# Data Engineering with R, R Markdown, Shiny and algorithms

Dmitry Gorodnichy
Chief Data Office
Data Science Division

StatsCan – Applied ML for Text Analysis CoP
25 March 2021, Ottawa

This presentation builds on the results presented at the GCData2021 conference Data Literacy Fest workshop "**Data Engineering Challenges and Solutions: Demo of Shiny".** Please visit  https://gccollab.ca/discussion/view/7407617 to view the workshop presentation/demonstration and results.

# Outline

- "Round table" (Slido)

- Data Cleaning/Linking Problems: General, Canada specific
- Vision for solution: Methodological approach to Data Engineering
- Tools of the trade: algorithms, R, R Markdown, Shiny
- Use Cases: Record linking, COVID, web scraping, topic extraction
- Next steps: Introducing GCcollab **Use R!** and GCcode **r4gc** groups

- Demos & Discussion

# Problems - general

- Dates :                            '20210820' vs. 'dob 20 Aug 2021'
- Names:                             'Dmitry Gorodnichy' vs. 'Dimitri Horodnytchyyi'
- Business Names:                AC, AirCanada, Air Canada Corp.
- Geographic Names:           Ottawa, Orleans, Orléans


- General Text :                   "<tag> ca$h 4 u !  Sooo… C O O L!   Cant believe it ☹ "
- Postal:                             "klo 0O1" vs "K100o1"
- Text matching:                  Phrase matching, topics/keywords detection

# 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

# Problem – Canada specific

- Hundreds of dataset related to Canada

- Hundreds of data scientists working with those datasets

- Million lines of codes???

- What about quality?

# Vision

- Build data-driven solution for entire GC community
  and also
  build Community of Practice

- that leverages Public Data (esp. Open Canada Data)
  and
  on Public knowledge (esp. R global community)


- so that it is
  of good quality, robust, transparent, scalable, reusable, documented,
  and sustainable over time

# Methodological approach to Data Engineering

**Software Engineering** - a sub-field of Computer Science  focused on developing "scientific and technological **knowledge**, **methods**, and **experience** to the design, **implementation**, **testing**, and **documentation** of <u>software</u>"  [IEEE Vocabulary]

- "**Software engineering** encompasses <u>**not just the act of writing code, but all of the tools and processes**</u> **an organization uses to build and maintain that code over time.** …

- **Software engineering** can be thought of as "**programming integrated over time.**" [Software Engineering at Google]

**Data Engineering** - a sub-field of [Science and Engineering] that is focused on  developing scientific and technological **knowledge, methods,** and **experience** to the design**, implementation, testing**, and **documentation** of <u>data-driven solutions</u>.

# Vision (cntd)

- Stream 1: development of knowledge
    - https://gccode.ssc-spc.gc.ca/r4gc/
    - https://gccollab.ca/groups/profile/7391537/enuse-rfruse-r
    - https://github.com/canada-ca

- Stream 2: development of codes
    - R packages
    - Testbeds
    - Toolkits
    - Use cases

# Taxonomy of DE tasks

Single-variable:

- General text cleaning and formatting

- Dates and time-stamps extraction, cleaning and formatting

- Canadian Postal code and municipality names recovery

Multiple-variable:

- De-duplication of entries

- Large-scale Records Linking

Relations-based:

- Entity resolution

- Text analysis and plagiarism detection

# Tools

- Algorithms
- R
- R Markdown
- Shiny

# Algorithms – means to automate, scale, re-use!

## Input: **X**

(<u>any</u> raw "noisy" unknown data)



## Output: **Y**

(<u>meaningful</u> "filtered" result/conclusion)

# R – fastest growing in popularity environment

| Jul 2020 | Jul 2019 | Change | Programming Language | Ratings | Change |
|----------|----------|--------|----------------------|---------|--------|
| 1 | 2 | ▲ | C | 16.45% | +2.24% |
| 2 | 1 | ▼ | Java | 15.10% | +0.04% |
| 3 | 3 | | Python | 9.09% | -0.17% |
| 4 | 4 | | C++ | 6.21% | -0.49% |
| 5 | 5 | | C# | 5.25% | +0.88% |
| 6 | 6 | | Visual Basic | 5.23% | +1.03% |
| 7 | 7 | | JavaScript | 2.48% | +0.18% |
| 8 | 20 | ▲▲ | R | 2.41% | +1.57% |
| 9 | 8 | ▼ | PHP | 1.90% | -0.27% |
| 10 | 13 | ▲ | Swift | 1.43% | +0.31% |
| 11 | 9 | ▼ | SQL | 1.40% | -0.58% |

# R Markdown – describes algorithm (from data to result)

```
---
title: "NLP analysis of TBS-ATI data"
# author: "Source: https://github.com/open-data/TBS-ATI-NLP_Exploration"
output:  html_document
---

```{r}
source("TBS-ATI-functions.R");
# library(ATIP) # Eventually this could be converted to package or
```

Top 9 departments

```{r}
dtATI <- readATI()

