# How to Start a Data Science Insurrection at an Organization that Would Prefer You Not

Jonathan Hersh (Chapman Argyros School of Business)

*R Stats DC*
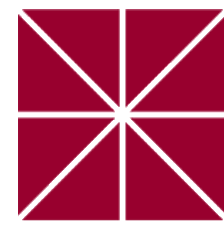
9/9/2018

# (Four Tips for Implementing Machine Learning Projects that Scale within Organizations)

Jonathan Hersh (Chapman Argyros School of Business)

*R Stats DC*

9/9/2018

# When you think of machine learning

# Not necessarily these organizations

facebook

Bai du 百度

Google

NETFLIX

amazon

THE WORLD BANK

United States Census Bureau

cfpb  Consumer Financial Protection Bureau

IDB  Inter-American Develo

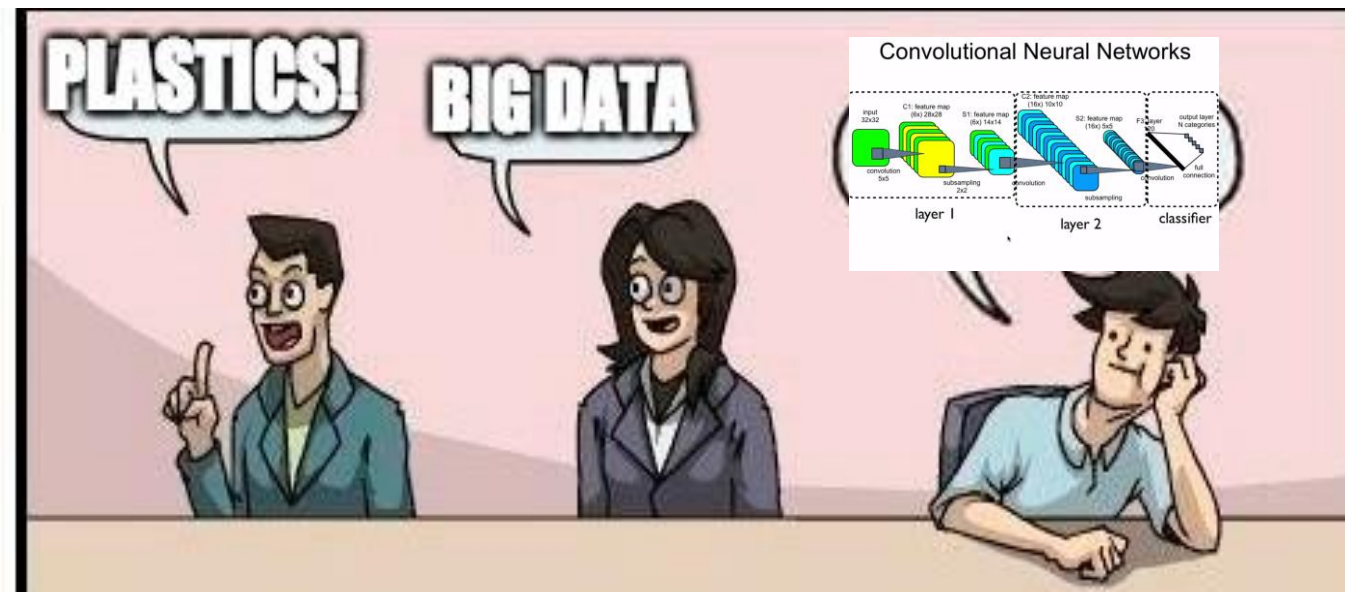UNITED STATES FEDERAL RESERVE SYSTEM

INTERNATIONAL MONETARY FUND

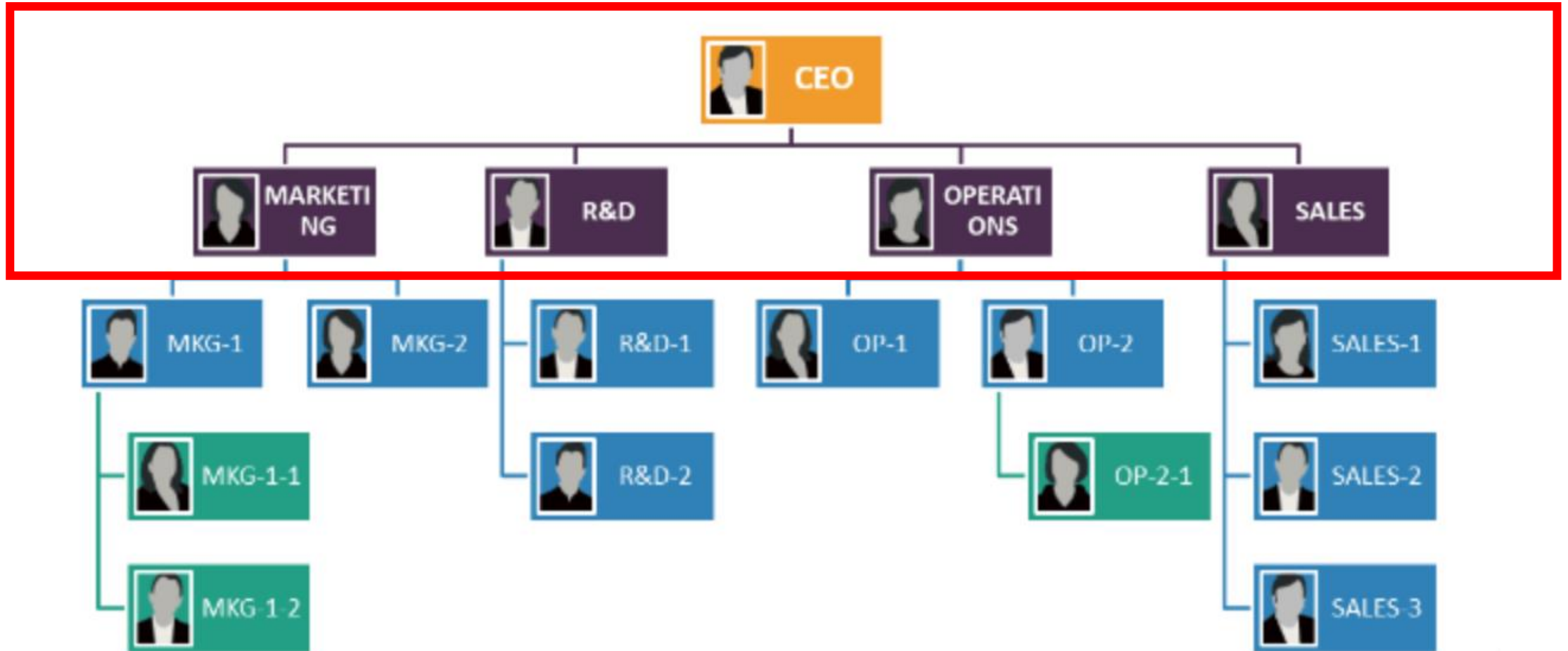# Who actually understands machine learning?

