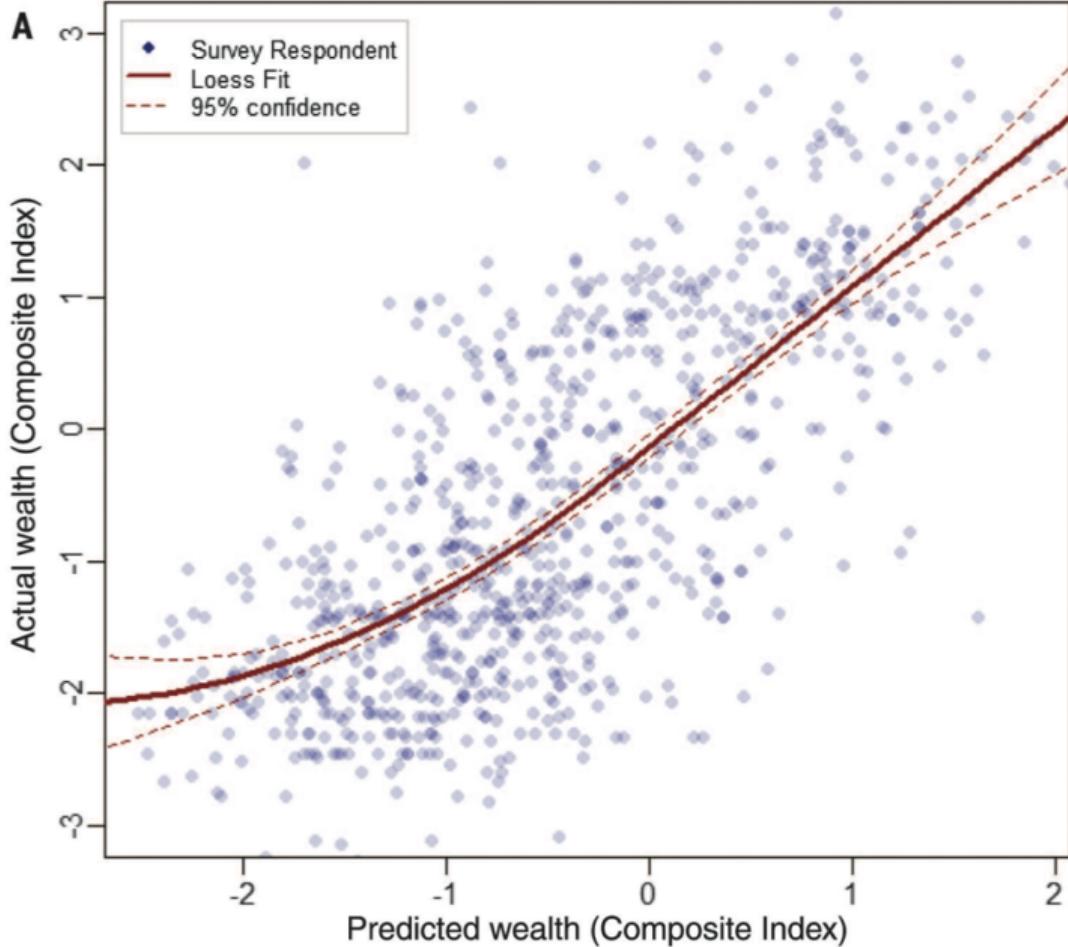
survey  
data

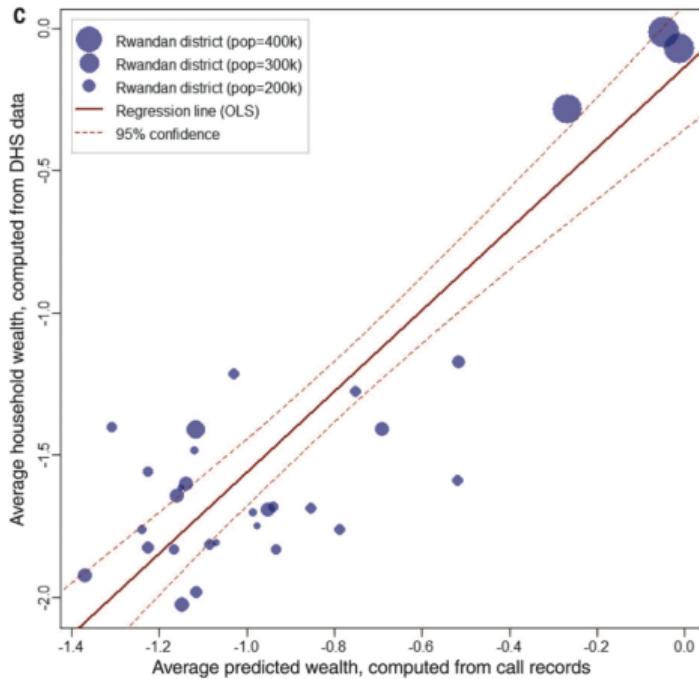


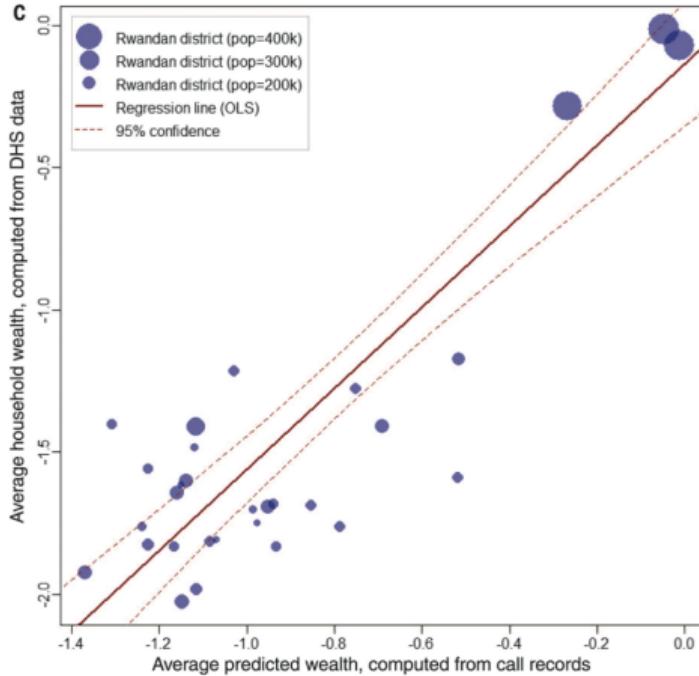












