# Data Quality Monitoring for the 2022 SCF

**Identifying Falsified Data & Introduction of NORC ProofPoint**

04.12.23

Jimmy Herdegen, Research Associate I
Kate Bachtell, Senior Research Director II
Jason Keller, Senior Data Scientist II
Cathy Haggerty, Vice President
Lisa Blumerman, Senior Vice President

NORC at the University of Chicago

# Agenda

# Introduction

# Problem Statement

**Problems Addressed**

• There is no one way falsified data can happen, which poses challenges on how to detect invalid cases

• Creating a framework to filter out falsified data and enhance data quality is critical for large scale surveys like the SCF

**Goal of Presentation**

• By learning lessons from previous cycles, and creating and implementing a set of metrics to filter cases through, we can be smarter and faster in detecting falsification

# Field interviewer falsification comes in many forms

## Survey Format

- Interviewers with invalid data often use filter branches in surveys to intentionally shorten an interview (Walzenbach, 2021)
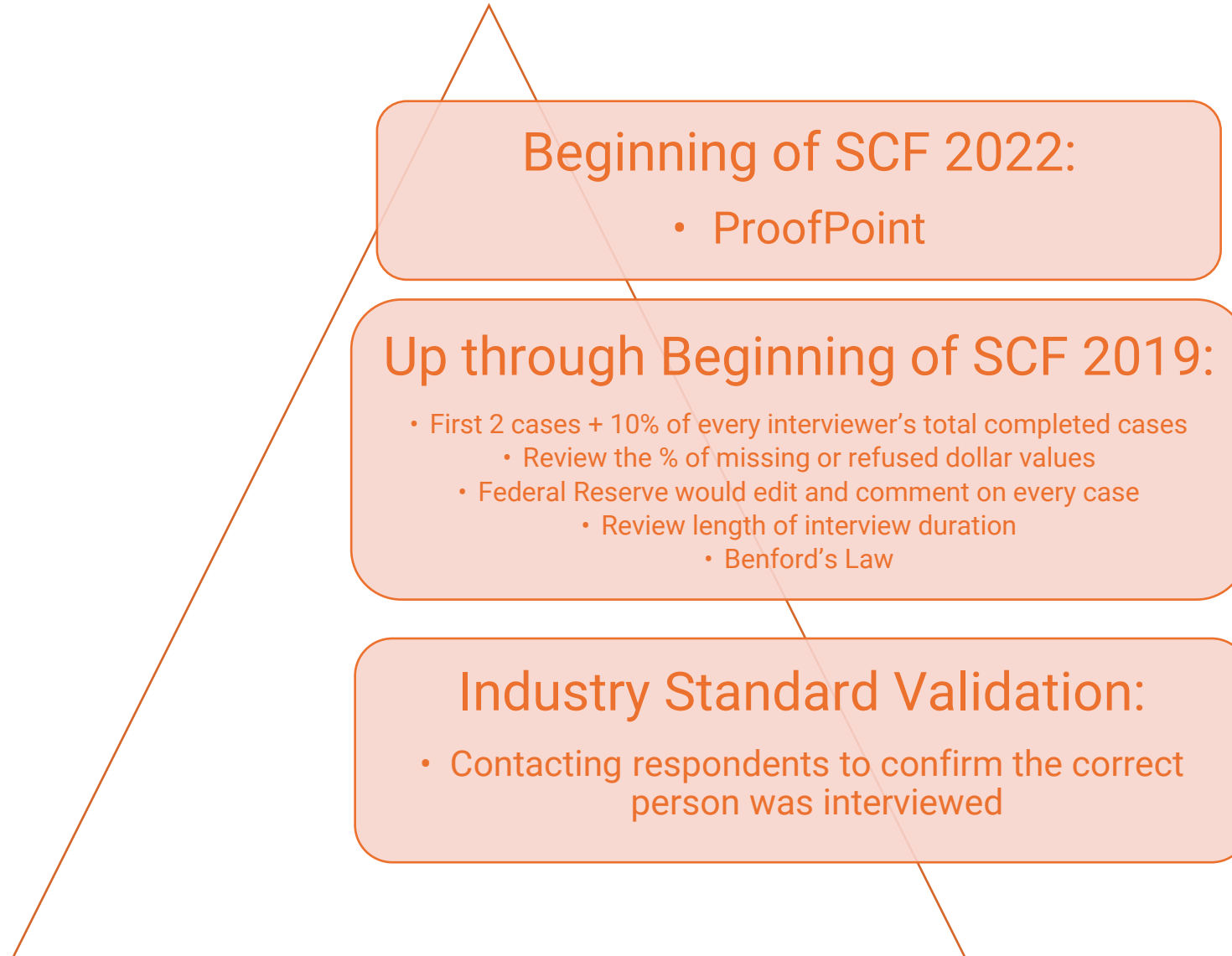- Fewer missing answers (SCF 2019)