**Tip 1:**

**Explain how your model works in language your boss's boss can understand**

# Use *their* language to explain what you do

# How did these people get to the top?

# Tip 2:

**Motivate machine learning projects by showing how current methods are insufficient**

# Previous Research (Engstrom, Hersh, Newhouse, 2017) Using Intermediate Features to Estimate Poverty in Sri Lanka
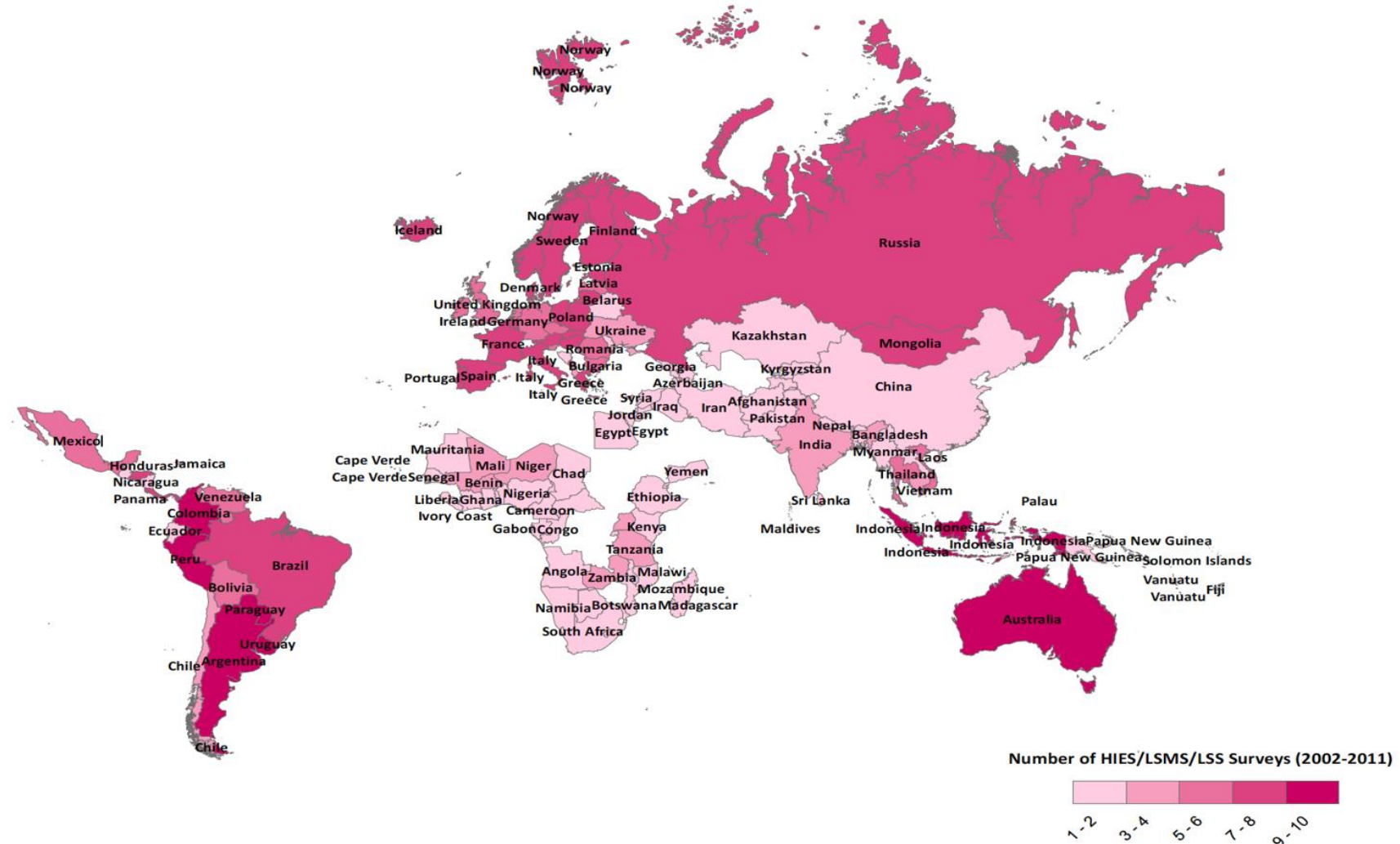
Satellite Image