- ▶ 10 times faster
- ▶ 50 times cheaper



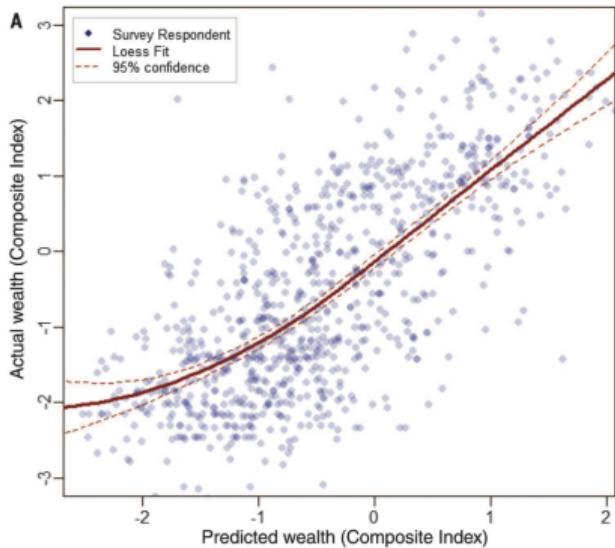
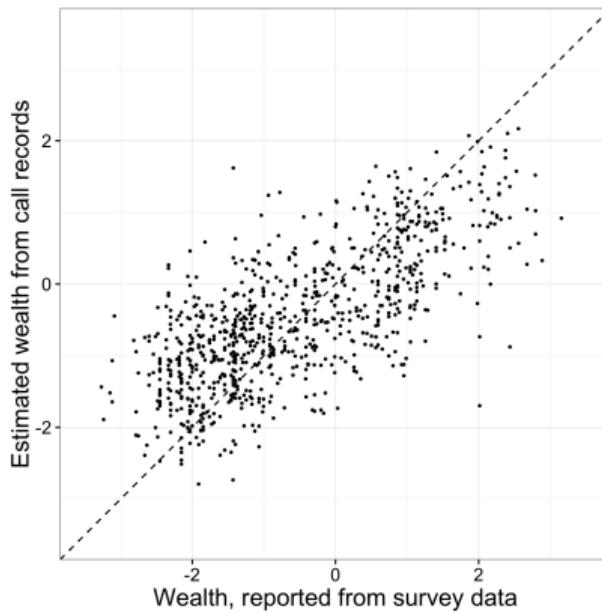
Readymades