## Survey Content

- Certain behavioral & attitudinal variables could predict falsified data (Menold, Kemper, 2014)

## Curb-Stoning

- Fabricating the entire interview at the time of the interview (Thissen, Myers, 2016)
- Omitting interviews by reporting them as refusals or unlocatable, when little effort was done to contact the respondent (Thissen, Myers, 2016)

# Detection Methods at the Beginning of the 2019 Cycle

NORC at the University of Chicago

# Multiple layers to validation procedure

## Beginning of SCF 2022:

- ProofPoint

## Up through Beginning of SCF 2019:

- First 2 cases + 10% of every interviewer's total completed cases
- Review the % of missing or refused dollar values
- Federal Reserve would edit and comment on every case
- Review length of interview duration
- Benford's Law

## Industry Standard Validation:

- Contacting respondents to confirm the correct person was interviewed

# Falsification Problems & Responses During the 2019 Cycle

# We discovered bad actors later in the field period

- Additional measures were implemented during the midpoint of data collection during the 2019 cycle:

  - Paradata Analysis:
    - Missing phone numbers via CAPI and NORCSuite
    - Looking at the record of calls (ROCs)

  - Reviewing GPS coordinates from 2 different apps associated with interviewer activity

  - Electronic signatures for payment receipts

  - CAPI data quality analysis:
    - Quex timings
    - Missing dollar values
    - Data quality flags
    - Feedback from the Federal Reserve

# Characteristics of Falsified Cases During the 2019 Cycle

**Shorter**

- Without genuine human interaction between the field interviewer and respondent, survey times will tend to be shorter

**Less Descriptive**

- Invalid case comments are shorter and not as descriptive

**Lack Contact Info**

- Interviews that were falsified were more likely not to include contact info for respondent

**Fewer Missing Answers**

- Field Interviewers would tend to provide fewer missing responses compared to valid data

# Examples of Respondent Signatures from Valid & Invalid Cases

**Example of Signatures from Valid Cases**

**Example of Signatures from Invalid Cases**

# Solution for SCF 2022: ProofPoint

# Interviewer Summary View

Field Manager ID
040531

Interviewer ID ▾

Dashboard | FI Summary ⋮ | FI Detail | Map | +

## FM Validation Summary

| Field Manager ID | Interviewer ID | Number of Cases | Recorded Email Address | Recorded Phone Number ▲ | Duration in Minutes | Dooblo Flag | TSheets Flag |
|---|---|---|---|---|---|---|---|
| 040531 | 011996 | 5 | 60% 🟨 | 80% 🟩 | 146 🟩 | 80% 🟥 | 0% 🟩 |
| 040531 | 011675 | 16 | 6% 🟥 | 81% 🟩 | 119 🟨 | 94% 🟥 | 0% 🟩 |
| 040531 | 012058 | 34 | 35% 🟥 | 85% 🟩 | 99 🟨 | 97% 🟥 | 0% 🟩 |
| 040531 | 002901 | 10 | 70% 🟨 | 90% 🟩 | 139 🟩 | 100% 🟥 | 0% 🟩 |
| 040531 | 002892 | 28 | 25% 🟥 | 93% 🟩 | 199 | 86% 🟥 | 11% 🟩 |
| 040531 | 011439 | 3 | 100% 🟩 | 100% 🟩 | 165 | 100% 🟥 | 0% 🟩 |
| 040531 | 011891 | 1 | 0% 🟥 | 100% 🟩 | 89 🟥 | 100% 🟥 | 0% 🟩 |
| 040531 | 010667 | 12 | 42% 🟥 | 100% 🟩 | 119 🟨 | 100% 🟥 | 58% 🟥 |
| 040531 | 011882 | 10 | 80% 🟩 | 100% 🟩 | 111 🟨 | 70% 🟥 | 0% 🟩 |
| 040531 | 010421 | 6 | 100% 🟩 | 100% 🟩 | 122 🟨 | 67% 🟥 | 0% 🟩 |
| 040531 | 011673 | 3 | 0% 🟥 | 100% 🟩 | 241 | 100% 🟥 | 0% 🟩 |
| 040531 | 011915 | 4 | 0% 🟥 | 100% 🟩 | 140 🟩 | 75% 🟥 | 0% 🟩 |
| 040531 | 011989 | 6 | 50% 🟥 | 100% 🟩 | 212 | 100% 🟥 | 17% 🟩 |
| 040531 | 012061 | 9 | 22% 🟥 | 100% 🟩 | 106 🟨 | 100% 🟥 | 0% 🟩 |
| 040531 | 011524 | 14 | 50% 🟥 | 100% 🟩 | 131 🟩 | 93% 🟥 | 36% 🟨 |
| 040531 | 009105 | 31 | 10% 🟥 | 100% 🟩 | 72 🟥 | 100% 🟥 | 23% 🟩 |
| 040531 | 012063 | 3 | 67% 🟨 | 100% 🟩 | 107 🟨 | 100% 🟥 | 0% 🟩 |
| 040531 | 002874 | 5 | 60% 🟨 | 100% 🟩 | 98 🟨 | 80% 🟥 | 0% 🟩 |