Buildings

How Most Data is Collected

# 57 Countries Have Zero or One Poverty Estimate 2002-2011
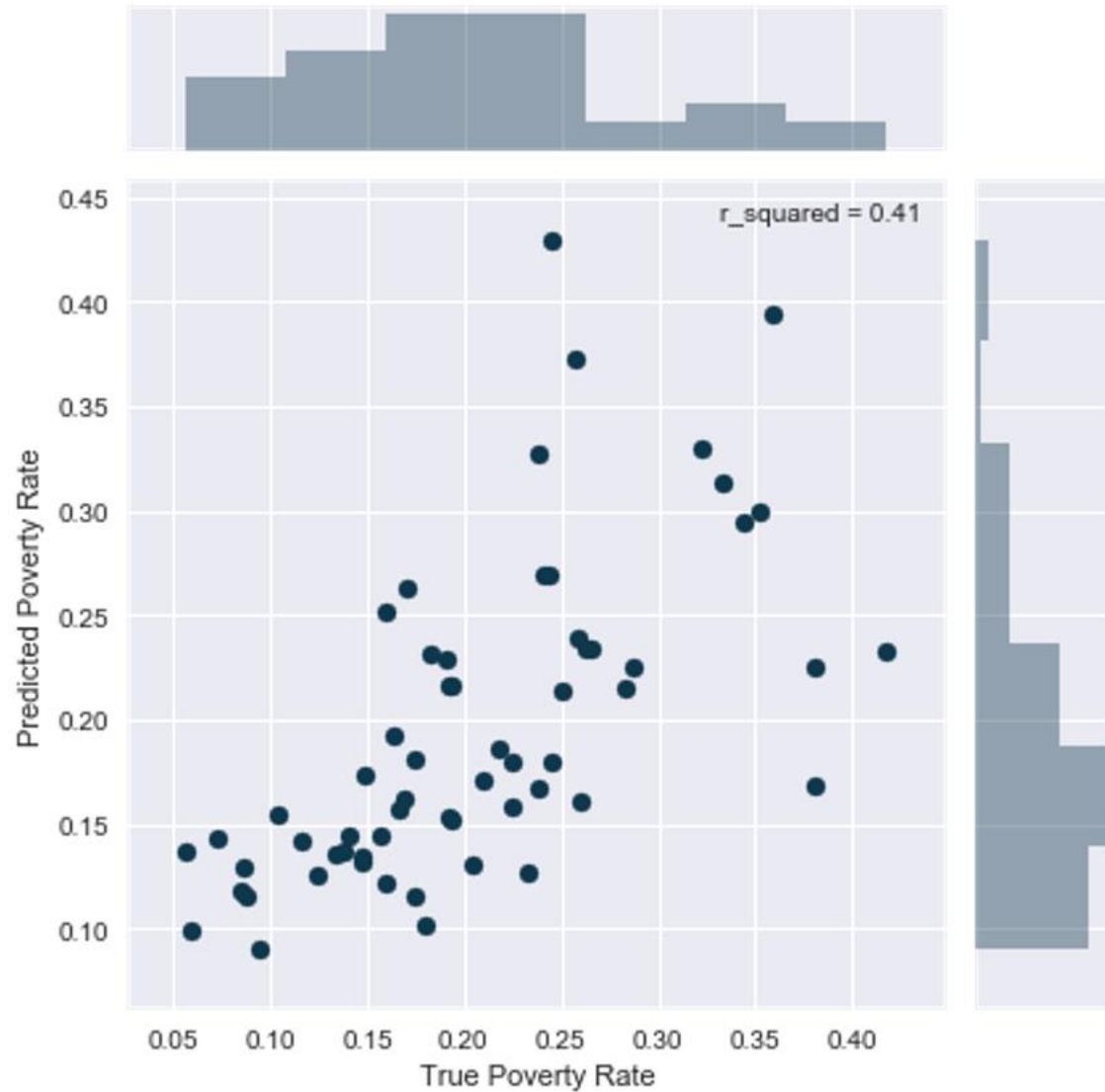


Number of Poverty Data Points, 2002 - 2011

# Tip 3:
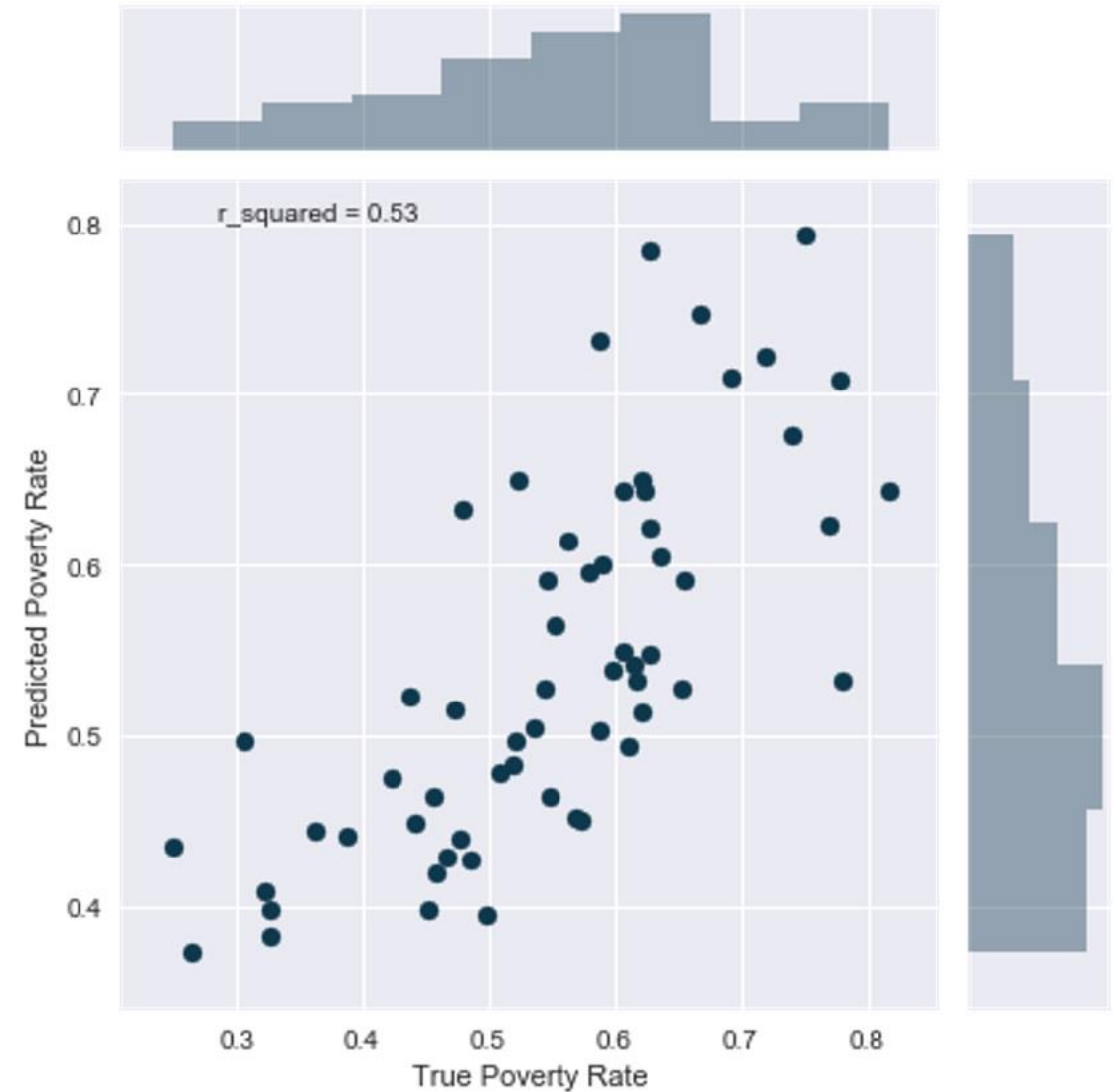
## Be honest about the limitations of your methods



Real world problem

# Urban Municipalities: Predicted vs True

# Tip 4:

## But always show the upside if it works out

# "Micro-Satellites" ~ Daily Revisit Rate

Every Road in Mexico

Legend

Fractional Road Coverage
- 0
- 0.02
- 0.04
- 0.06
- 0.08 +
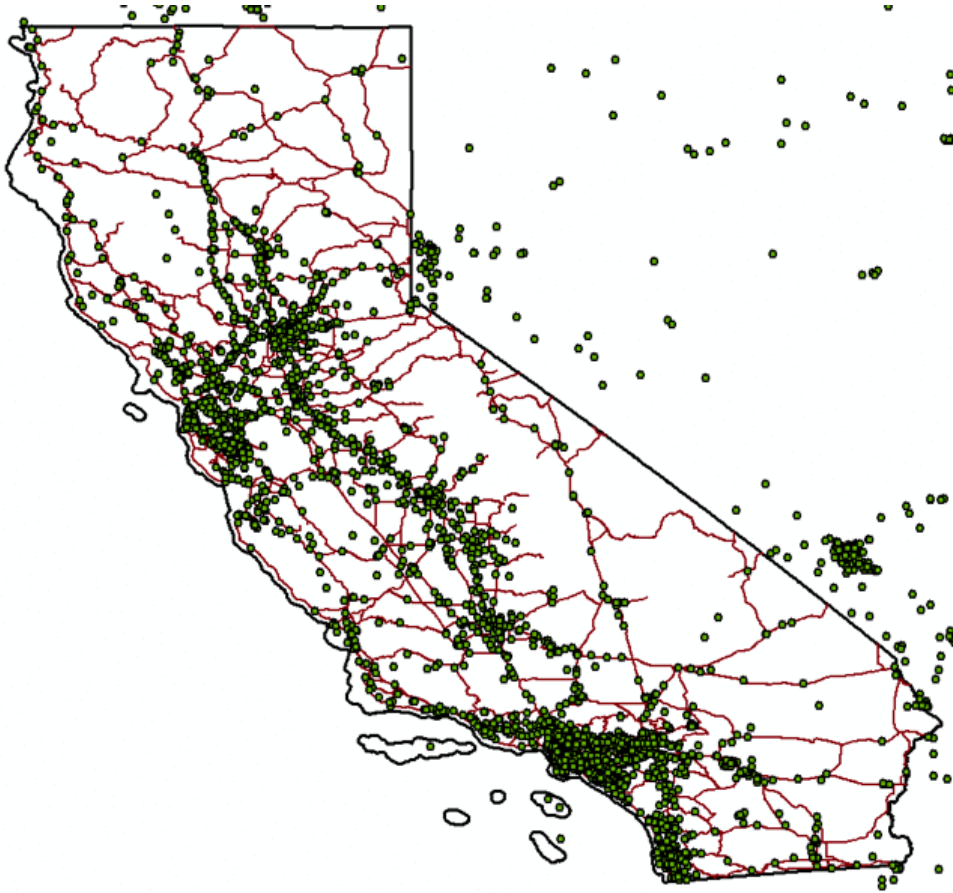
Monterrey

Guadalajara

Mexico City

# Shifting gears: do mobile phones cause car accidents?

- Research Question: do internet enabled mobile phones increase traffic accidents?

