# Slido Poll:
## www.sli.do
## (code: CBSA1)

# Demo URL:
https://rCanada.shinyapps.io/demo

# Workshop layout

- "Round table" poll (Slido)

- Common data quality problems
- Data engineering as a solution
- Common tools for common problems
- Discussion (Slido)

# Common data quality problems

- Facing same challenges turning "records" into "data"

  - Dates :                          '20210820' vs. 'dob 20 Aug 2021'

  - Names:                         'Dmitry Gorodnichy' vs. 'Dimitri Horodnytchyyi'

  - Business Names:          AC, AirCanada, Air Canada Corp.

  - Geographic Names:       'Ottawa Airport', 'YOW', Ottawa International Airport

# Poor data quality impedes interoperability

- Good interoperability allows various data to be linked and enriched
- Probabilistic (approximate, fuzzy) matching is used to link "noisy" data
    - All words need to be compared to each other
    - Various techniques in data linkages include: using edit metrics, look-up tables, q-grams, phonetic, heuristics, …
- However, probabilistic matching has its share of challenges as well
    - How to assign threshold?
    - How to measure quality?
    - Lost nuances?
        - E.g., Bell Canada vs. Shell Canada
- No perfect solution

# Data engineering to address data quality

- Data engineers develop techniques to standardize and organize data to help address data integrity, e.g.,

  - 'Ottawa Airport' → 'YOW'
  - 'YOW' → 'YOW'
  - 'Ottawa International Airport → 'YOW'

- On average, 80% of efforts of data scientists goes to address data engineering issues

# Common Tools for Common Problems

- In GoC, we are working on the same set of data engineering problems
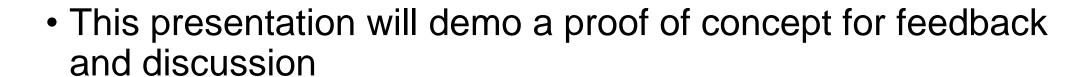
    - Standardizing various data fields
    - Cleaning, linking and searching data

    so we can carry out analysis

- Often, data scientists end up "reinventing the wheel"

- We need to build common data engineering tools for common GoC data engineering problems

# Vision for Solution

- We need a set of 'libraries' that are built and maintained by GoC data science community that is
  - Open
  - Free
  - Available to any data scientist who needs them
  - GCcode helps us to do that
  - R already has many libraries, supported by global community


- This presentation will demo a proof of concept for feedback and discussion

# Discussion

- Our works has just started. Much more ahead.
  - We build on Public Data (esp. Open Canada Data) <u>and</u> Public knowledge (esp. R global community)
  - We build solution (for entire GC community) <u>and also</u> we build Community of Practice
  - Codes and resources: [https://gccode.ssc-spc.gc.ca/r4gc/](https://gccode.ssc-spc.gc.ca/r4gc/)

- Planned milestones:
  - rCanada Package, Testbed App, Toolkit App: 2021-2022
  - Use cases (for on-going Agency needs):  Spring - Winter 2021

- We need your help!
  - curating & organizing DE challenges and public domain solutions (codes/papers)
  - curating & organizing public domain Data-sets
  - testing &  benchmarking

# Thank you !

Dmitry.Gorodnichy@cbsa-asfc.gc.ca

# Demo time

https://rCanada.shinyapps.io/demo