

# Collaboration en science des données au sein du gouvernement du Canada

## Développement de bibliothèques R pour les tâches communes avec Données ouvertes au Canada

Jonathan Dench, analyste de la recherche, Division des résultats, Secrétariat du Conseil du Trésor du Canada  
Dmitry Gorodnichy, chercheur scientifique, Dirigeant principal des données, Agence des services frontaliers du Canada  
Patrick Little, conseiller, Gouvernement ouvert, Secrétariat du Conseil du Trésor du Canada  
Joseph Stinziano, analyste scientifique, Agence canadienne d'inspection des aliments

Slido: r4gc

Symposium international sur les questions de méthodologie de  
Statistique Canada de 2021





# Aperçu

- Raison d'être
- Vision
- Pourquoi R (pour la collaboration en sciences des données)
- Plateformes de collaboration du gouvernement du Canada (GC) (pour développer le savoir technique)
- Résultats clés (jusqu'ici)
- Quelles sont les prochaines étapes?
- Annexes : démonstrations et renseignements techniques



# Raison d'être

- Au GC, nous travaillons à la résolution des mêmes problèmes de science des données
  - Nous travaillons avec les mêmes données (p. ex. géospatiales, StatCan, ouvert.canada.ca)
  - Nous créons de nombreux outils de visualisation, d'analyse et de production de rapports semblables
  - Nous relevons bon nombre des mêmes défis en matière d'ingénierie et d'extraction de données
- Défis
  - Souvent, les scientifiques spécialistes des données finissent par « réinventer la roue » et ne sont pas en mesure de suivre le rythme du développement rapide d'outils en science des données.
  - Le manque de collaboration et d'examen par les pairs entraîne des risques, comme l'inefficacité et la production de solutions sous-optimales
- Nous pouvons en faire beaucoup plus si nous tirons parti du travail de chacun!
  - Discussion lors de la conférence du GC sur les données en 2021, [Atelier sur l'ingénierie des données](#)



# Vision

Pour assurer l'adoption d'approches normalisées et uniformes de la science des données à l'échelle du gouvernement du Canada, nous devons :

1. Développer et maintenir nos compétences et notre base de connaissances
2. Élaborer des codes et des outils pour les problèmes communs de science des données
  - Contribution, examen et tenue à jour par la communauté des sciences des données du GC
  - Ouvert et gratuit – accessible à tout scientifique spécialiste des données qui en a besoin

En tirant parti des éléments exemplaires et déjà disponibles au sein du GC :

1. Plateformes de collaboration : **gccode, gccollab, gcwiki, github**
2. Environnement de programmation : **R**



# Pourquoi R?

1. Graphiques avancés avec *ggplot2* et ses extensions
2. Production automatisée de rapports, de tutoriels et de manuels avec *RMarkdown*
3. Élaboration simplifiée de progiciel avec *devtools*
4. Élaboration et déploiement d'interfaces interactives et de tableaux de bord simplifiés avec *Shiny*
5. « Le meilleur pour le géocalcul »
6. Conception commune ordonnée, partagée entre les progiciels
7. Répertoire de progiciels organisé, testé par des pairs sur le CRAN
8. RStudio IDE (Integrated Development Environment) sur ordinateur de bureau et nuage (rstudio.cloud)
9. Prise en charge complète et interopérabilité avec Python à partir du même IDE
10. Mouvement mondial dirigé par RStudio pour l'éducation et l'avancement en R (rstudio.com)

<https://geocompr.robinlovelace.net/intro.html#why-use-r-for-geocomputation>

<https://gccollab.ca/discussion/view/7404883/why-r>



# Plateformes de collaboration

Restreint au GC :

- <https://gccode.ssc-spc.gc.ca/r4gc/>
- <https://wiki.gccollab.ca/UseR!>
- <https://gccollab.ca/groups/profile/7391537/Use-R>
  - [Sciences des données avec R – « déjeuner-formation » Rencontres du vendredi](#)



Pour le grand public :

- <https://github.com/open-canada>
  - Documentation NON CLASSIFIÉE pour les déjeuners-formations
  - Applications (p. ex. <https://open-canada.github.io/Apps/atip>)
- Affichage CRAN (idéal pour les progiciels finis)





**GitLab** Projects ▾ Groups ▾ More ▾

**r4gc**

**r4gc**

Group ID: 7049

▾

Group of Data Scientists and collaborators working with R for sharing Codes and Knowledgebase. See GCcollab group for more details:  
<https://gccollab.ca/groups/profile/7391537/Use-R>

**Subgroups and projects**

Shared projects Archived projects

> **GC packages** Owner

Packages we are building from "raw" codes. Join the effort ! See /packages101 and <https://gccollab.ca/di...>

0 5 1

> **Codes** Owner

Various "raw" R codes contributed by GC community.

0 7 2

▾ **Resources** Owner

Books, tutorials, presentations, blogs on R. NB: This folder contains subfolders that are visible only to log...

0 4 1

**IntroSpatialAnalysis**

Tutorials and Codes for Geo/Spatial coding and visualization in R. See <https://gccol...>

0 1 day ago

**howTo**

This is a project that will be developed as a series of RMarkdown documents, typic...

0 2 months ago

**meetings**

This is what we discuss at our weekly meet-ups

0 4 months ago

**gccode101**

You have questions on how to use GCCode or git? - Here you'll find the answers!