- Joint work with Matt and Bree Lang at UC Riverside – "Car Accidents and 3G Coverage: New Evidence Using Cell Phone Towers"
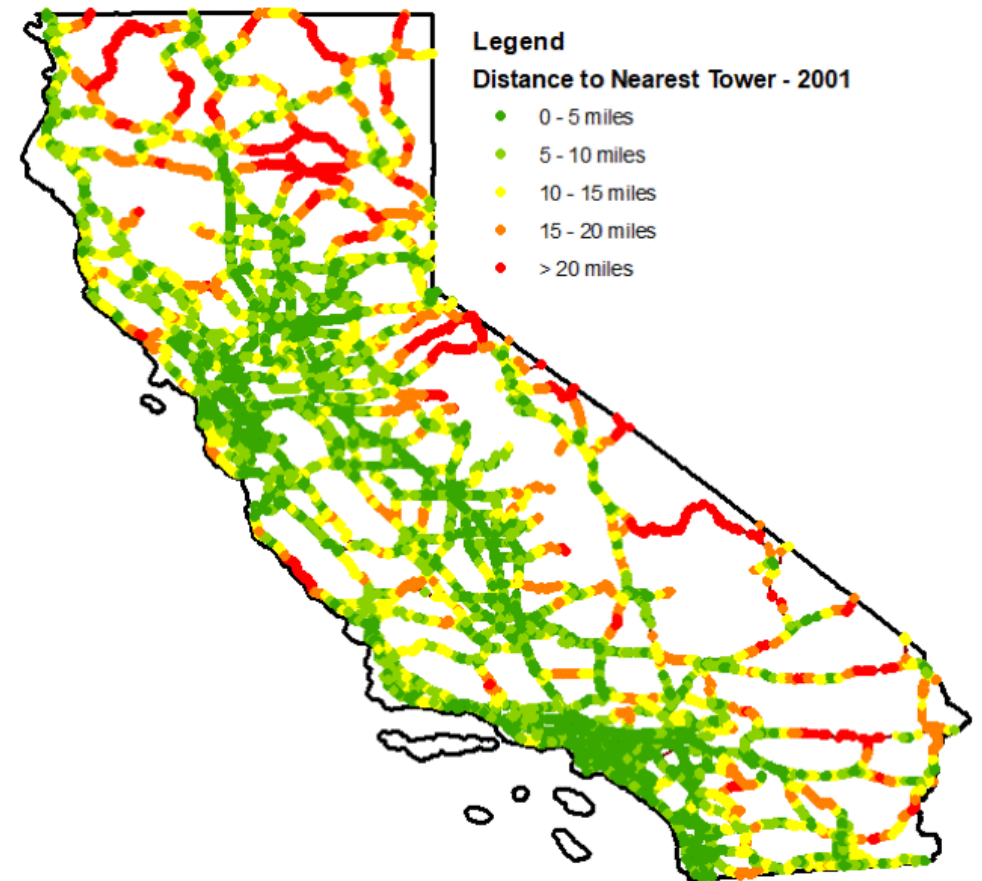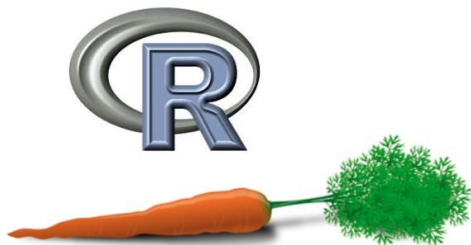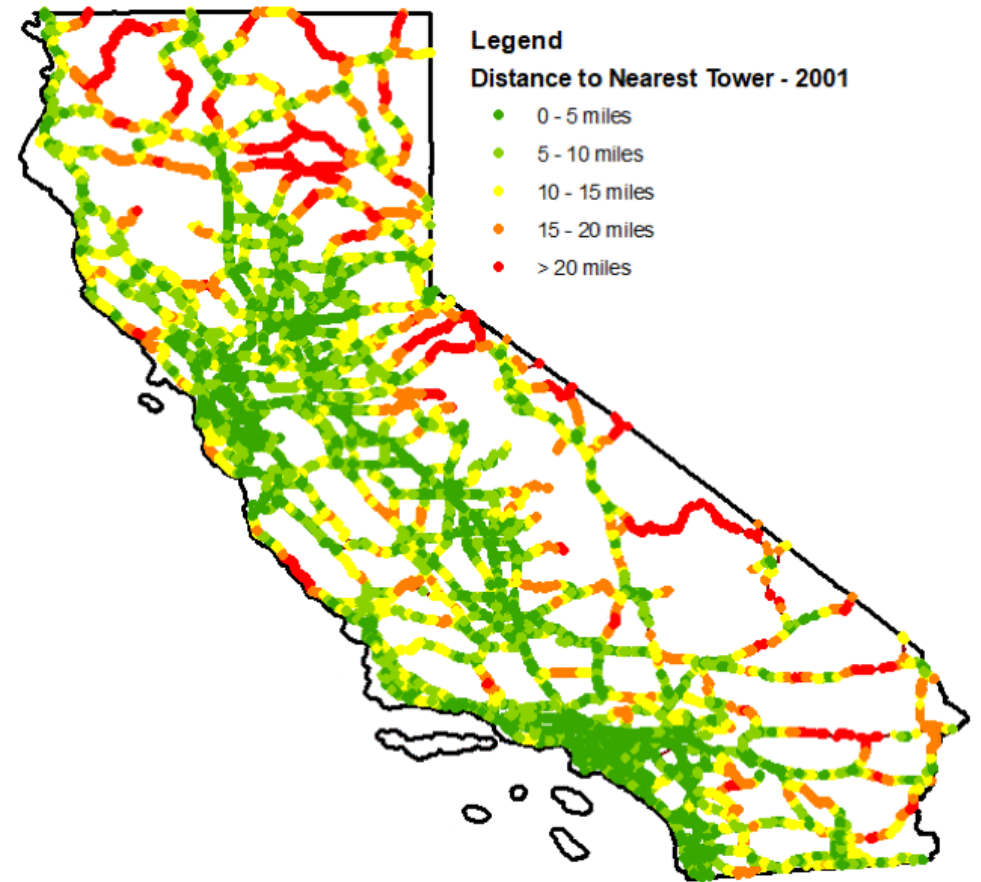
# Event study: growth of 3G

### Cell Tower Growth
### 2001 - 2011



### Road distance to
### nearest tower
### 2001 - 2011



Legend
**Distance to Nearest Tower - 2001**
- 0 - 5 miles
- 5 - 10 miles
- 10 - 15 miles
- 15 - 20 miles
- > 20 miles