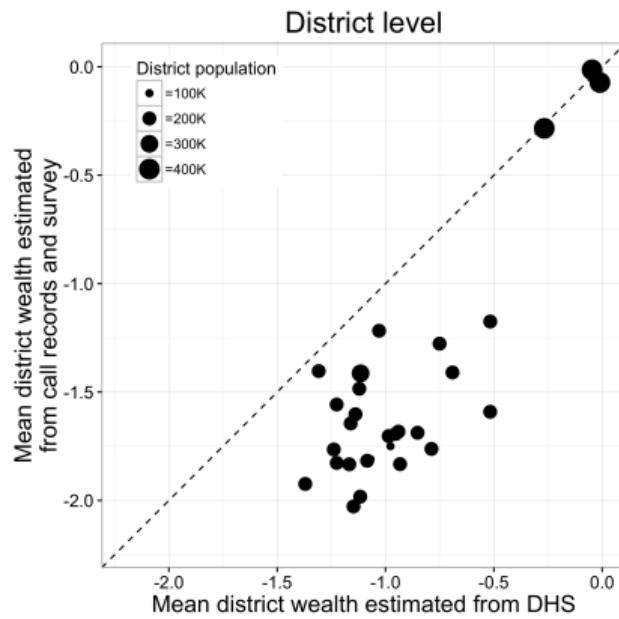
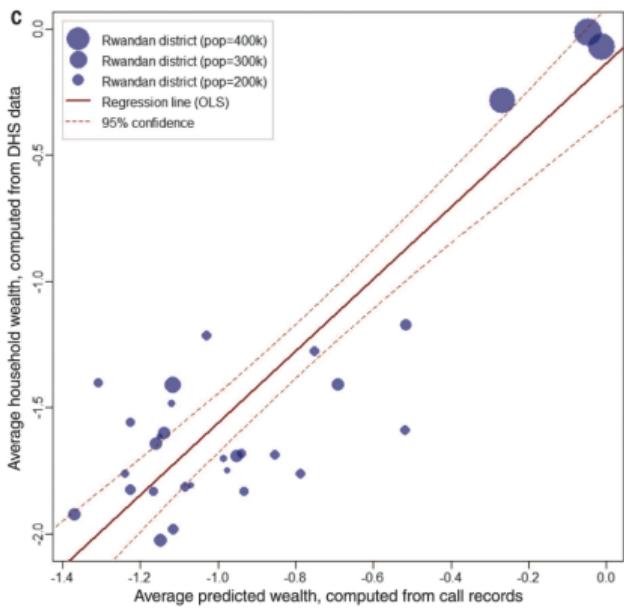
+



Custommades

## Individual level





# Predicting poverty and wealth from mobile phone metadata

Joshua Blumenstock,<sup>1\*</sup> Gabriel Cadamuro,<sup>2</sup> Robert On<sup>3</sup>

## Calling for Better Measurement: Estimating an Individual's Wealth and Well-Being from Mobile Phone Transaction Records

Joshua E. Blumenstock  
University of Washington  
Seattle, WA  
joshblum@uw.edu

## Behavioral Modeling for Churn Prediction: Early Indicators and Accurate Predictors of Customer Defection and Loyalty

Muhammad Raza Khan<sup>1</sup>, Joshua Manoj<sup>2</sup>, Aniket Singh<sup>3</sup>, Joshua Blumenstock<sup>4</sup>  
<sup>1</sup>Information School, University of Washington, Seattle, WA, USA  
Email: [mraza@uw.edu](mailto:mraza@uw.edu), [joshuam@uw.edu](mailto:joshuam@uw.edu), [aniketing@uw.edu](mailto:aniketing@uw.edu), [joshblum@uw.edu](mailto:joshblum@uw.edu)