4 1 month ago

7





GCwiki

Connecting people  
and ideas[Main Page](#)  
[Browse categories](#)  
[Random page](#)  
[Help](#)[Actions/Tools](#)[Special pages](#)[Tools](#)[Related changes](#)  
[Printable version](#)  
[Permanent link](#)[Print/export](#)[Download as PDF](#)  
[Printable version](#) [GCaccount](#) [GCcollab](#) [GCmessage](#)Page [Discussion](#)[Read](#)[View source](#)[View history](#)

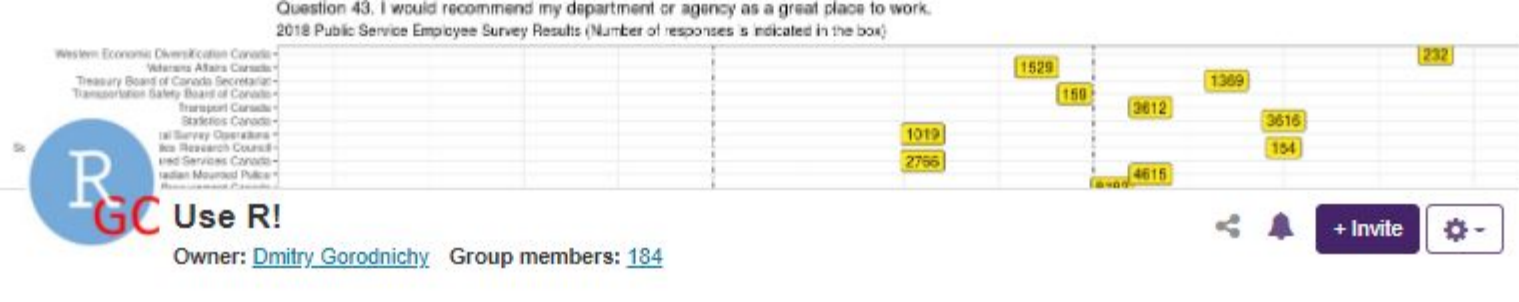
## UseR!

[Data Science Communities of Practice - UseR!](#)

**This page provides the list of discussions organized by the GCcollab's [Use R! group](#). Please consider contributing to those discussions by joining the [Use R! group](#) and participating in group's weekly "[Lunch and Learn Data Science with R](#)" meetups.**

- General topics:
  - [Why R?](#)
  - [Best way to start \(and keep learning\) R](#)
  - [Events and Forums for R users](#)
  - [From Excel to R](#)
  - [R with Python \(and other languages/tools\)](#)
  - [Efficient programming in R \(coding style, memory-efficient coding, collaboration-ready codes, source control\)](#)
  - **[data.table](#)** for efficient data processing
  - [Reading various kinds of data in R](#)
  - [Open R codes for GC: on GCcode and GitHub](#)
  - [RStudio news and tricks](#)
- Specialized topics:
  - **[ggplot2](#)** and its extensions for data visualization
  - **[Shiny](#)** for Interactive Data Visualization, Analysis and Web App development
  - **[R Markdown](#)** for automated and reproducible data science
  - [Record Linking](#) and other Data Engineering tasks in R
  - [Geo/Spatial](#) coding and visualization in R
  - [Text Analysis](#) in R
  - [Machine Learning and Modeling](#) in R
- Webinars and Tutorials (NB: you need to join the "[Lunch and Learn Data Science with R](#)" meetups group to access recordings of these sessions)
  - [30 Jul 2021: Geospatial data tools in R \(code\)](#)
  - [16 Jul 2021: Dual Coding - Python and R unite !](#)
  - [9 Jul 2021: Exploring ggplots \(recording, code\)](#)
  - [2 Jul 2021: Parsing GC Tables \(code\)](#)
  - [25 Jun 2021: Using the Open Government Portal API within R \(recording, code on \[github.com/open-canada\]\(#\)\)](#)
  - [21 Apr 2021: Analyzing PSES results using R and Shiny](#)
  - [16 Apr-15 May 2021: Building R packages \(recording, code\)](#)





[Activity](#) [Discussion](#) [Files](#) [Blog](#) [More-](#)

## Discussion topics

Add discussion topic

### [Best way to start \(and keep\) learning R](#)

The number of resources and ways to learn R is enormous. Some of us had tried many of them until we found the ones that we believe are the best ones. Share them here! Here's how my personal recommendation - quoted...

[4 Replies](#)

Started by [Dmitry Gorodnichy](#) 2021-03-05 22:52

[3 likes](#)

### [Shiny for Interactive Data Visualization, Analysis and Web App development](#)

This discussion thread is dedicated to Shiny package - a RStudio-curated tool for developing and deploying Interactive Data Visualization and Analysis tools and applications. Share your experiences, tricks, tools and questions...

[6 Replies](#)

Started by [Dmitry Gorodnichy](#) 2021-03-05 23:22

[3 likes](#)

### [Geo/Spatial coding and visualization in R](#)

There's much effort to across many GC departments to link and visualize geo-data. This discussion is the place to share your results, ideas or problems related to the problem. Below is a great resource to start, which also...

[4 Replies](#)

## Discussion

### [Reading \(all sorts of\) data in R - efficiently!](#)

My favourite methods for reading / writing "regular" .csv files has been 'data.table::fread() / fwrite()' - the fastest and automated in many ways. Now there's another one - with package 'vroom' -...

