# Who is Bogus? Using One-Sided Labels to Identify Fraudulent Firms from Tax Returns

**Shekhar Mittal**[1]    **Ofir Reich**[2]    **Aprajit Mahajan**[3]

[1]UCLA Anderson (GEM)

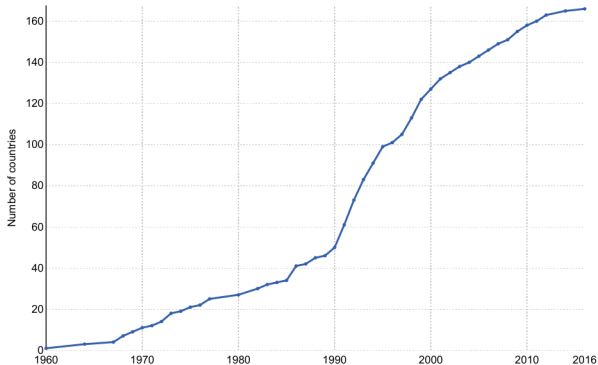[2]CEGA

[3]UC Berkeley (Dept. of ARE) & CEGA

June 21, 2018

ACM COMPASS 2018: Fraud & Security

# Rapid Increase in Value Added Tax Adoption Since 1960



Number of countries having implemented Value Added Taxes, 1960 to 2016

Source: OECD – Consumption Tax Trends 2016

OurWorldInData.org/taxation/ • CC BY-SA

- ▶ 1 country in 1960 → 50 in 1990 → 160 in 2015.
- ▶ Tax levied at each stage of production or distribution (contra sales tax).

# Evasion under VAT: Bogus Firms

- VAT requires buyer & seller to independently report each transaction.
  - Opposing incentives should reduce scope for collusion and evasion.
  - Whether this occurs, particularly in emerging economies, is an open question. (Limited) evidence.
    - In related work, we show that digitization enabled cross-checking of buyer and seller reports increased collections but only from better monitored firms.[1]
  - VAT example
- Alternative evasion strategy – "Bogus" firms.
  - Bogus firms are shell firms created to enable firms to lower tax bills.
  - Create (fake) paper trails of transactions with genuine firms.
  - Bogus firms example
- Precise extent and magnitudes largely unknown.
  - Media reports estimate the loss, in Delhi alone $\approx$ \$300m.[2]
- Commonly reported in many VAT systems.
  - Still relevant for the newly launched Goods and Services Tax (India).
  - Early conversations in Mexico, Dominican Republic, and Zambia.

---

[1]Mittal and Mahajan (2017)
[2]India today article, TOI article, BS article

# Detecting Bogus Firms: Current Practice

- Physical inspections gold standard, but resource intensive.
    - Audit resources limited (particularly in low-income countries).
- Key problem: How to identify firms for inspection?
    - More of them with less effort
- Officials in the central office create a list of "risky" firms.
    - Based on (limited set of) variables: low (VAT deposited/turnover), high turnover, high revisions, invalid address.
- Local inspectors sent out for inspections.
    - Firms deregistered ("cancelled") if inspection fails.

## Our Work
We apply a random forest classifier to the value added tax (VAT) returns from Delhi (India) to increase tax compliance by identifying bogus firms which can be further targeted for physical inspections.

# Highlights

- One sided labels and in-sample predictions ⇒ Cross-validation.
- Precision, recall, F1 score not ideal ⇒ Focus on top recommendations.
- Non-RCT evaluation (for now) ⇒ Point-in-time simulation.
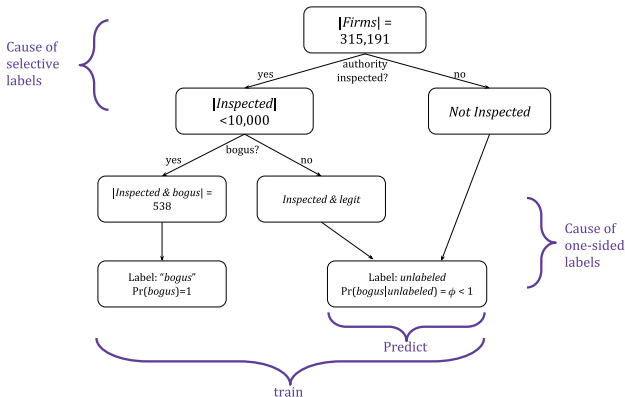- Multiple firm-quarter observations but class timeless ⇒ Aggregate the predictions.