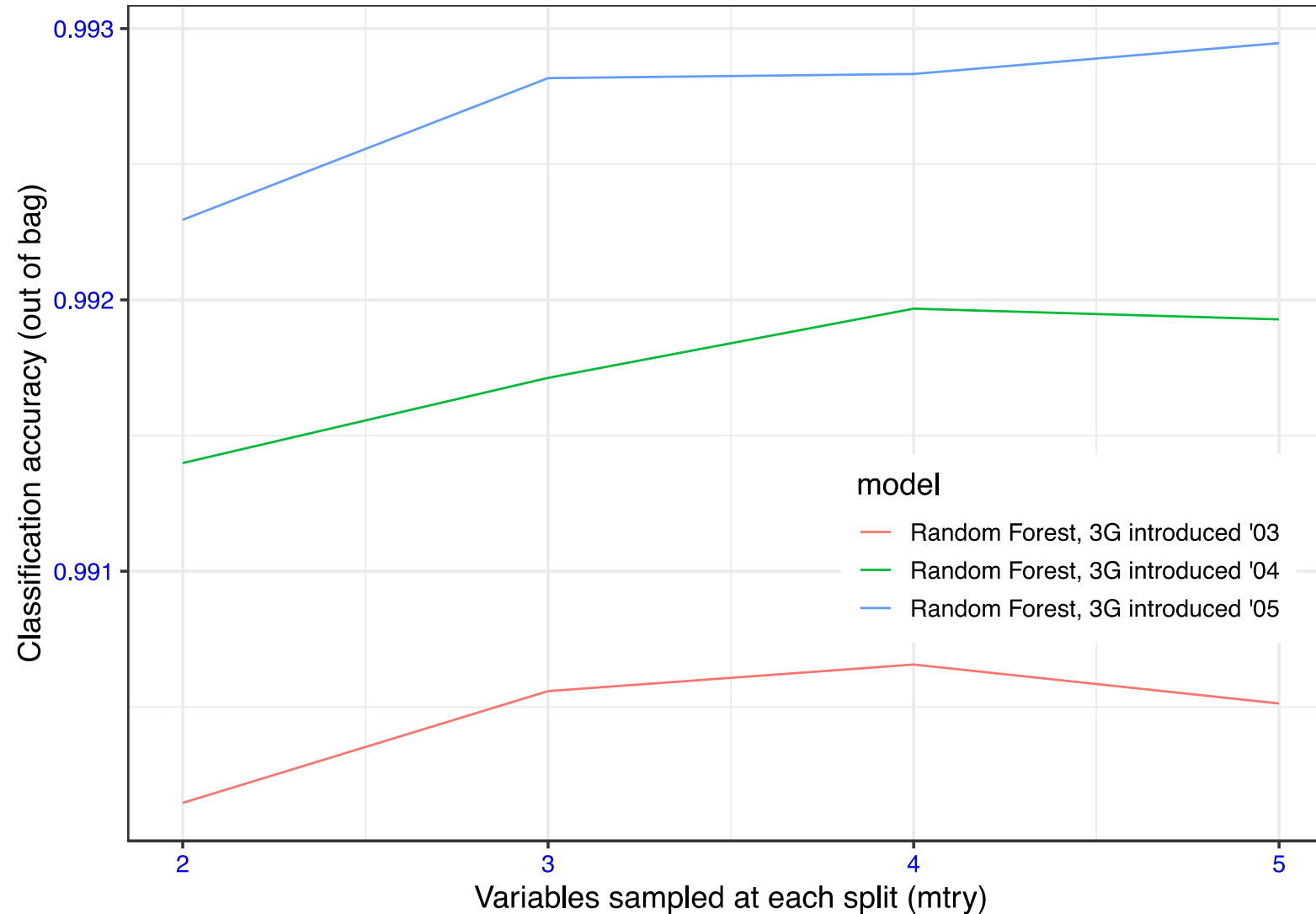# Problem: Only only know if a road has 3G access in 2016

- But: we know closest tower characteristics 2001 – 2016

- Solution: build random forest model to predict 3G coverage 2001 – 2016 based on closest tower characteristics



**Legend**
**Distance to Nearest Tower - 2001**
- 0 - 5 miles
- 5 - 10 miles
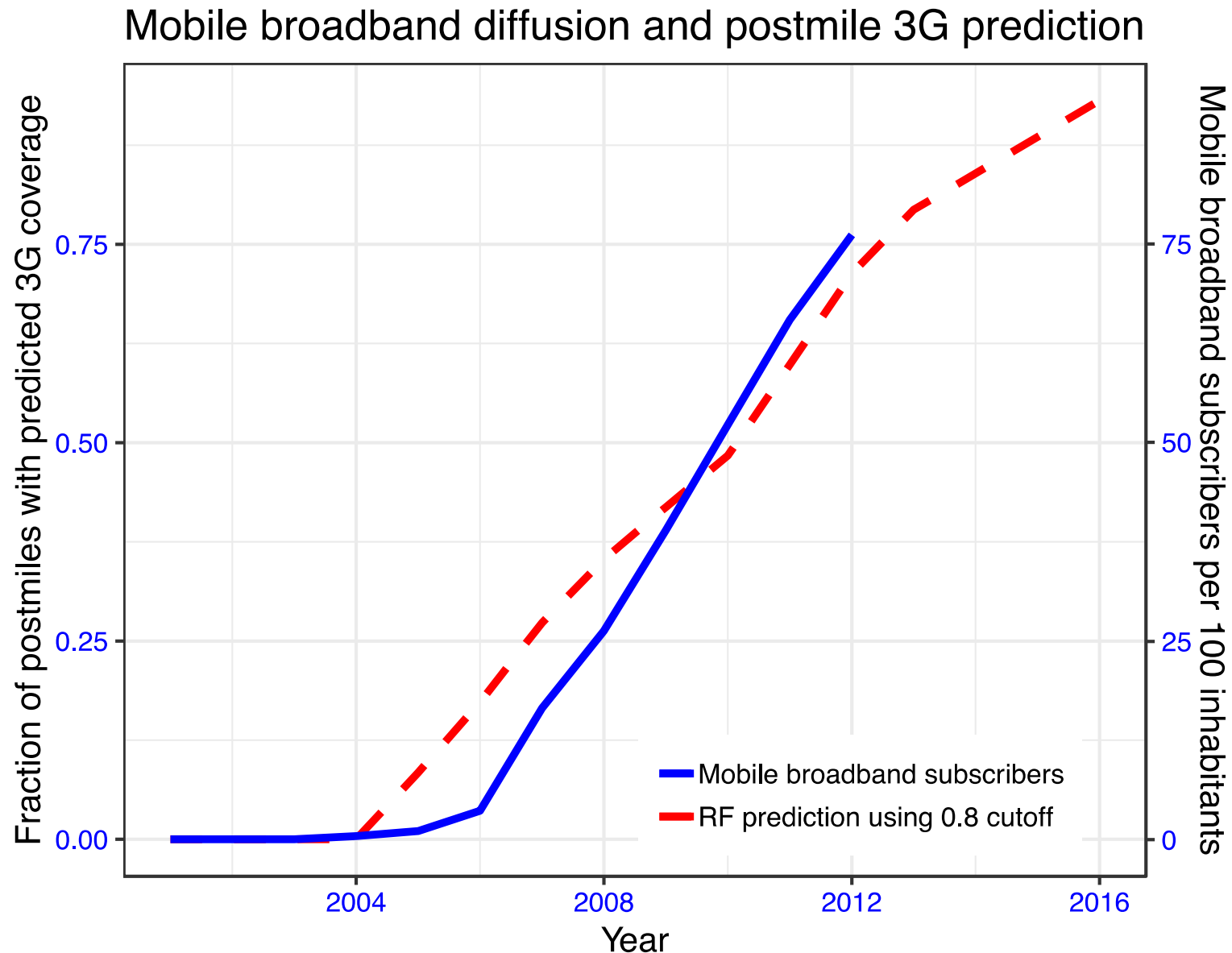- 10 - 15 miles
- 15 - 20 miles
- > 20 miles