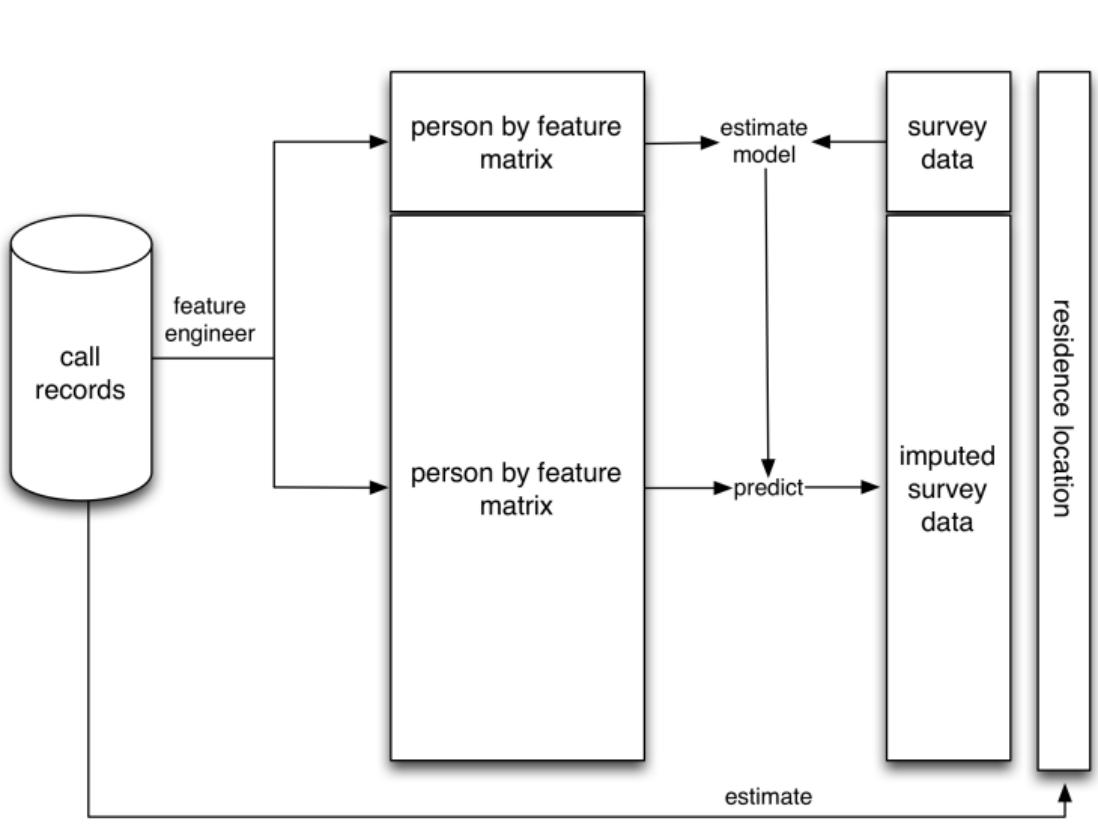
# Combining satellite imagery and machine learning to predict poverty

Neal Jean,<sup>1,2\*</sup> Marshall Burke,<sup>3,4,5\*</sup> Michael Xie,<sup>1</sup> W. Matthew Davis,<sup>4</sup>  
David B. Lobell,<sup>3,4</sup> Stefano Ermon<sup>1</sup>

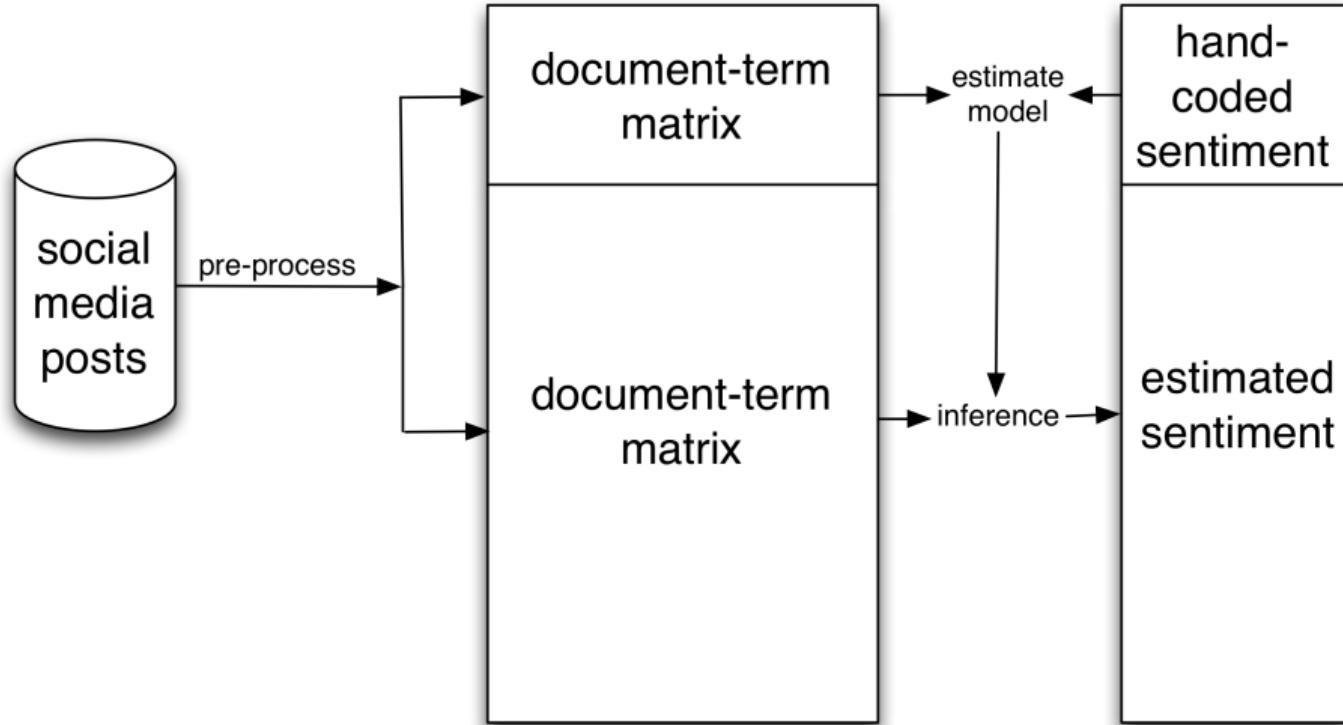
This paper is amazing and surprising. First a digression . . . .