# Delhi and Taxation

- Population: 16.8 million (2011 census)
- Real GSDP of NCT of Delhi for 2015-16: ₹4,560 billion (US$ 71 billion)
  - Tax to GSDP ratio: 5.7%
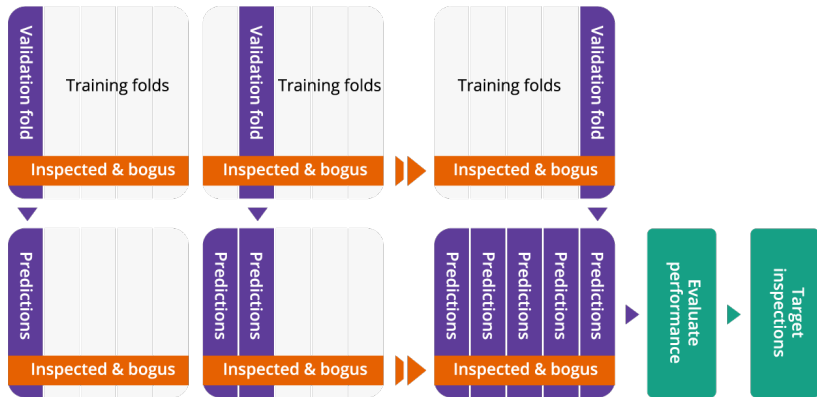- VAT accounts for 52.4% of total government revenues

# Data Description

- Entire universe of registered firms from Delhi, India (315,191 firms, unbalanced).
- 3 years of quarterly VAT returns - 2012-13, 2013-14, 2014-15.
- 3 years of quarterly firm level interactions - 2012-13, 2013-14, 2014-15.
- Firm profile data.
- Bogus firm data: 531 bogus firms identified (2012-2015).

- ▶ Inspection is based on the tax authority's discretion and so biased (selective labels).
- ▶ Class labels are known only for firms both inspected and found to be bogus, not for the rest (one-sided labels).
- ▶ We use all for training, but want to predict for those firms still unlabeled (in-sample predictions).
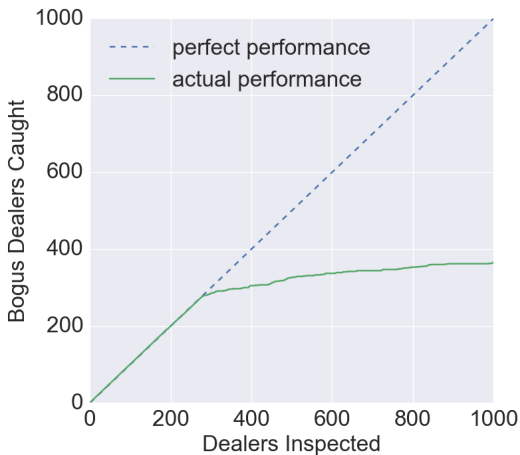
# Cross-validated Prediction Procedure



Multi-Period Model          Wide Model

# Random Forest Model Performance on Top 1000 Recommendations



▶ Results similar when we control for revenue size.

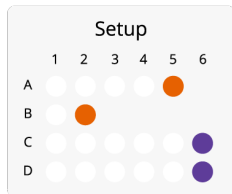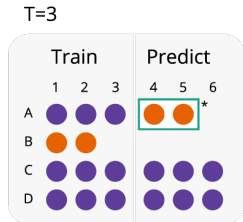Different Classifiers    Different Feature Sets    All Recommendations    Important Features
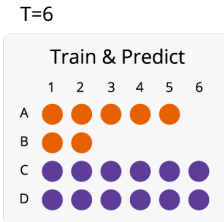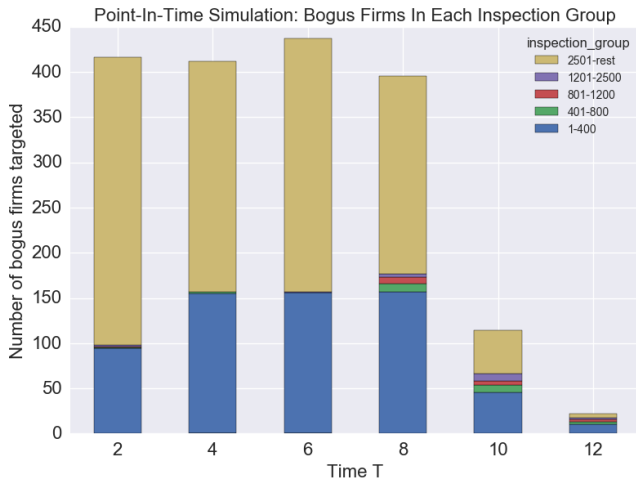
**Point-in-time simulation**