- Summary of each interviewer's quality metrics

- Percentage of critical contact info:
  - Email address
  - Phone number

- Percentage of validation passed:
  - Email
  - Phone
  - Mail

- Average interview duration time

- Average falsification score

# Interviewer Detail

- Displays case-level metrics for every interviewer

- Indicates if email address or phone number were captured

- Interview duration time

- Distance between interviewer and respondent's home

- How a case has been validated

- Case's falsification score

# Falsification Score Calculation
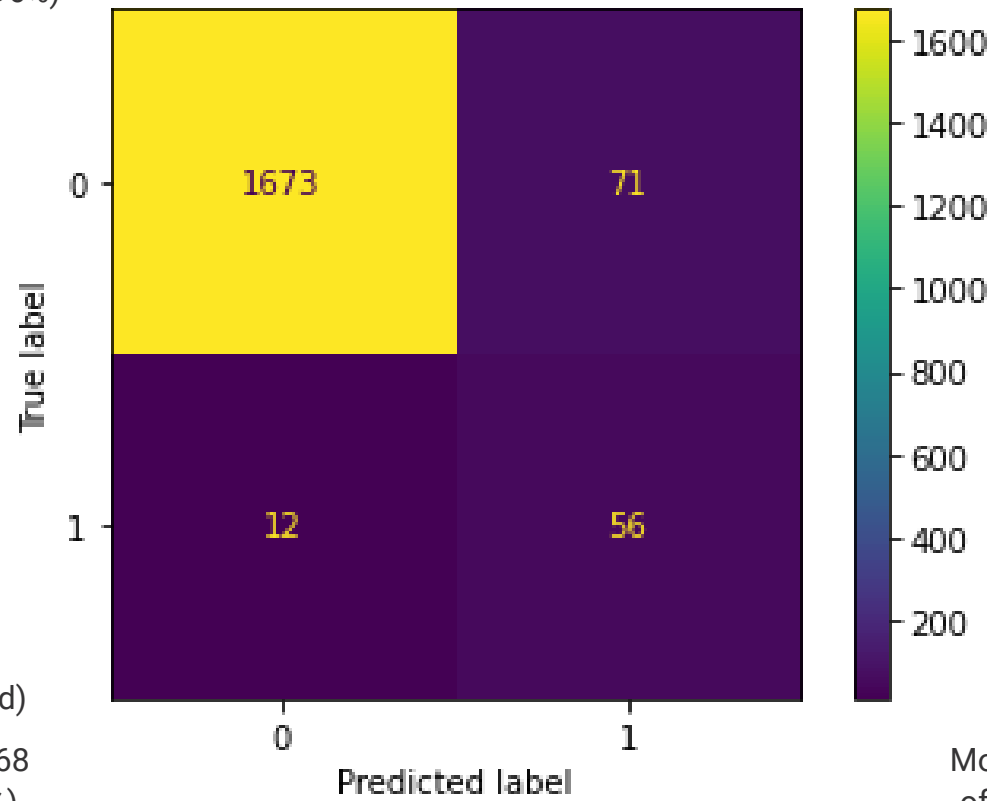
**True Negatives**

(Not falsified and not flagged)

Model correctly identifies 1,673 out of 1,744 valid interviews (96%)

**False positives**

(Not falsified but flagged)

Model incorrectly identifies 71 out of 1,744 valid interviews as falsified (4%)