# Cross-validate and select optimal 'mtry'



Parameter selection RF model for postmile 3G coverage prediction
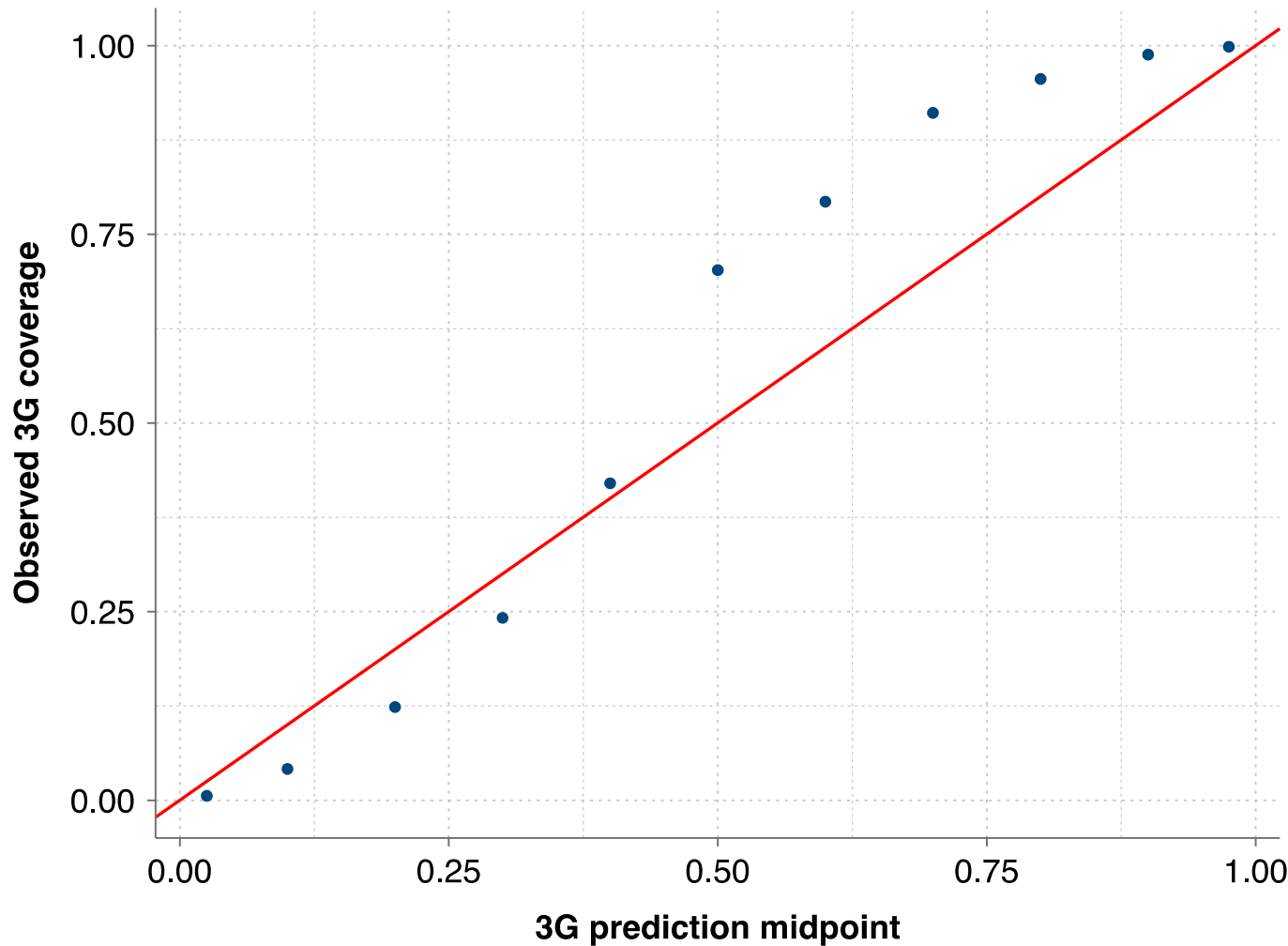Training data, 80% sample. Assuming 3G introduced as shown

# Model Calibration Plot


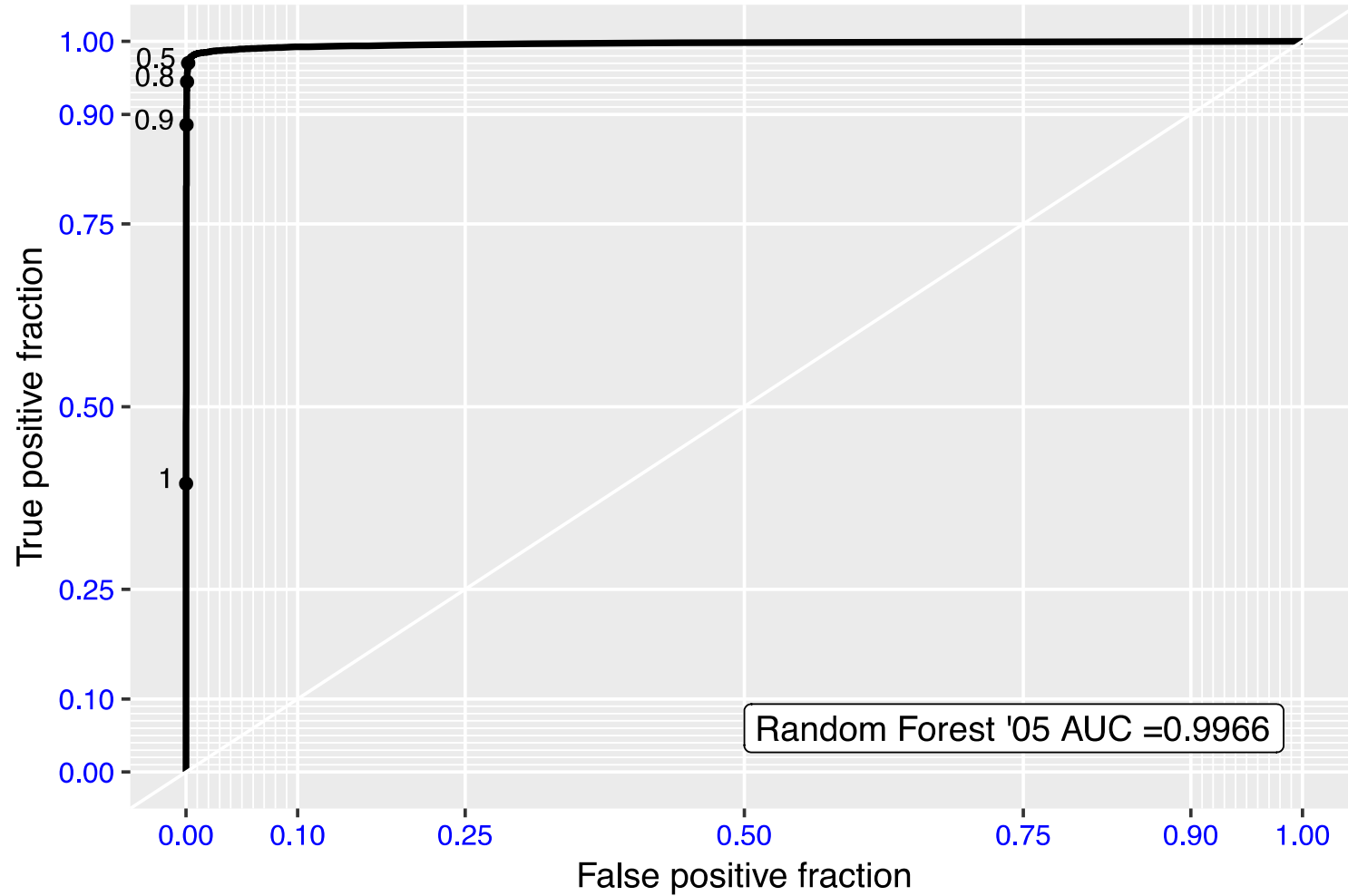
**Observed versus actual road 3G coverage**