Roll back observations to the state of knowledge at time period T

*Potential revenue saved

- In a real-world scenario, do not have access to all returns by all firms and not required to predict retroactively.
- Blind the model to information obtained after time "T" - do not consider "future" tax returns after T.
- Only use bogus firms that were already classified by time T.

Point-In-Time Simulation: Bogus Firms In Each Inspection Group

Revenue Saved

# Conclusions and Challenges

- Used digitized tax returns to create a ML tool to identify potentially fraudulent firms.
  - Next: Tax authority inspects most suspicious firms (create training data).
  - Finally: Compare revenue implications against current practices.
- Challenge: Firms will respond to better targeting – e.g. by creating more bogus firms faster.
  - ML tool will require regular updating (more training data).
  - Real world example of adversarial ML.
- Interest from many tax authorities, potentially useful tool in the hands of high level officials.

# Thanks!

- We thank GoNCTD, IGC, CEGA, EDI, and JPAL for support.
- This project was funded with UK Aid from the UK Government.
- Starting as a PostDoc at Berkeley School of Information (Josh Blumenstock).

Shekhar Mittal
shekhar.mittal.1@anderson.ucla.edu
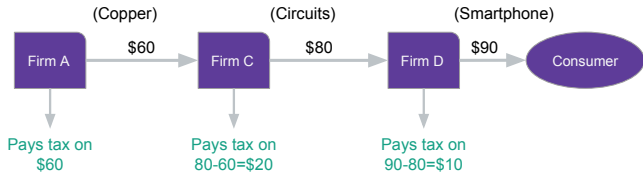UCLA

Ofir Reich
CEGA

Aprajit Mahajan
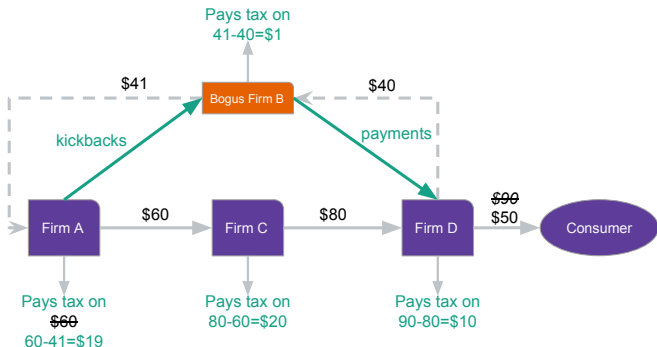aprajit@berkeley.edu
UC Berkeley & CEGA

FUNDED BY

UKaid
from the British people

IGC
International
Growth Centre

ECONOMIC
DEVELOPMENT
& INSTITUTIONS

# VAT: Example



(Copper)    (Circuits)    (Smartphone)

Firm A → $60 → Firm C → $80 → Firm D → $90 → Consumer

Pays tax on $60

Pays tax on 80-60=$20

Pays tax on 90-80=$10

**Government receives tax on $90 value added.**

Back

# Bogus Firms: Example



Pays tax on
41-40=$1

$41
Bogus Firm B
$40

kickbacks
payments

Firm A — $60 → Firm C — $80 → Firm D — ~~$90~~ $50 → Consumer

Pays tax on
~~$60~~
60-41=$19

Pays tax on
80-60=$20

Pays tax on
90-80=$10

**Government receives tax on $50 value added. Surplus is divided between offenders.**

▶ Firms A, C and D not necessarily in the same chain.

▶ Bogus firm can make sales to any firm which needs input credits.