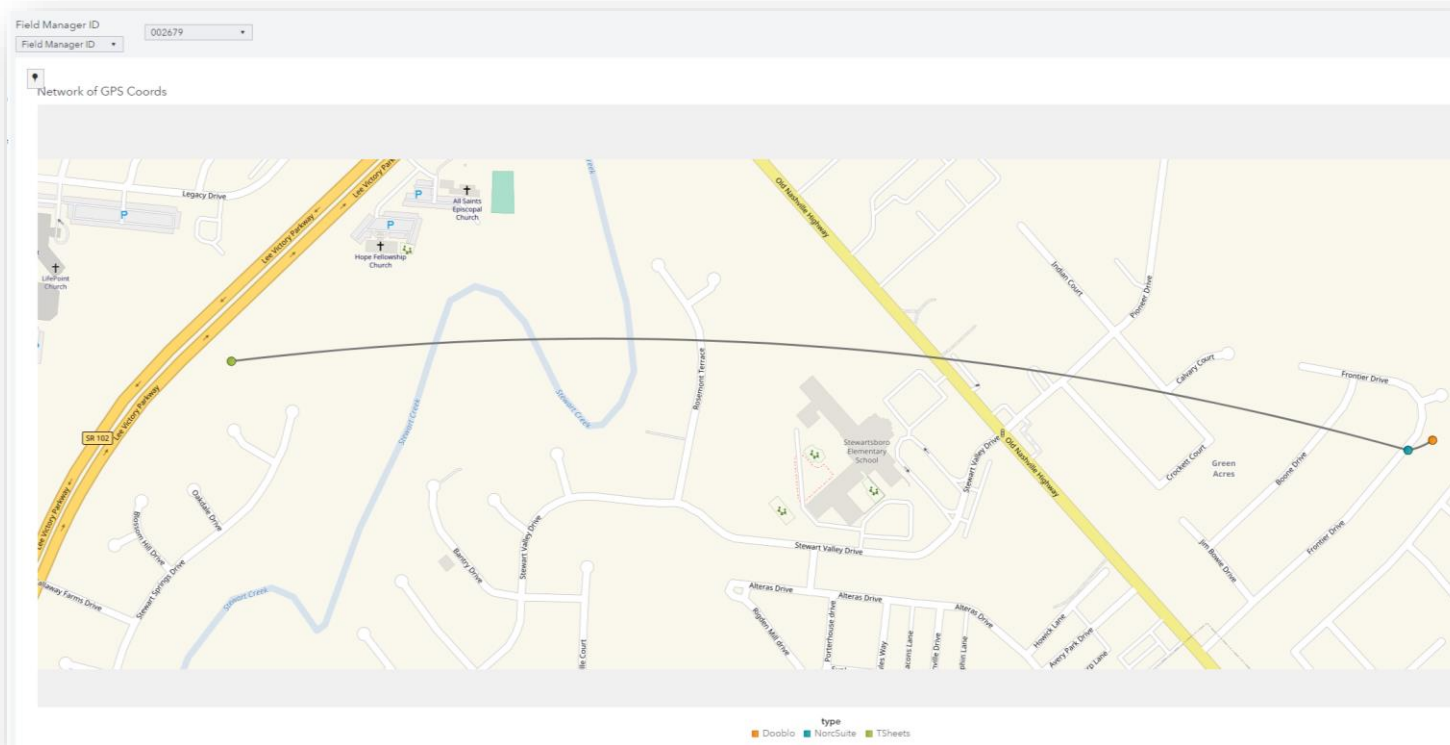**False Negatives**

(Falsified and not flagged)

Model misses 12 out of 68 falsified interviews (18%)

**True Positives**

(Falsified and flagged)

Model correctly identifies 56 out of 68 falsified interviews (82%)

# Proofpoint has been effective in isolating potential falsifiers, though improvements will strengthen it



- Falsifiers were identified sooner compared to 2019
  - Falsification scores were able to identify some falsified cases, though they were unable to capture others

- GPS locating helpful for reviewing in-person cases

- ProofPoint's data complements phone and mail validation efforts

- Data helped enlighten the average for key metrics
  - Email and phone number recorded
  - Survey duration time
  - Phone validation rate

# Looking Ahead

# How can we take full advantage of Proofpoint for the future?

Examine
falsification score

Retrain model

Better Proofpoint
integration

# Thank you.

**Jimmy Herdegen**
Research Associate I
Herdegen-Jimmy@norc.org

**Kate Bachtell**
Senior Research Director II
Bachtell-Kate@norc.org

**Jason Keller**
Senior Data Scientist II
Keller-Jason@norc.org

**Cathy Haggerty**
Vice President
Haggerty-Cathy@norc.org

**Lisa Blumerman**
Senior Vice President
Blumerman-Lisa@norc.org

✳ Research You Can Trust™

✳ NORC at the University of Chicago