Supervised learning:

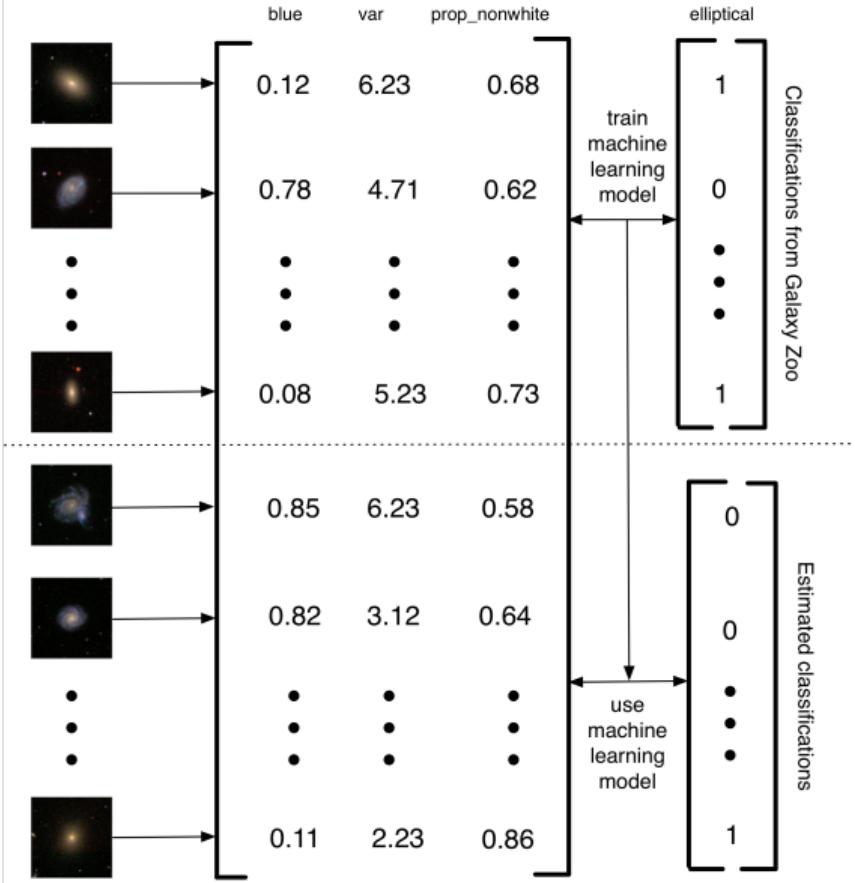
Lots of input-output pairs; goal is to develop a function that will predict the output from the input



See Chapter 3 of Salganik (2016)



See Chapter 2 of Salganik (2016)



See Chapter 5 of Salganik (2016)

Supervised learning:

Lots of input-output pairs; goal is to develop a function that will predict the output from the input

What if rather than engineering the features you could “learn” them automatically?

# Deep learning

Yann LeCun<sup>1,2</sup>, Yoshua Bengio<sup>3</sup> & Geoffrey Hinton<sup>4,5</sup>

<http://dx.doi.org/10.1038/nature14539>

# Combining satellite imagery and machine learning to predict poverty

Neal Jean,<sup>1,2\*</sup> Marshall Burke,<sup>3,4,5\*</sup>† Michael Xie,<sup>1</sup> W. Matthew Davis,<sup>4</sup>  
David B. Lobell,<sup>3,4</sup> Stefano Ermon<sup>1</sup>

# Combining satellite imagery and machine learning to predict poverty

Neal Jean,<sup>1,2\*</sup> Marshall Burke,<sup>3,4,5\*</sup>† Michael Xie,<sup>1</sup> W. Matthew Davis,<sup>4</sup>  
David B. Lobell,<sup>3,4</sup> Stefano Ermon<sup>1</sup>

## Artificial Intelligence Is Predicting Human Poverty From Space

August 18, 2016 // 02:00 PM EST

<http://dx.doi.org/10.1126/science.aaf7894>

[https://motherboard.vice.com/en\\_us/article/artificial-intelligence-is-predicting-human-poverty-from-space](https://motherboard.vice.com/en_us/article/artificial-intelligence-is-predicting-human-poverty-from-space)