Test data, 20% sample

# Test Set: ROC Accuracy Plot



ROC Curve, predicting highway segment 3G coverage

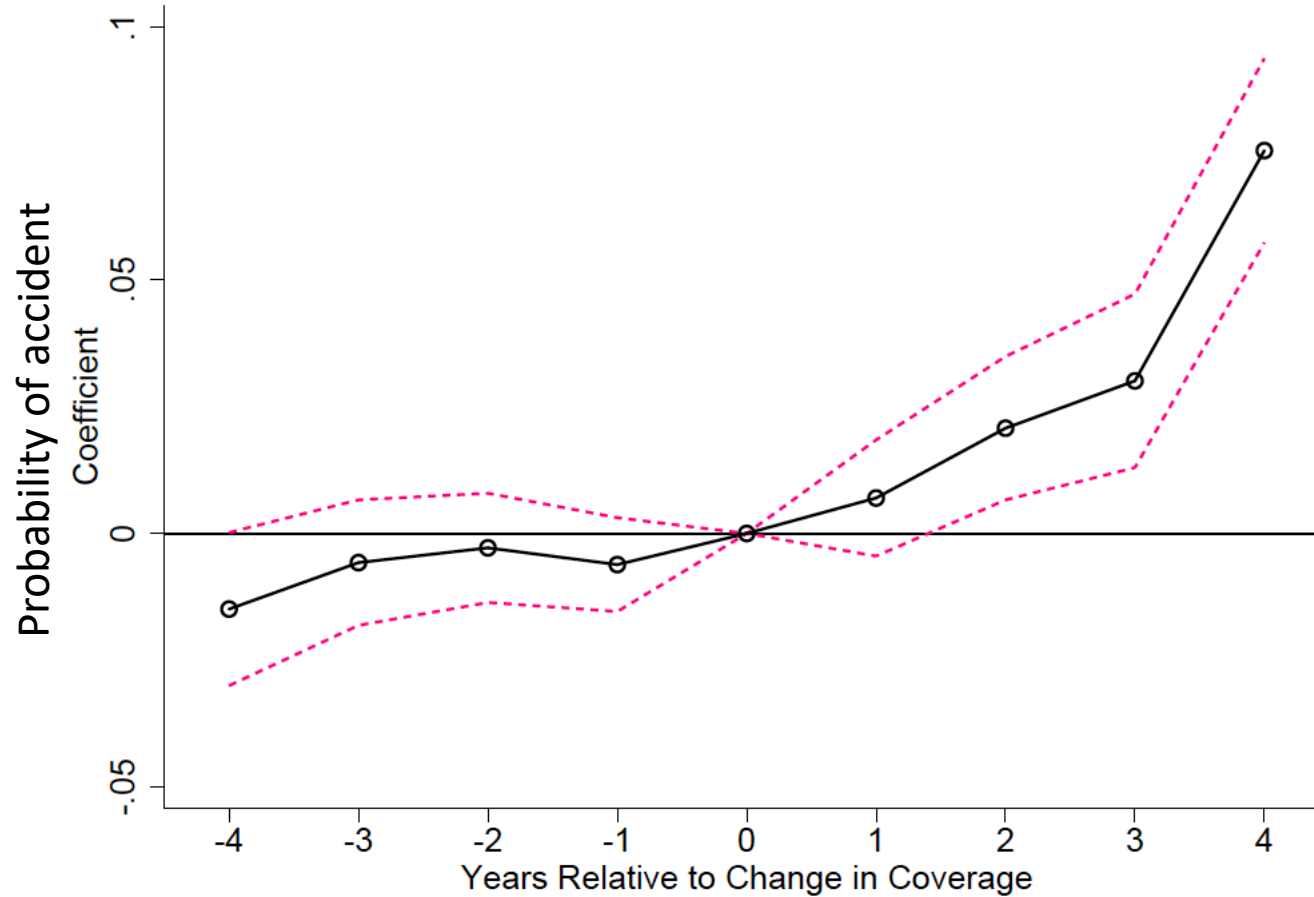Test data, 20% sample (N = 63728). Assumed 3G introduced '05

Random Forest '05 AUC =0.9966

# Event Study Fixed Effect Poisson Model

$$\mathrm{E}[Accident_{it}|X_{it}]$$
$$= \exp\left(\sum_{k=-4}^{-1} \theta_k S_{it+k} + \sum_{j=1}^{4+} \theta_j S_{jt+j} + \ln(Road\ Traffic_{it}) + \gamma_i + \tau_t + v_{it}\right)$$

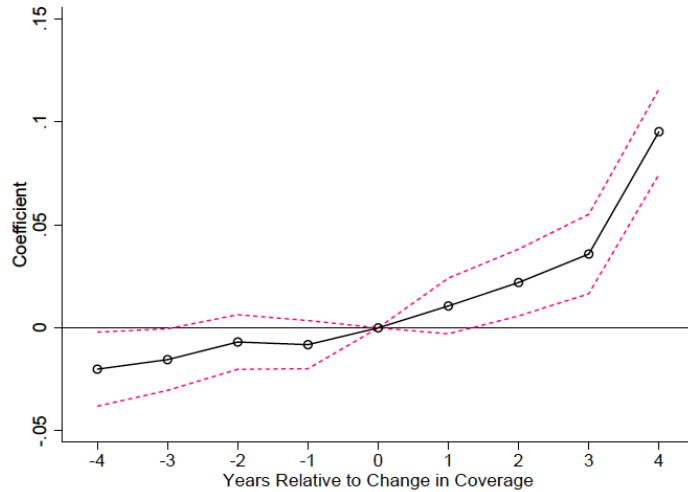- $\theta_j$ impact of mobile internet coverage on accident probability
- $S_{jt}$ when predicted 3G coverage = 1
- $\gamma_i$ road segment fixed effect
- $\tau_t$ time fixed effect
- $v_{it}$ error