[2 Replies](#)

Started by [Dmitry Gorodnichy](#) 2021-04-22 12:40

[2 likes](#)

### [Excel -> R](#)

There was a keen interest expressed at last Friday meetup on transitioning from Excel to R. Incidentally, there was an RStudio Community Meet-up focused exactly on this topic: Meetup: Making the Shift from Excel to R:...

[2 Replies](#)

Started by [Dmitry Gorodnichy](#) 2021-04-14 16:38

[3 likes](#)

### [R communities in GC](#)

Roughly sorted by the level of group activity  
GCConnex (GC...





# 'Lunch and Learn' Data Science with R: Friday Meet-ups

Owner: [Dmitry Gorodnichy](#) Group members: [59](#)

[+ Invite](#)

[Activity](#) [Discussion](#) [Files](#) [Blog](#) [More -](#)

## 'Lunch and Learn' Data Science with R: Friday Meet-ups's files

New file folder

Upload a file

### Folder structure

↳ [Main folder](#)

### Did you know?

You can drag and drop files on to the folders to organize them!

[Main folder](#)

- ☐  [VIDEO & NOTES: 23 April-15 May, 2021. Building R packages - Sessions 1-4](#)  
By [Dmitry Gorodnichy](#) - 17 May 2021 @ 3:33pm - [Download](#)
- ☐  [VIDEO & NOTES: Meetup 28 May 2021. - Lookup table function w. data.table, delivering packages, new...](#)  
By [Dmitry Gorodnichy](#) - 1 June 2021 @ 4:49pm - [Download](#)
- ☐  [VIDEO & NOTES: Meetup 4 June 2021 - Utility functions and Converting Shiny to Exe](#)  
By [Dmitry Gorodnichy](#) - 4 June 2021 @ 8:31pm - [Download](#)
- ☐  [VIDEO & NOTES: Meetup 2021-06-11. - parse\\_gcTable\(\), api.canada.ca, shiny in aws, best PSES...](#)  
By [Dmitry Gorodnichy](#) - 11 June 2021 @ 8:19pm - [Download](#)
- ☐  [VIDEO & NOTES: Meetup 2021-06-18. How to dynamically assign\(\)](#)  
By [Dmitry Gorodnichy](#) - 18 June 2021 @ 5:34pm - [Download](#)
- ☐  [VIDEO & NOTES: Meetup 2021-07-25. Using API for working with Open Government Data within R](#)  
By [Dmitry Gorodnichy](#) - 25 June 2021 @ 5:28pm - [Download](#)
- ☐  [VIDEO: Lunch and Learn \(2021-07-09\). Automating advanced common visualizations with ggplot\(\)](#)  
By [Jonathan Dench](#) - 9 July 2021 @ 6:30pm - [Download](#)



# Prochaines étapes

- Le travail est en cours (et le sera toujours!)
- Il reste beaucoup à faire. Nous avons besoin de vous!
  - gérer les problèmes de données et les solutions du domaine public (codes/documents)
  - conservation des ensembles de données du domaine public
  - essai et analyse comparative
  - tutoriels, cas d'utilisation
- Joignez-vous à la communauté : Joignez-vous à GCcollab/Groupes de codes GC
- Coordonnées :
  - [Jonathan.Dench@tbs-sct.gc.ca](mailto:Jonathan.Dench@tbs-sct.gc.ca)
  - [Dmitry.Gorodnichy@cbsa-asfc.gc.ca](mailto:Dmitry.Gorodnichy@cbsa-asfc.gc.ca)
  - [Patrick.Little@tbs-sct.gc.ca](mailto:Patrick.Little@tbs-sct.gc.ca)
  - [Joseph.Stinziano@inspection.gc.ca](mailto:Joseph.Stinziano@inspection.gc.ca)



# Annexes : résultats clés (jusqu'ici)

- Code GC 101 for les employés du GC : <https://gccode.ssc-spc.gc.ca/r4gc/resources/gccode101>
- Progiciels R 101 pour les employés du GC :  
<https://gccode.ssc-spc.gc.ca/r4gc/gc-packages/packages101>
- Mode d'emploi : Tutoriels *rmarkdown* / *learnr* interactifs pour la résolution de divers problèmes
- Analyse et visualisation géospatiale : cas d'utilisation créé dans *rmarkdown*
- Ingénierie des données : progiciel et application pour les correspondances partielles, le couplage d'enregistrements et la déduplication <https://rCanada.shinyapps.io/demo>
- Applications interactives Shiny : pour AIRP, SAFF, COVID-19, temps d'attente à la frontière :  
<https://open-canada.github.io/Apps/atip> (~/[pses](#), ~/[covid](#), ~/[border](#))
- Travailler avec l'API du Portail du gouvernement ouvert dans R (à l'aide de *ckanr* et de *adobeanalyticsr*)
- Automatisation des scripts R à exécuter avec GitHub Actions



Les diapositives ci-dessous ne seront pas présentées,  
et ne servent qu'à titre de référence.



# Progiciels R 101

- Progiciel clé
  - *devtools* dispose d'une série de fonctions clés pour la configuration d'un progiciel, en particulier des structures de fichiers et de répertoire
- Code d'essai
  - La rédaction d'essais est une compétence clé pour assurer un code robuste et reproductible
    - L'objectif est de s'assurer que chaque étape d'une fonction fonctionne correctement avec un exemple reproductible
    - Par exemple, est-ce que le résultat de la fonction X est une liste?
  - Les progiciels *testthat* et *testthis* facilitent la rédaction des essais
- Principales considérations relatives aux progiciels R du GC
  - Octroi de licences
  - Qu'est-ce qui peut être soumis au CRAN? Quelles sont les répercussions légales?