# owners = ati%>%group_by(owner)%>% count()%>% ungroup() %>% top_n(9, n) %>% pull(owner)
aStrOwners <- dtATI[, .N, by=owner] %>% .[order(-N)] %>% .[1:9, owner]
```

...
```

# R Shiny – enables interactive testing and dashboards

```
---
title: "Data Engineering Testbed"
# author: "Source: https://gccode.ssc-spc.gc.ca/gorodnichy/rCanada"
output:   flex_dashboard
runtime: shiny
---

```{r}
source("rCanada-functions.R");
#library(rCanada)
dtNames <- as.data.table(lexicon::common_names ) %>% setnames("Name")
```

```{r de_1_dates.Rmd, child = 'de_1_dates.Rmd'}
```

```{r}
r.dtNames <- reactive({
    dtNames [, dist:=stringdist( Name, input$typedName, input$metric] })
```

...
```

# Use Cases

- Records deduplication and linking: https://rCanada.shinyapps.io/demo

- Web crawling:  …/demo/#section-web-crawling

  - Dates extraction

  - Finding nicknames and names variants

- UofT COVID data: https://rCanada.shinyapps.io/covid

- CBSA BWT data: https://itrack.shinyapps.io/border

- TBS PSES data: https://itrack.shinyapps.io/PSES

- TBS ATIP data: https://rCanada.shinyapps.io/TBS-ATI-NLP

# Steps for Record Linking / Deduplication

1. Data preparation: feature selection and preparation

2. Perform pair-wise comparison

3. Set constraints:
   soft vs. hard constraints, inter- vs. intra- class relationships

4. String similarity metrics (stringsim):
   q-grams vs. edit steps vs. heuristics vs. soundex

5. Algorithms:
   automated vs. semi-automated

6. Quality & Precision metrics: Accuracy vs. Precision/Recall


Ref:

# Steps for Text Analysis / Topic Extraction

- Load thesaurus and stop-words

- Words as token / unigrams

- Compute Top words, N/Total, bigram and n-grams

- Compute TF-IDF (term frequency – inverse document frequency)

- Compute correlations
  - Visualize bigram / n-grams relationship (ggraph, wordcloud)

- Topic modeling (w. topicmodel / textmineR):
  - Compute DTM (document term matrix)
  - Compute LDA (Latent Dirichlet Allocation)
  - Visualize dominant topics (ggplot, wordcloud)

Refs: S. Silge, D, Robinson,"Text Mining with R: a Tidy Approach", tidytextmining.com (github.com/dgrtwo/tidy-text-mining)
https://gccollab.ca/discussion/view/7404441/text-analysis-in-r

# Next steps

- The works has just started. Much more ahead.
  - [GCcollab group "Use R!"](#)
  - [GCcode group "r4gc"](#): [https://gccode.ssc-spc.gc.ca/r4gc/](https://gccode.ssc-spc.gc.ca/r4gc/)
  - Codes, Apps
- We need your help!
  - curating DE challenges and public domain solutions (codes/papers)
  - curating public domain Data-sets
  - testing & benchmarking
- Join GCcollab / GCcode groups.
- Join Friday "Lunch and Learn R" meet-ups
- Contact: Dmitry.Gorodnichy@cbsa-asfc.gc.ca

# Time for Demo and Discussion!

# Appendices

- Records cleaning, deduplication and linking: https://rCanada.shinyapps.io/demo
(Leveraging various R packages for data cleaning and linking)

- NLP topic modeling in TBS ATIP data: https://rCanada.shinyapps.io/TBS-ATI-NLP
(Leveraging the work of TBS and various R packages for text mining)

text2date() : converts text to a date using various decision logics.

### Test it:

Enter dates, any way you want, and observe how they get automatically converted to `YY MM DD` format.

[ 7 jul 35 ]    [ Reset table ]

### Result:

7 jul 35 --> 2035-07-07

| text | YY | MM | DD |
|---|---|---|---|
| 7jul35 | 2035 | 7 | 7 |
| 1935.08..7 | 1935 | 8 | 7 |
| DOB 12/26/2010... | 2010 | 12 | 26 |
| 26/12/1930 | 1930 | 12 | 26 |
| 7.VI.35 | 2035 | 6 | 7 |
| 7 jul35 | 2035 | 7 | 7 |
| 7 jul 35 | 2035 | 7 | 7 |

---

text2timestamp() : extracts automatically timestamp from free-form text

### Test it:

Enter a timestamp any way you want and observe how it gets converted to the same canonical timestamp `YY-MM-DD hh:mm:ss` format.