Back

Back

# Wide Model
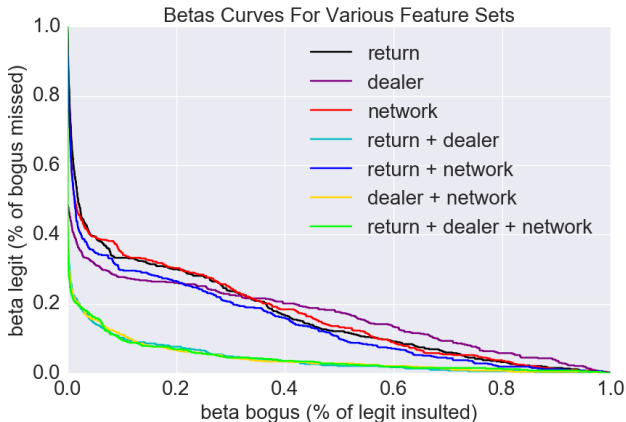


▶ Entry and exit of firms will result in the dataset having a lot of NULL values

Back

# Comparison of Different Classifiers



Different Classifiers: Bogus Firms In Each Inspection Group

Back

Betas Curves For Various Feature Sets

# Model performance on all recommendations

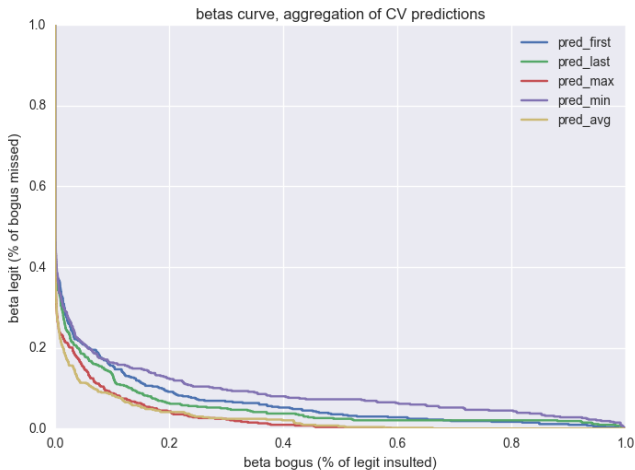| Inspection Group | Firms Inspected | Total Bogus Firms Caught | Bogus Firms Caught/Inspection |
|---|---|---|---|
| 1 - 400 | 400 | 305 | 0.76 |
| 401 - 800 | 400 | 48 | 0.12 |
| 801 - 1200 | 400 | 24 | 0.06 |
| 1201 - 2500 | 1300 | 29 | 0.02 |
| 2501 - rest | 313229 | 132 | 0.00 |

Back

# Point-in-time simulation performance for the 1-400 inspection group

| T | Total Bogus Firms Caught | Bogus Firms Caught/Inspection | Revenue Gained by Inspecting Entire Group (USD Millions) | Revenue Gained per Inspection (USD 000s) | Total Bogus Firms in the Sample | Revenue Lost from All Bogus Firms (USD Millions) |
|---|---|---|---|---|---|---|
| 2 | 94 | 0.24 | 19.44 | 48.60 | 416 | 49.40 |
| 4 | 155 | 0.39 | 43.19 | 107.97 | 412 | 108.38 |
| 6 | 156 | 0.39 | 25.48 | 63.70 | 437 | 63.84 |
| 8 | 157 | 0.39 | 9.38 | 23.46 | 395 | 26.43 |
| 10 | 46 | 0.11 | 1.70 | 4.24 | 114 | 4.52 |
| 12 | 10 | 0.02 | 0 | 0 | 22 | 0 |