Live demo:

<https://www.google.com/maps/place/Kigali,+Rwanda/@-1.9546259,30.0345059,26517m/data=!3m2!1e3!4b1!4m5!3m4!1s0x19dca4258ed8e797:0xf32b36a5411d0bc8!8m2!3d-1.9705786!4d30.1044288>

But, most people had been using night lights



[https://www.nasa.gov/multimedia/imagegallery/image\\_feature\\_2480.html](https://www.nasa.gov/multimedia/imagegallery/image_feature_2480.html)

Prior research:

Nightlights + survey data to estimate wealth in places without surveys

Jean et al. (2016):  
Day pictures + Nightlights + survey data to estimate wealth in places without surveys

## Predicting poverty

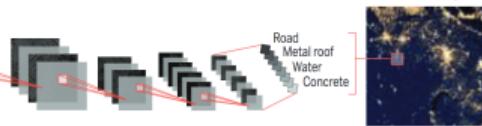
Satellite images can be used to estimate wealth in remote regions.

Neural network learns features in satellite images that correlate with economic activity

Daytime satellite photos capture details of the landscape



Convolutional Neural Network (CNN) associates features from daytime photos with nightlight intensity

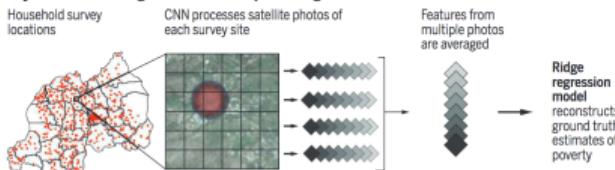


Satellite nightlights are a proxy for economic activity

Daytime satellite images can be used to predict regional wealth

Household survey locations

CNN processes satellite photos of each survey site



Features from multiple photos are averaged

Ridge regression model reconstructs ground truth estimates of poverty

- ▶ Start with CNN pretrained on ImageNet (e.g. hamsters and weasels)

## Predicting poverty

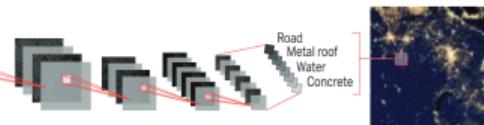
Satellite images can be used to estimate wealth in remote regions.

Neural network learns features in satellite images that correlate with economic activity

Daytime satellite photos capture details of the landscape



Convolutional Neural Network (CNN) associates features from daytime photos with nightlight intensity



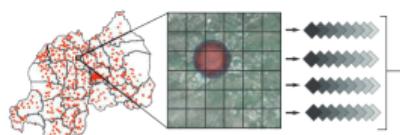
Satellite nightlights are a proxy for economic activity



Daytime satellite images can be used to predict regional wealth

Household survey locations

CNN processes satellite photos of each survey site



Features from multiple photos are averaged

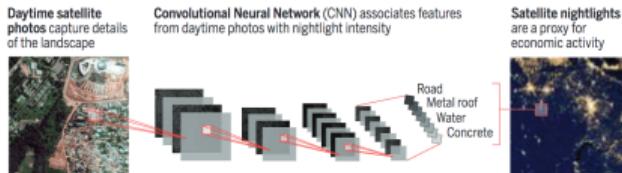
Ridge regression model reconstructs ground truth estimates of poverty

- ▶ Start with CNN pretrained on ImageNet (e.g. hamsters and weasels)
- ▶ Train CNN to predict nightlights from day pictures (lots of training data)

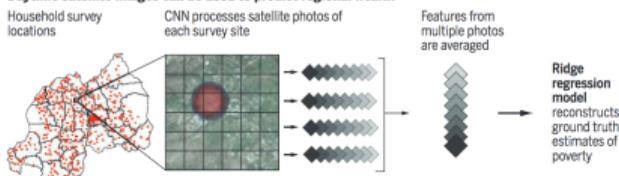
## Predicting poverty

Satellite images can be used to estimate wealth in remote regions.