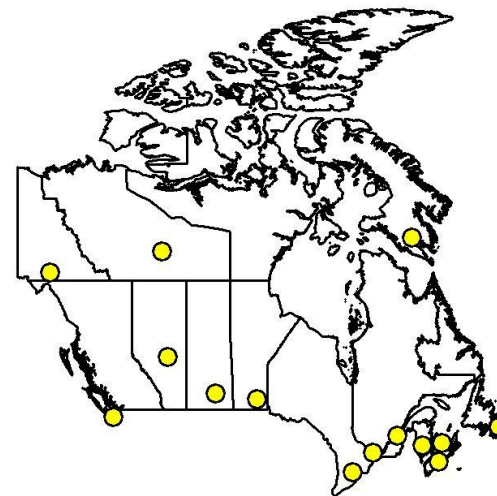




# Analyse géospatiale en R

- Guide et tutoriels

- « Applied Spatial Data Analysis with R », (2008) Roger Bivand et coll.
- « Geocomputation with R » (2021) (<https://geocompr.robinlovelace.net/>)
- Préparation d'une série d'ateliers et de code annoté pour le groupe R4GC.







# Travailler avec l'API du Portail du gouvernement ouvert (1)

- CKAN est un progiciel très utilisé pour alimenter les catalogues du portail de données ouvertes (data.gov, ouvert.canada.ca, data.gov.uk, etc.)
- CKAN offre une API qui peut être utilisée pour extraire des ensembles de données et des métadonnées du système, mais aussi pour créer, mettre à jour et gérer des ensembles de données.
- L'utilisation du progiciel *ckanr* offre une bonne expérience pour le développeur pour l'utilisation de l'API CKAN dans R.

ckanr v0.6.0.92

Search...

rOpenSci: The *ckanr* package

repo status Active CRAN ERROR R-check failing downloads 1030/month CRAN 0.6.0

ckanr is an R client for the CKAN API.

Description

CKAN is an open source set of tools for hosting and providing data on the web. (CKAN users could include non-profits, museums, local city/county governments, etc.).

ckanr allows users to interact with those CKAN websites to create, modify, and manage datasets, as well as search and download pre-existing data, and then to proceed using in R for data analysis (stats/plotting/etc.). It is meant to be as general as possible, allowing you to work with any CKAN instance.

Get started: <https://docs.ropensci.org/ckanr/>



# Travailler avec l'API du Portail du gouvernement ouvert (2)

## *Utiliser ckanr*

Que permet-il de faire?

Fonction	Commande API	Fonction CKANR
Obtenir de l'information à propos du système	<b>action/status_show</b>	ckan_info()
Dresser une liste des organisations qui publient des données	<b>action/organization_list</b>	organization_list()
Obtenir une liste d'ensembles de données sur le portail	<b>action/package_list</b>	package_list()
Récupérer des métadonnées pour un ensemble de données	<b>action/package_show/{id}</b>	package_show()
Chercher des ensembles de données	<b>action/package_search?q={some thing-to-search-for}</b>	package_search()
Créer un nouvel ensemble de données	<b>action/package_create</b>	package_create()
Mettre à jour une ressource existante	<b>action/resource_patch()</b>	resource_patch()

Exemple d'un cas d'utilisation : Quels sont les ensembles de données sur la COVID-19 qui sont disponibles sur le portail?

- Navigateur Web :  
[https://open.canada.ca/data/api/action/package\\_search?q="COVID"](https://open.canada.ca/data/api/action/package_search?q='COVID')
- Batch : **curl --verbose https://open.canada.ca/data/api/action/package\_search?q="COVID"**
- ckanr :

```
library(ckanr)
ckanr_setup(url="https://open.canada.ca/data")
search_results<-package_search(q="COVID", as= "table")
View(search_results$results)
```



# Analyse Web dans R avec Adobe Analytics

- Le GC utilise Adobe Analytics pour mesurer l'utilisation de Canada.ca ainsi que de plusieurs applications Web autonomes.
- Le progiciel *adobeanalyticsr* permet à un analyste d'extraire des données d'Adobe Analytics pour créer des rapports d'analyse Web dans R.
- Cette fonction peut être utilisée pour générer des extraits de données simples, mais également pour créer des rapports Rmd ou alimenter des applications Shiny.

adobeanalyticsr

## R Client for Adobe Analytics API 2.0

Connect to the Adobe Analytics API v2.0, which powers Analysis Workspace. The package was developed with the analyst in mind and will continue to be developed with the guiding principles of iterative, repeatable, timely analysis. New features are actively being developed and we value your feedback and contribution to the process. Please submit bugs, questions, and enhancement requests as [issues in this Github repository](#).