Back

betas curve, aggregation of CV predictions

# Interpreting Features: Gaming Measures



Ratio of money deposited to turnover

Firms grouped in deciles(10%) of ratio of money deposited to turnover
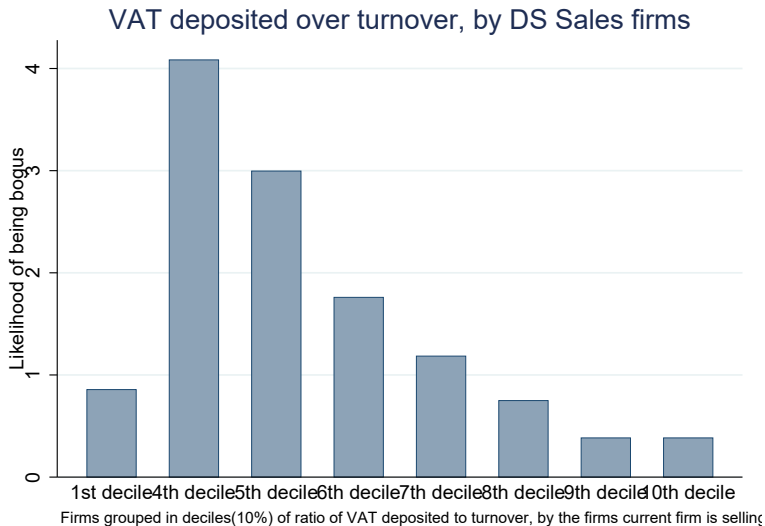
▶ Bogus firms likely to have ratio in middle indicates that they know tax authority monitors extreme values so they make sure they are not in extremes.
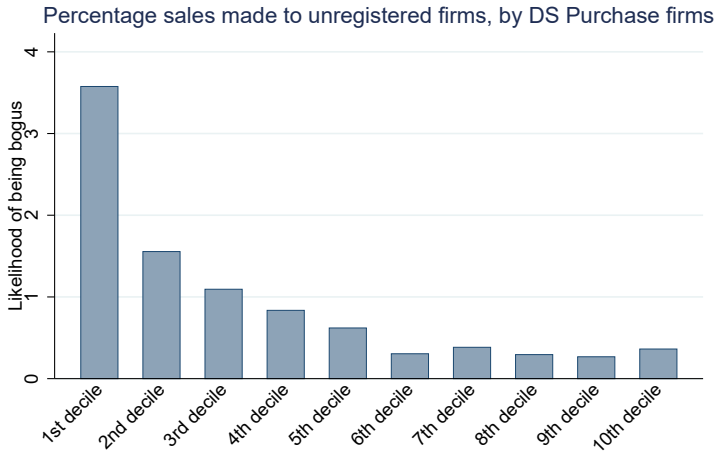
Back to Results

Pagerank (purchases)

Firms grouped in deciles(10%) of pagerank from the purchased from data (2A)

Back to Results

## VAT deposited over turnover, by DS Sales firms



Firms grouped in deciles(10%) of ratio of VAT deposited to turnover, by the firms current firm is selling
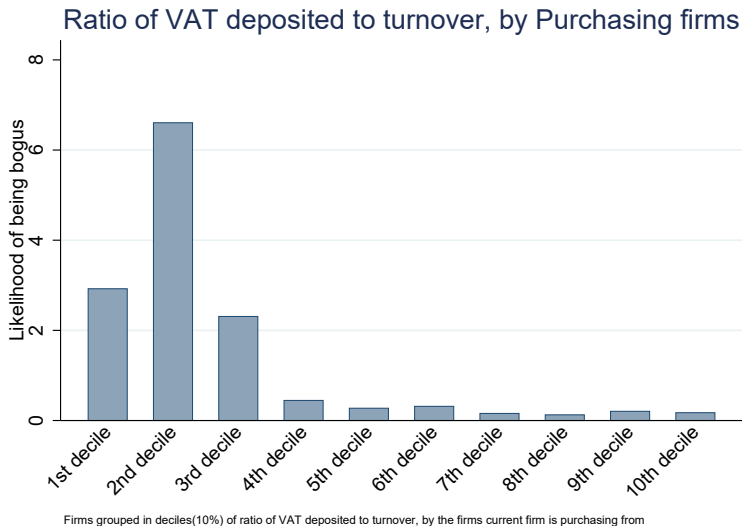
Back to Results

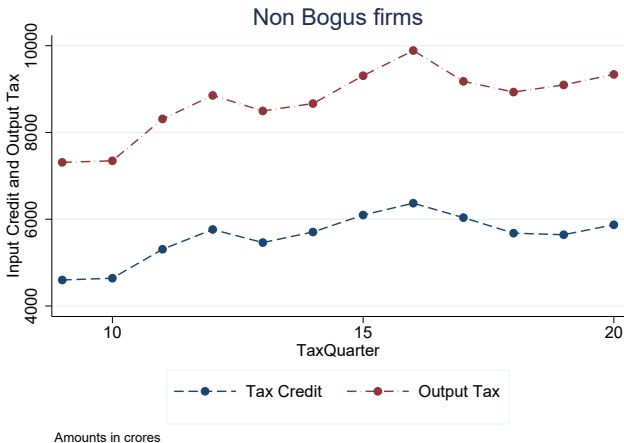Percentage sales made to unregistered firms, by DS Purchase firms