Neural network learns features in satellite images that correlate with economic activity

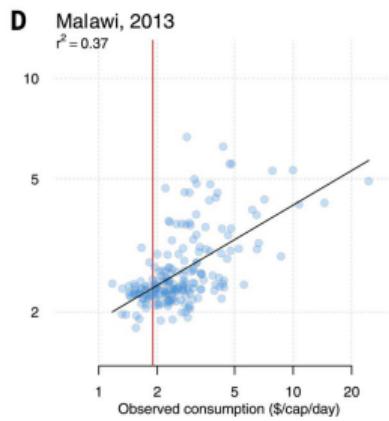
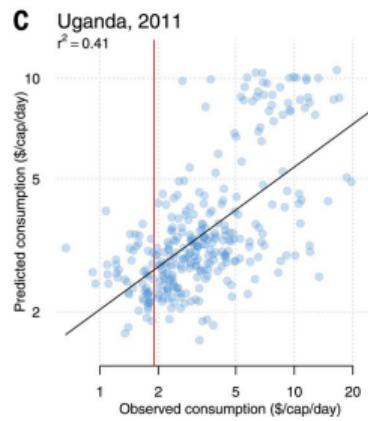
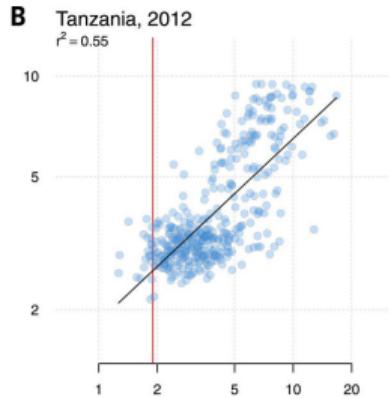
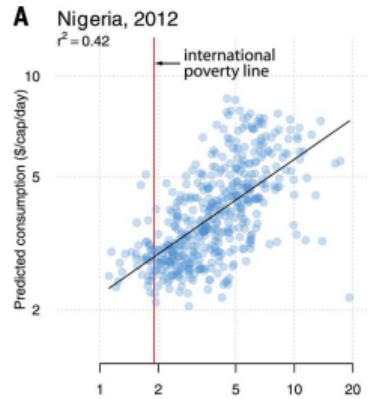


Daytime satellite images can be used to predict regional wealth

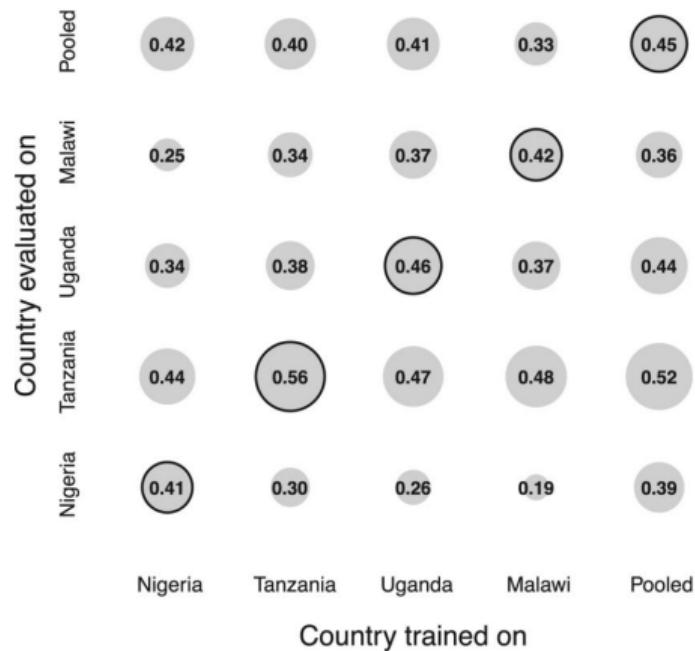


- ▶ Start with CNN pretrained on ImageNet (e.g. hamsters and weasels)
- ▶ Train CNN to predict nightlights from day pictures (lots of training data)
- ▶ Take features from CNN and train ridge regression to predict cluster mean survey response

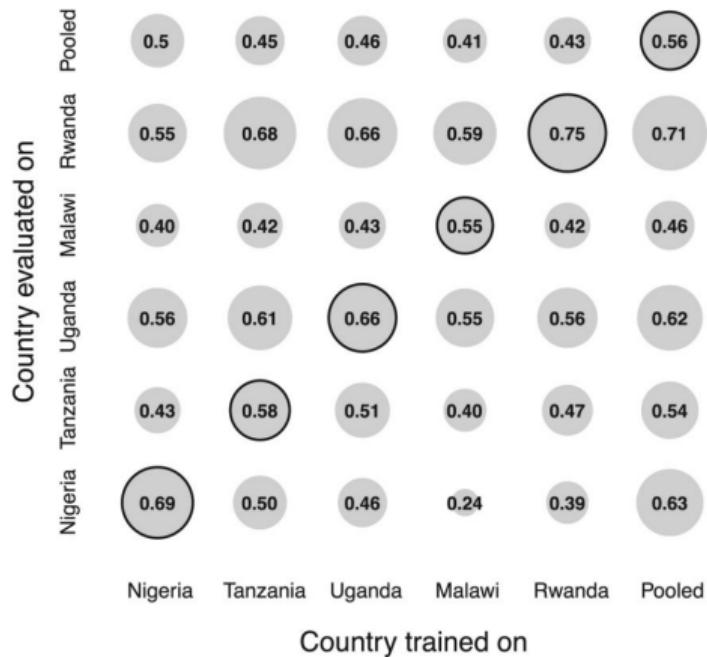
<http://dx.doi.org/10.1126/science.aah5217>