# Adobeanalyticsr – utilisation de base

- S'authentifier dans Adobe Analytics à l'aide d'un jeton OAuth utilisant la fonction `aw_token()`
- Utilisez la fonction `aw_freeform_table` pour créer un rapport basé sur les paramètres fournis
- Les fonctions `aw_get_metrics`, `aw_get_dimensions`, `aw_get_segments` peuvent être utilisées pour obtenir les paramètres disponibles.
- Analysez ou visualisez vos données dans R

```
topPages<-aw_freeform_table(  
  date_range = c("2021-04-01", "2021-04-28"),  
  company_id = Sys.getenv("AW_COMPANY_ID"),  
  rsid = Sys.getenv("AW_REPORTSUITE_ID"),  
  dimensions = c("prop65", "evar11"),  
  metrics = c("pageviews", "visits", "event25"),  
  search = "MATCH 'OG-GO'",  
  top = c(20)  
)
```

```
## Estimated runtime: 16.8sec./0.28min.
```

```
## 1 of 21 possible data requests complete. Starting the next 1 requests.
```

```
## A total of 20 rows have been pulled.
```

```
names(topPages)<-c("App Name", "Page Name", "pageviews", "visits", "downloads")  
kable(topPages)
```

App Name	Page Name	pageviews	visits	downloads
OG-GO	Open Government Portal	76222	24853	3
OG-GO	Open Government	19418	15608	0
OG-GO	Search Grants and Contributions	15383	3950	0
OG-GO	Search Government Contracts over \$10,000	14397	3159	0
OG-GO	Completed Access to Information Requests	12511	2898	0
OG-GO	blank page title	12187	5943	0
OG-GO	Canada Base Map Transportation (CBMT) - Open Government Portal	11204	9643	968



## Automatisation des scripts R à exécuter dans GitHub Actions

- GitHub Actions est une plateforme gratuite axée sur le flux de travail conçue pour automatiser les tâches de développement de logiciels telles que les CI/CD.
- GitHub Actions utilise des conteneurs Docker qui peuvent être configurés pour exécuter une myriade de systèmes d'exploitation et de logiciels différents, y compris R.
- Cela permet à l'utilisateur d'exécuter un script R basé sur un calendrier cron, ou d'autres événements tels qu'une modification du script.
- GitHub Actions est très utile pour automatiser les rapports ou d'autres charges de travail R.



# Comment exécuter un script R dans GitHub Actions

Actions Projects Wiki Security Insights Settings

master OpenGov\_R\_Scripts / .github / workflows / main.yml

View runs

Go to file

...

PatLittle Update main.yml

Latest commit bfa7251 on Feb 3 History

1 contributor

25 lines (22 sloc) 590 Bytes

Raw

Blame

🖨

✎

🗑

```
1 on:
2   push:
3   schedule:
4     # * is a special character in YAML so you have to quote this string
5     - cron: '0 11 * * *'
```

On push – s'exécute chaque fois que vous modifiez le code dans le répertoire  
On schedule – fonctionne selon les heures cron

```
6 jobs:
7   build:
8     runs-on: ubuntu-latest
9     steps:
10      - uses: actions/checkout@v2
11      - run: |
12        echo $PWD
```

```
14 - uses: r-lib/actions/setup-r@master
15   with:
16     r-version: '4.0.2'
```

r-lib setup-r configurera votre conteneur avec Ubuntu et installera la version de R que vous spécifiez

```
17 run: |
18   Rscript PD-Count.r
19   git config user.name github-actions
20   git config user.email github-actions@github.com
21   git add .
22   git commit -m "generated"
23   git push
```

Rscript est le script R que vous voulez exécuter

git. add  
git. commit  
git. push

git. add, git. commit, git. push enregistrera la sortie de tout  
fichier que vous créez/modifiez dans le répertoire GitHub



# Ingénierie des données

## Nettoyage des enregistrements, déduplication et couplage

<https://rCanada.shinyapps.io/demo>

Tire parti du travail de l'ASFC, de divers progiciels R pour le nettoyage et le couplage des données, et du cadre Shiny de RStudio

Cas d'utilisation compris :

- Recherche Web : [.../demo/#section-web-crawling](https://rCanada.shinyapps.io/demo/#section-web-crawling)
  - Extraction de dates
  - Trouver des surnoms et des variantes de noms





# Difficultés liées au couplage de données

- Dates : '20210820' par rapport à 'dob 20 Aug 2021'
- Noms : 'Dmitry Gorodnichy' par rapport à 'Dimitri Horodnytychyyi'
- Noms d'entreprise : AC, AirCanada, Air Canada Corp.
- Noms géographiques : Ottawa, Orleans, Orléans
- Texte général : "<tag> ca\$h 4 u ! Sooo... C O O L! Cant believe it 😞 "
- Poste : "klo 0O1" par rapport à "K100o1"
- Correspondance de texte : Correspondance de phrases, détection de sujets/mots clés

`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.

7jul35

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

Dates de nettoyage



searchName(name) : find similar names

Find similar names, using a variety of string similarity metrics. For definitions of all metrics.

Type a name:

Dmitry

String similarity metric:

jaccard

Metric threshold



Dates

Postal

**Names**

for definitions

**Correspondance de nom  
approximative  
(correspondance  
partielle/probabiliste)**

Result

Search and Save:

	Name	osa	lv	hamming	lcs	qgram	cosine	jaccard	jw	soundex
	<char>	<num>	<num>	<num>	<num>	<num>	<num>	<num>	<num>	<num>
1:	myrtie	5	5	5	8	2	0.167	0.286	0.306	1
2:	myrtis	5	5	5	8	2	0.167	0.286	0.306	1
3:	timmy	4	4	Inf	7	3	0.228	0.333	0.411	1
4:	demetria	4	4	Inf	6	4	0.355	0.375	0.278	0
5:	demetrice	5	5	Inf	7	5	0.473	0.375	0.296	0
6:	meredith	7	7	Inf	8	4	0.355	0.375	0.403	1
7:	merideth	7	7	Inf	8	4	0.355	0.375	0.403	1
8:	meridith	7	7	Inf	8	4	0.225	0.375	0.403	1
9:	myrtice	6	6	Inf	9	3	0.228	0.375	0.337	1
10:	armida	5	5	6	8	4	0.423	0.429	0.444	1
11:	marita	5	5	6	6	4	0.423	0.429	0.347	1
12:	marti	4	5	Inf	7	3	0.270	0.429	0.261	1
13:	marty	3	4	Inf	5	3	0.270	0.429	0.261	1
14:	mertie	5	5	5	8	4	0.423	0.429	0.306	1
15:	mindy	3	3	Inf	5	3	0.270	0.429	0.300	1
16:	mirta	3	4	Inf	5	3	0.270	0.429	0.261	1
17:	misty	3	3	Inf	3	3	0.270	0.429	0.178	1
18:	myriam	6	6	6	8	4	0.278	0.429	0.444	1
19:	myrta	4	5	Inf	7	3	0.270	0.429	0.411	1
20:	trinity	5	5	Inf	7	5	0.261	0.429	0.357	1
21:	trudi	6	6	Inf	7	3	0.270	0.429	0.544	1
22:	trudy	5	5	Inf	5	3	0.270	0.429	0.544	1
23:	yadira	5	5	5	6	4	0.423	0.429	0.333	1
24:	demetrius	5	5	Inf	7	5	0.385	0.444	0.296	0
	Name	osa	lv	hamming	lcs	qgram	cosine	jaccard	jw	soundex

☐ Use speed-optimized matching (experimental)

6789 YV

5







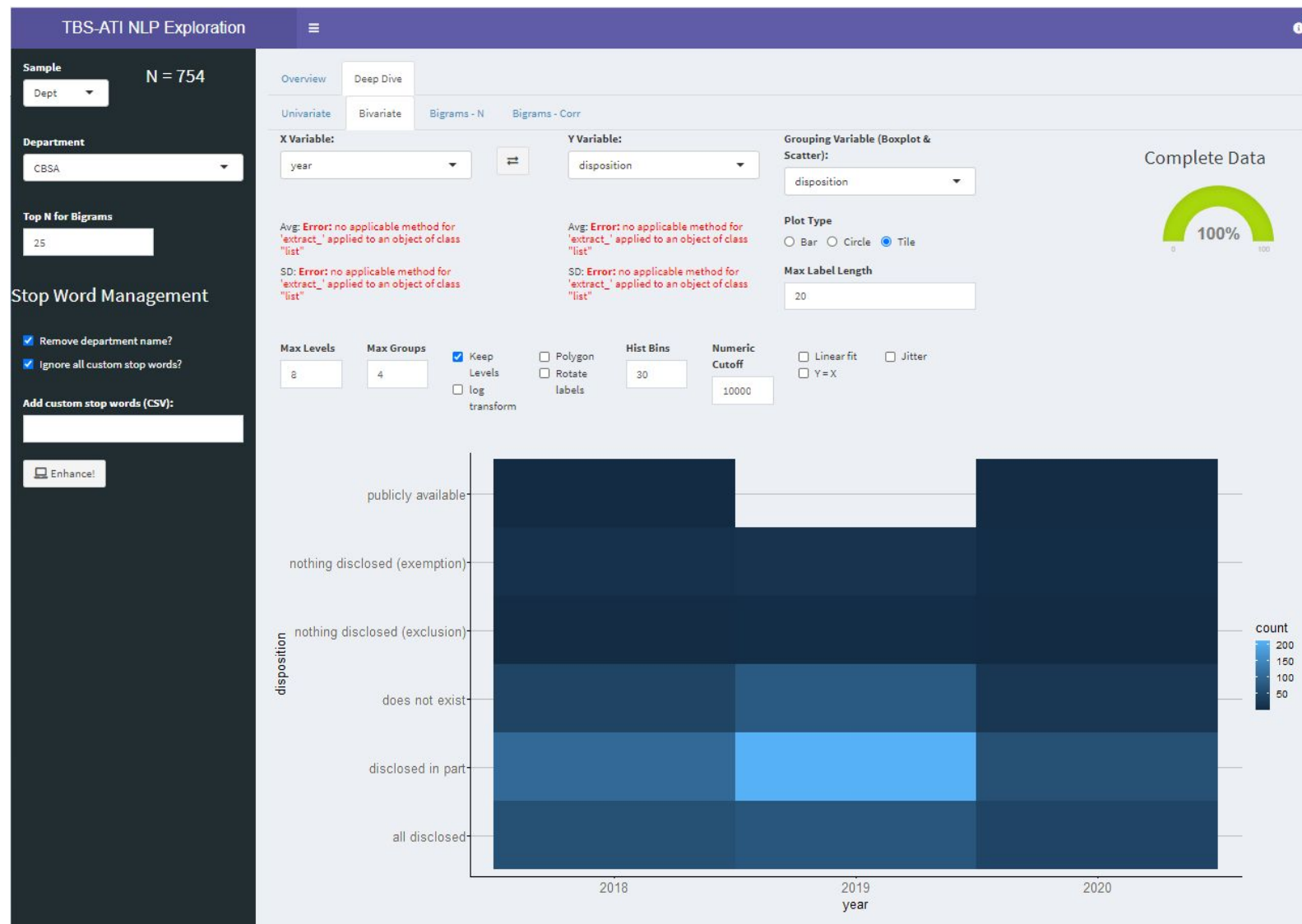
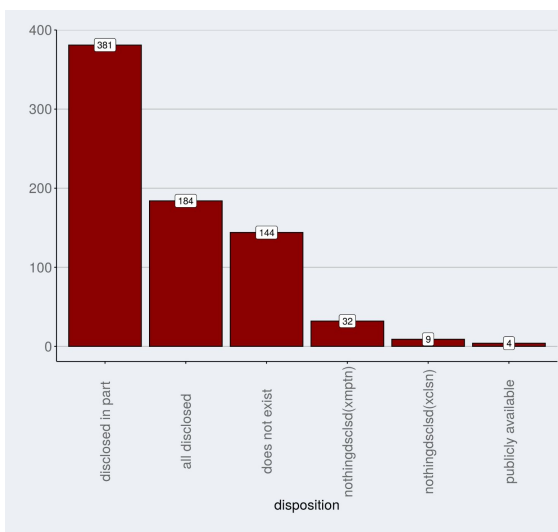
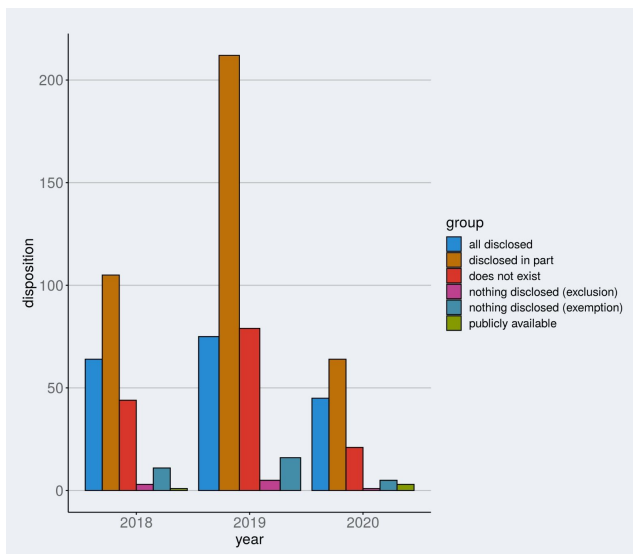
## **Modélisation NLP par sujet dans les données du SCT sur l'AIPRP**

