# Using Python for Record Linkage: Entrepreneurship, Research and Development, and Lobbying in the Unmanned Aerial Vehicle Industry

Russell J. Funk
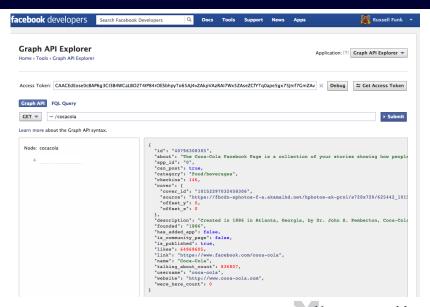funk@umich.edu

May 20, 2013

UNIVERSITY OF MICHIGAN

We'll learn how to. . .

1. pull data from disparate online sources
2. link messy data with python using fuzzy matching
3. use python to build a data set ready for analysis

# Motivation. . .

# Motivation. . .

What is the correlation between lobbying expenditures and research and development contracts for small businesses from Department of Defense?

**Quick background. . .**

- ▶ **Small Business Innovation Research Program (SBIR)**—requires that all federal agencies with extramural research budgets in excess of $100 million reserve 2.5% for contracts or grants to small businesses
- ▶ **Small Business Technology Transfer Program (STTR)**—similar to SBIR, but smaller, and emphasizes funding partnerships between small businesses and nonprofit organizations

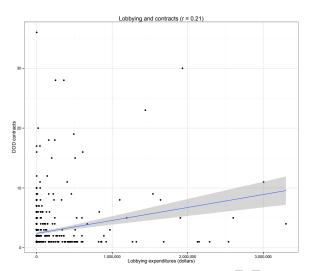Check **README.md** in the **/data** folder for instructions.

How can we link data records across sources without common unique identifiers?

| | |
|---|---|
| SYSTEMS & PROCESSES ENGINEERING CORP. | systems and processes engineering corp |
| SYSTEMS & PROCESSES ENGINEERING CORP. | systems and processes engineering corp |
| SYSTEMS & PROCESSES ENGINEERING CORP.(SP | systems and processes engineering corp |
| SYSTEMS & PROCESSES ENGINEERING CORP.(SP | systems and processes engineering corp |
| Systems & Processes Engineering Corporat | systems and processes engineering corp |
| Systems & Processes Engineering Corporat | systems and processes engineering corp |
| SYSTEMS AND PROCESSES ENGINEERING CORP.(SP | systems and processes engineering corp |

# Overview of the project directory. . .

```
matching/
    data/
        AwardsToExcel.csv --- DOD SBIR/STTR awards
        lob_indus.txt     --- OpenSecrets.org lobbying expenditures
    results/              --- A folder to store results
    __init__.py           --- Makes Python treat the directory as a package
    match.py              --- Main script to run the match
    preprocess.py         --- Extra functions to help with processing
```

# Finding (A)... but what does it mean?



Lobbying and contracts (r = 0.21)

# Finding (B). . . but what does it mean?



Lobbying and awards (r = 0.25)