[ 2021-03-17 19:14:08 ]

### Result:

2021-03-17 19:14:08 --> 2021-03-17 19:14:08

| text | TIMESTAMP |
|---|---|
| 2010-04-14 22:00 | 2020-10-04 14:22:00 |
| 2010-04-14 10pm | 2020-10-04 14:10:00 |
| 2010-04-14-04-35-59 | 2010-04-14 04:35:59 |
| 2010-04-01-12-00-00 | 2010-04-01 12:00:00 |
| 20/2/06 11:16:16.683 | 2020-02-06 11:16:16 |
| 20100101120101 | 2010-01-01 12:01:01 |
| 2009-01-02 12-01-02 | 2009-01-02 12:01:02 |
| 2009.01.03 12:01:03 | 2009-01-03 12:01:03 |
| 2009-1-4 12-1-4 | 2009-01-04 12:01:04 |
| 2009-1, 5 12:1, 5 | 2009-01-05 12:01:05 |
| 200901-08 1201-08 | 2009-01-08 12:01:08 |
| 20090107 120107 | 2009-01-07 12:01:07 |
| 10-01-10 10:01:10 and p format: AM | 2010-01-10 10:01:10 |
| Created on 10-01-11 at 10:01:11 PM | 2010-01-11 22:01:11 |

# Performance of Dates & Timestamps recognition

**Performance of various string similarity metrics**

# Record deduplication

# Record linking

# Appendices

- Records cleaning, deduplication and linking:
  https://rCanada.shinyapps.io/demo
  (Leveraging various R packages for data cleaning and linking)

- NLP topic modeling in TBS ATIP data:
  https://rCanada.shinyapps.io/TBS-ATI-NLP
  (Leveraging the work of TBS and various R packages for text mining)
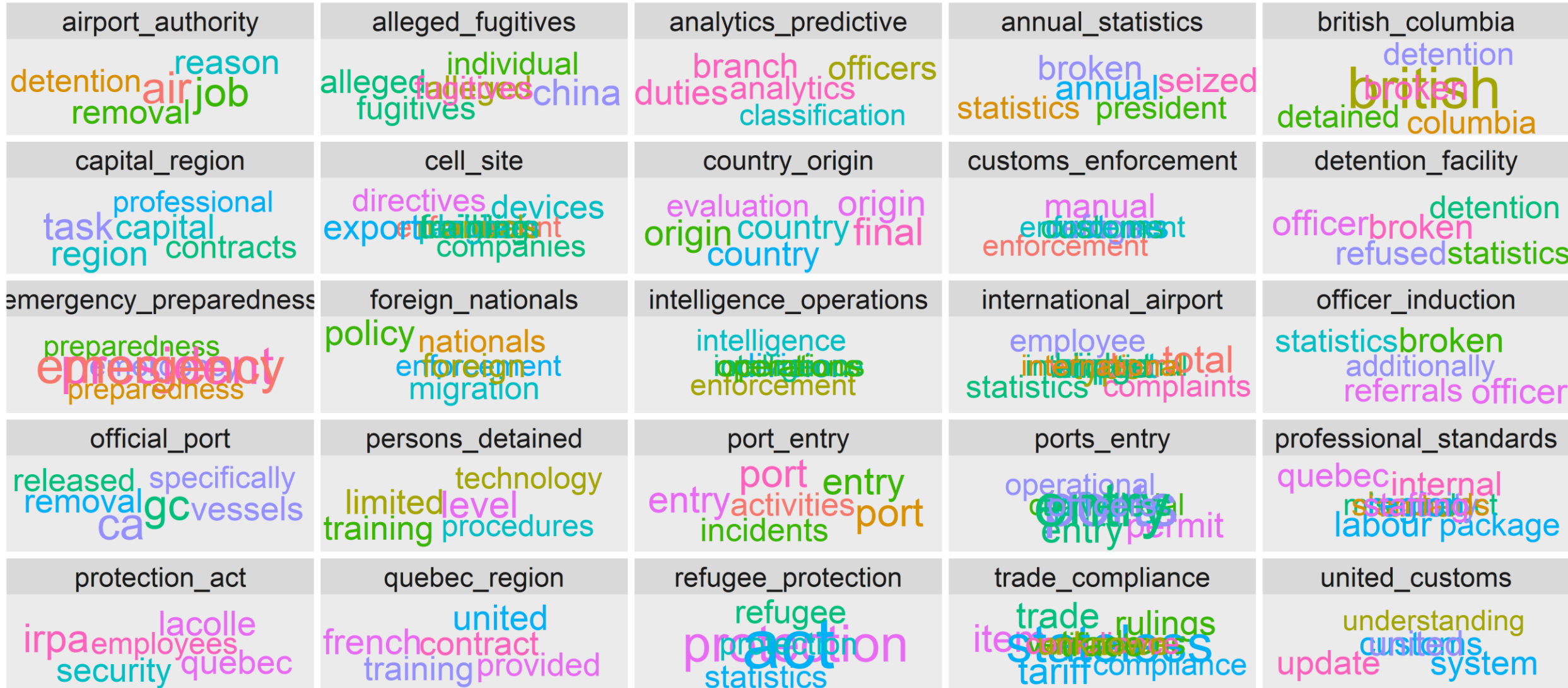
# Univariate and bivariate analysis of dataset variables

# 1-grams (single words)

# Topic modeling (30 main topics): wordcloud

**Topic modelling:**

**Graph / Network view**