# Event study: impact of predict 3G road coverage on traffic accident probability
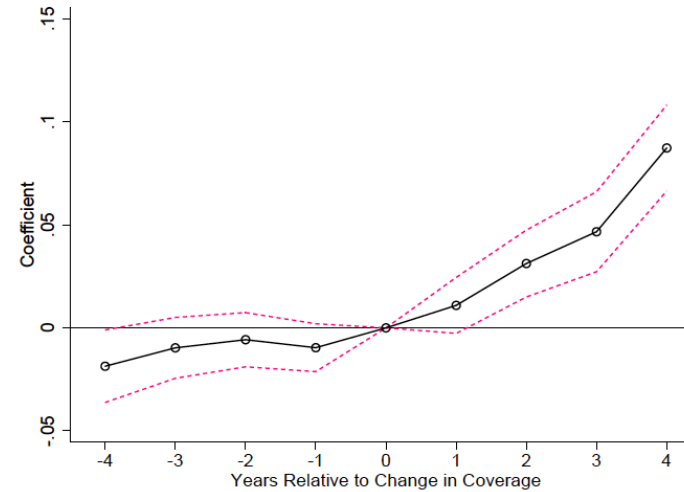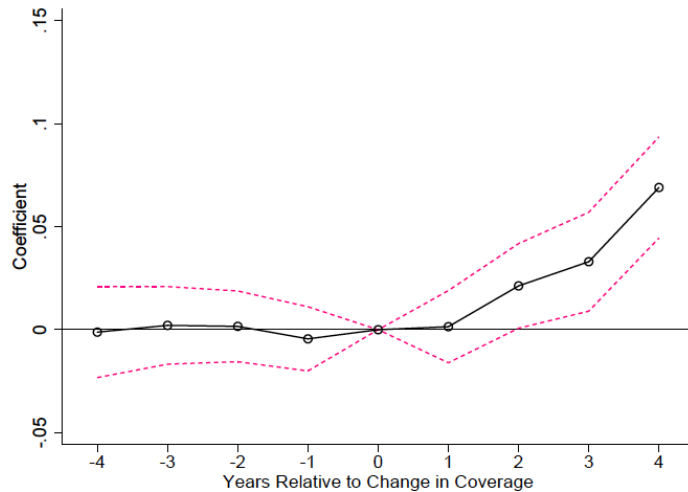


(c) ≥0.80 Threshold, N=40,877

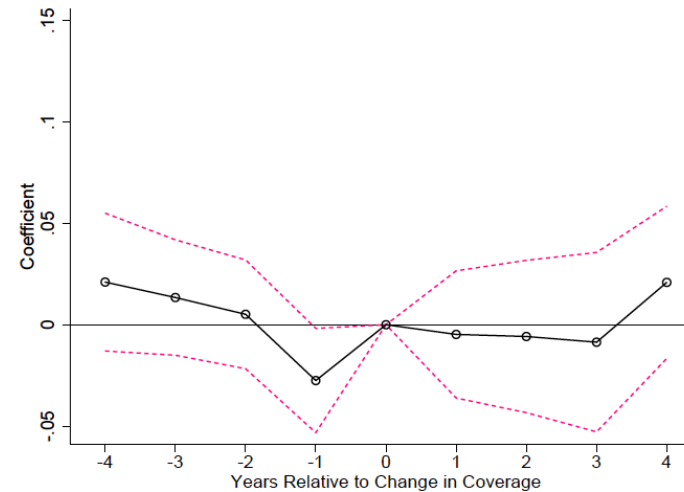# Younger Drivers More Affected by 3G road Access
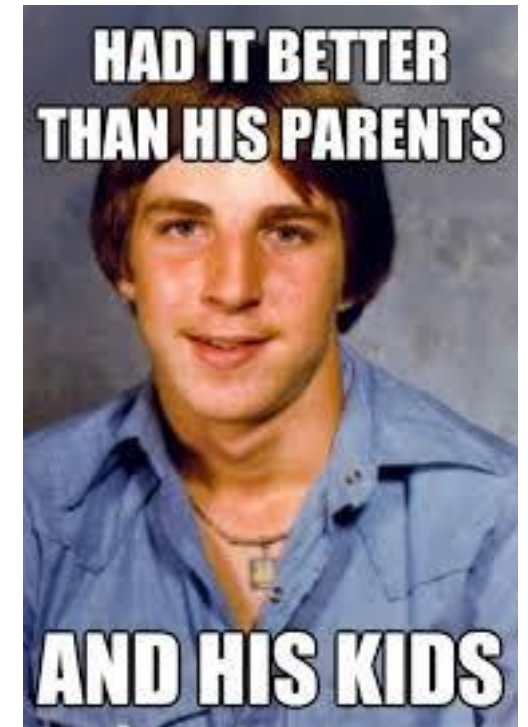


(a) 29 Years and Younger, N=36,880

(b) Between 30 and 49 Years, N=37,377

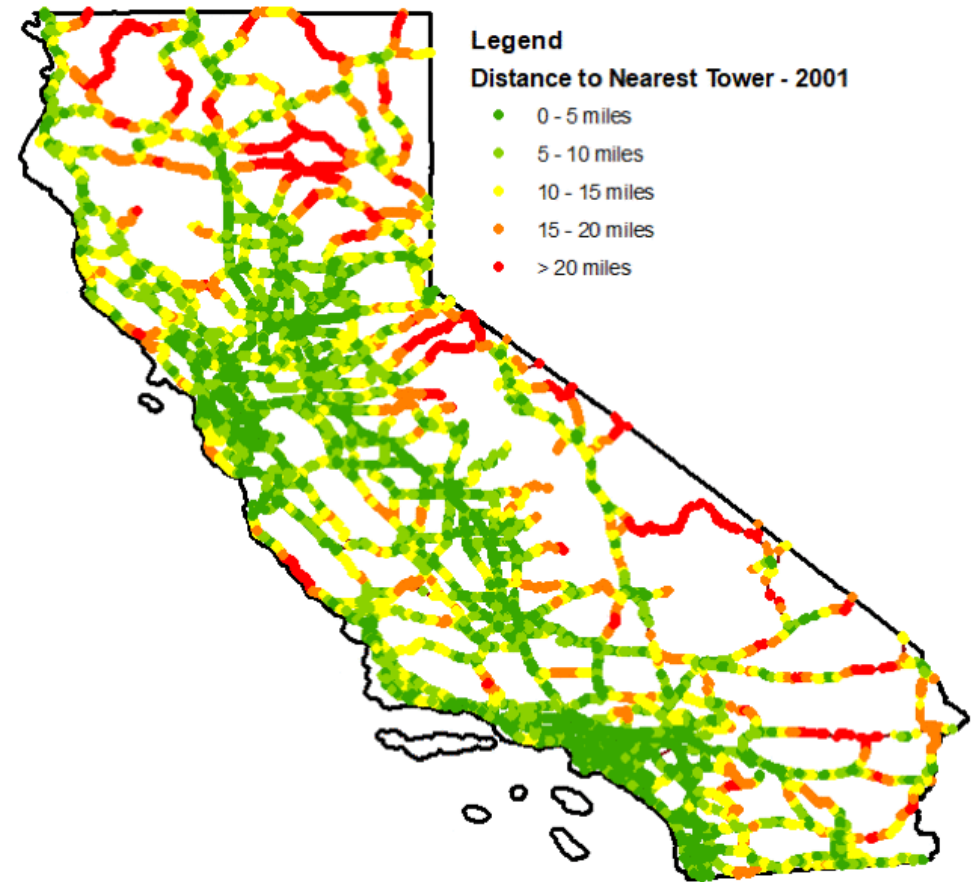(c) Between 50 and 64 Years, N=33,857

(d) 65 Years and Older, N=27,494

# Conclusion

- Accident rates increase 1.1 percent a road gets access to 3G coverage

- Internet connected mobile phones cause 3,305 accidents per year in California

- Further evidence you can embed machine learning in causal inference models



Legend
**Distance to Nearest Tower - 2001**
- 0 - 5 miles
- 5 - 10 miles
- 10 - 15 miles
- 15 - 20 miles
- > 20 miles

# Comments/suggestions appreciated!

Jonathan Hersh

Argyros School of Business, Chapman University

hersh@chapman.edu

(Please talk to me if you're interested in teaching machine learning @ Chapman)