<https://open-canada.github.io/Apps/atip>

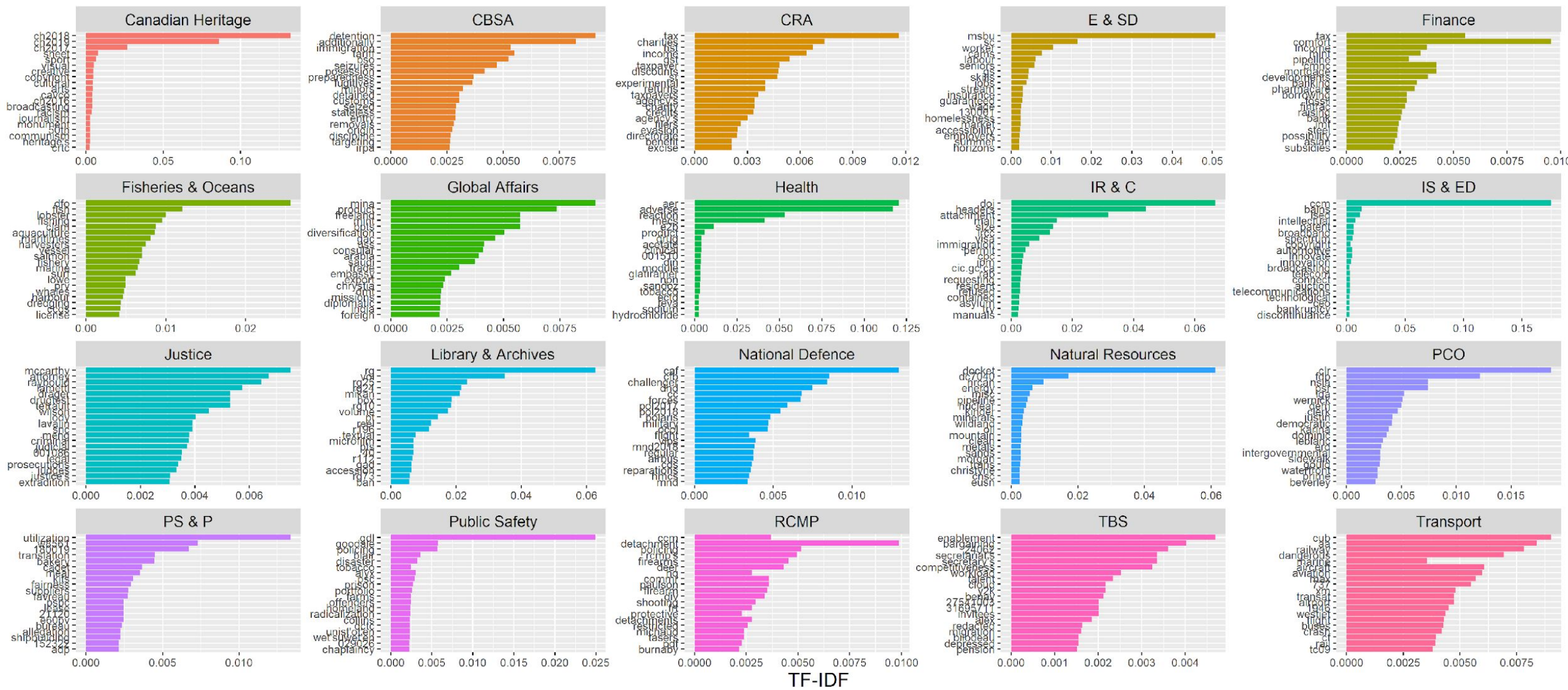
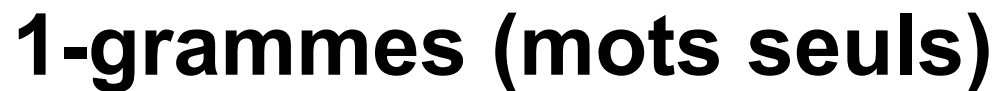
Tire parti du travail du SCT, de divers progiciels R pour l'exploration de texte et du cadre Shiny de RStudio