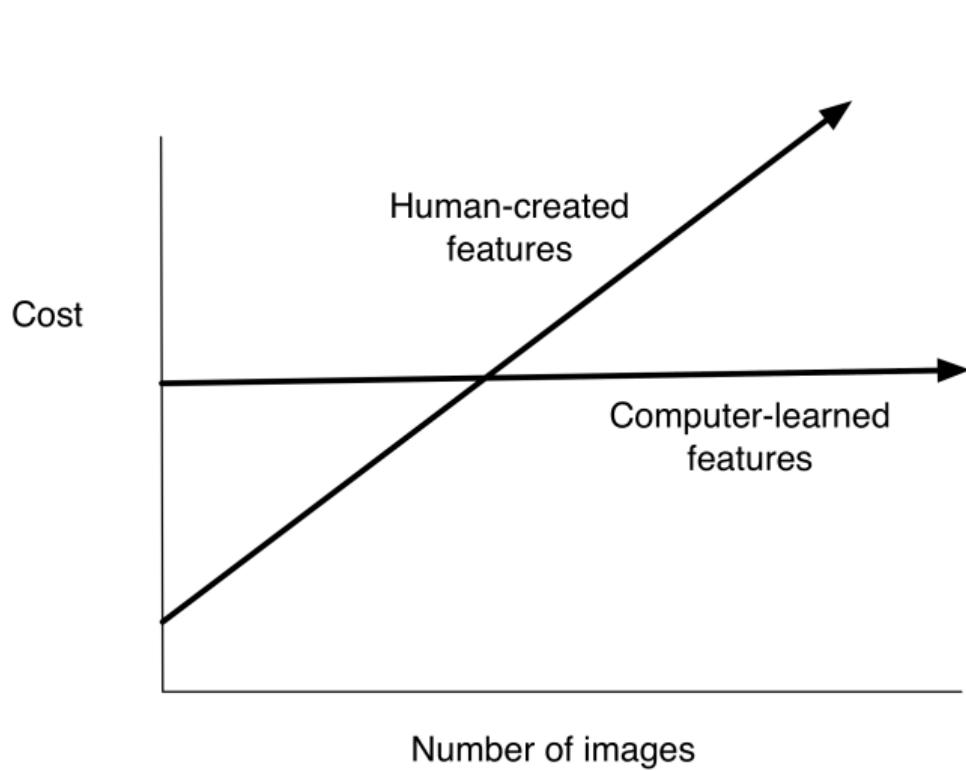


## A Consumption expenditures



## B Assets





 [nealjean / predicting-poverty](#)

[Watch](#) 22   [Star](#) 110   [Fork](#) 60

[Code](#)   [Issues 1](#)   [Pull requests 0](#)   [Projects 0](#)   [Wiki](#)   [Insights](#)

Combining satellite imagery and machine learning to predict poverty

 18 commits    1 branch    0 releases    4 contributors    MIT

Branch: [master](#) [New pull request](#)   [Create new file](#) [Upload files](#) [Find file](#) [Clone or download](#)

 <a href="#">imthexie</a>	select middle of pixel	Latest commit 975fddc on Mar 27
 <a href="#">data/input</a>	Clean replication code	10 months ago
 <a href="#">figures</a>	Fixing cluster prefix in fig_utils.py	7 months ago
 <a href="#">model</a>	Clean replication code	10 months ago
 <a href="#">scripts</a>	select middle of pixel	3 months ago
 <a href="#">.gitignore</a>	Clean replication code	10 months ago
 <a href="#">LICENSE</a>	MIT License	6 months ago
 <a href="#">README.md</a>	Update README.md	8 months ago
 <a href="#">requirements.txt</a>	Clean replication code	10 months ago

<https://github.com/nealjean/predicting-poverty>

## Wrap-up:

- ▶ Surveys and big data are compliments not substitutes

## Wrap-up:

- ▶ Surveys and big data are compliments not substitutes
- ▶ Sometime we do “enriched asking” and sometimes “amplified asking” (role of big data source is different in both cases)

## Wrap-up:

- ▶ Surveys and big data are compliments not substitutes
- ▶ Sometime we do “enriched asking” and sometimes “amplified asking” (role of big data source is different in both cases)
- ▶ Learn more: see “what to read next” in Ch 3 of Bit by Bit.