Firms grouped in deciles(10%) of percentage sales made to unregistered firms,
By the firms current firm is purchasing from

Ratio of VAT deposited to turnover, by Purchasing firms
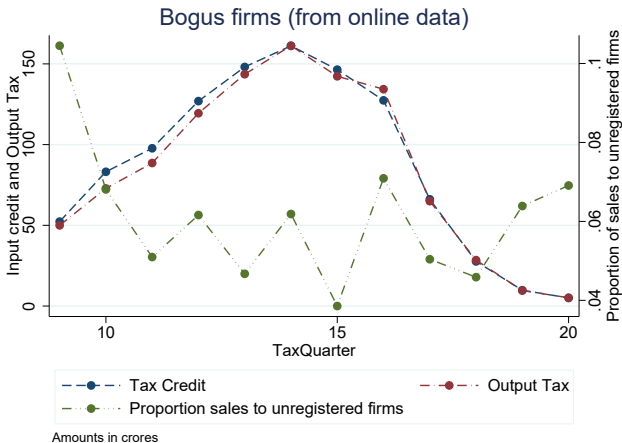
Firms grouped in deciles(10%) of ratio of VAT deposited to turnover, by the firms current firm is purchasing from

# How genuine firms look



Non Bogus firms

Amounts in crores

▶ Total output tax reliably larger than input tax credit.

# Size of problem: From explicit data



Bogus firms (from online data)

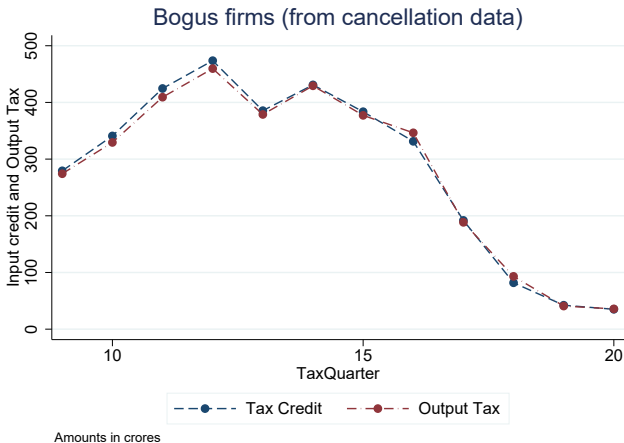Amounts in crores

- ▶ Input credit claimed weakly greater than output tax declared
- ▶ From the limited sample, revenue loss between ₹4-6 billion, annually
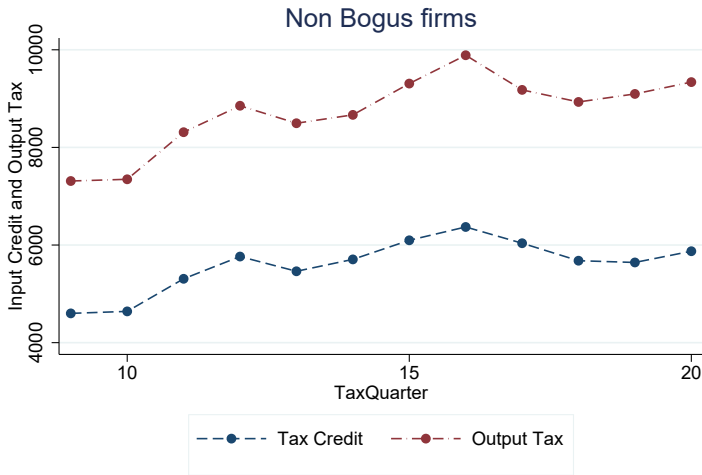- ▶ Drop in later quarters due to missing data

From cancellation records

Bogus firms (from cancellation data)

Amounts in crores

- From the much bigger sample, revenue loss around ₹15 billion, annually
- Drop in later quarters due to missing data

Back

Non Bogus firms

Amounts in crores