# Analyse univariée et bivariée des variables de l'ensemble de données



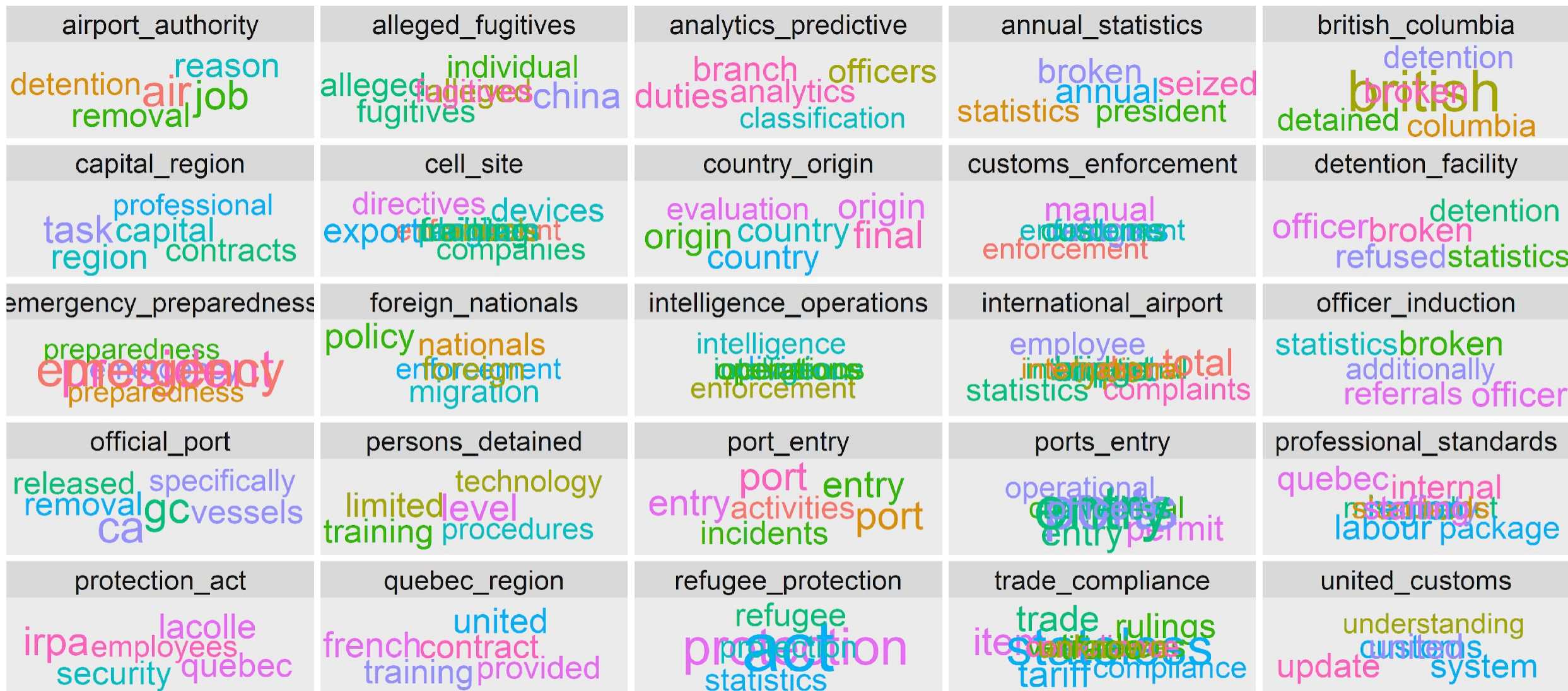








## Modélisation par sujet (30 sujets principaux): wordcloud





 Enhance!

Bigrams - Corr

0.3

entrie

Search:

	item1	item2	correlation
1	preparedness	emergency	0.971129452586342
2	emergency	preparedness	0.971129452586342
3	fugitives	alleged	0.936519829277792
4	alleged	fugitives	0.936519829277792
5	fugitives	control	0.909927934323004
6	control	fugitives	0.909927934323004
7	arrangements	supply	0.885947712418301
8	standing	supply	0.885947712418301
9	supply	arrangements	0.885947712418301
10	supply	standing	0.885947712418301
11	alleged	control	0.851958886147721
12	control	alleged	0.851958886147721
13	china	control	0.844085451931497
14	control	china	0.844085451931497
15	fugitives	individual	0.832662742876612
16	individual	fugitives	0.832662742876612
17	china	fugitives	0.831695559343427
18	fugitives	china	0.831695559343427

Showing 1 to 18 of 308 entries

[Previous](#)

1

2

Next

