

Analyse de données textuelles

Svetlana Zhuk

1. Introduction

En 1982, la Loi constitutionnelle de 1982 est adoptée par le Canada sans le consentement du Québec. Brian Mulroney, député du Parti progressiste-conservateur, promet au Québec une stratégie visant à l'intégrer dans l'Acte constitutionnel et à reconnaître son caractère distinct. En 1984, Brian Mulroney remporte les élections et forme le 33e Parlement. Notre objectif est d'analyser le premier mois de travail de la 33e législature et d'identifier les sujets les plus discutés. Notre question de recherche est donc la suivante :

- Quelles ont été les principales discussions menées par Brian Mulroney au cours du premier mois de la 33e législature ?

Pour répondre à cette question, nous utilisons des ensembles de données parlementaires contenant des transcriptions de débats fournies par [LIPAD](#). Le mois qui nous intéresse est novembre 1984.

2. Les données et nettoyage

Les données utilisées sont des débats transcrits du Parlement du Canada, ce qui nous permet d'effectuer une analyse des données textuelles. Pour mieux comprendre les discussions politiques, nous utiliserons le [lexicoder dictionary](#).

Ce dictionnaire facilite l'analyse automatisée des textes et est une source ouverte et gratuite. Pour permettre une analyse plus nuancée qui tienne compte du contexte, nous avons également développé la catégorie "Québec" afin d'examiner la fréquence des discussions liées à la situation du Québec au Parlement.

Dans les sections suivantes, nous détaillons les processus de nettoyage des données, de fusion des ensembles de données, d'intégration des dictionnaires et de représentation graphique.

Package version: 3.3.1

Unicode version: 14.0

ICU version: 71.1

Parallel computing: 8 of 8 threads used.

See <https://quanteda.io> for tutorials and examples.

Skipping install of 'clessnverse' from a github remote, the SHA1 (aa697496) has not changed :
Use `force = TRUE` to force installation

DISCLAIMER: As of July 2023, `clessnverse` is no longer under active development.
To avoid breaking dependencies, the package remains available "as is" with no warranty of any

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.4
v forcats    1.0.0      v stringr    1.5.1
v ggplot2    3.4.4      v tibble     3.2.1
v lubridate  1.9.3      v tidyr      1.3.1
v purrr      1.0.2
-- Conflicts ----- tidyverse_conflicts() --
x ggplot2::%>%() masks crayon::%>%()
x dplyr::filter() masks stats::filter()
x dplyr::lag() masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

```
#dictionary
```

```
lexicoder_en <- dictionary(file = "~/Dropbox/fas_1001_Zhuk/_tp/_tp3/dictionary/policy_agenc
```

```
#new_dictionary_Québec
```

```
quebec_dictionary <- list(quebec = c("sign", "patriation", "referendum", "veto", "federali
  dictionary())
```

```
#fusionner deux dictionnaires
```

```
lexicoder_en_m <- lexicoder_en |> stack()
quebec_dictionary_m <- quebec_dictionary |> stack()
new_dictionary <- bind_rows(lexicoder_en_m, quebec_dictionary_m) |>
  unstack(values ~ ind) |>
  dictionary()
```

```

Rows: 258 Columns: 15
-- Column specification -----
Delimiter: ","
chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl (1): basepk
lgl (1): speakerriding
date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 23 Columns: 15
-- Column specification -----
Delimiter: ","
chr (11): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl (1): basepk
lgl (2): subsubtopic, speakerriding
date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 316 Columns: 15
-- Column specification -----
Delimiter: ","
chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl (1): basepk
lgl (1): speakerriding
date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 450 Columns: 15
-- Column specification -----
Delimiter: ","
chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl (1): basepk
lgl (1): speakerriding
date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 329 Columns: 15
-- Column specification -----
Delimiter: ","

```

```
chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl (1): basepk
lgl (1): speakerriding
date (1): speechdate
```

```
i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 389 Columns: 15
```

```
-- Column specification -----
Delimiter: ","
chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl (1): basepk
lgl (1): speakerriding
date (1): speechdate
```

```
i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 263 Columns: 15
```

```
-- Column specification -----
Delimiter: ","
chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl (1): basepk
lgl (1): speakerriding
date (1): speechdate
```

```
i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 397 Columns: 15
```

```
-- Column specification -----
Delimiter: ","
chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl (1): basepk
lgl (1): speakerriding
date (1): speechdate
```

```
i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 285 Columns: 15
```

```
-- Column specification -----
Delimiter: ","
chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl (1): basepk
lgl (1): speakerriding
```

date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.

i Specify the column types or set `show_col_types = FALSE` to quiet this message.

Rows: 355 Columns: 15

-- Column specification -----

Delimiter: ","

chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...

dbl (1): basepk

lgl (1): speakerriding

date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.

i Specify the column types or set `show_col_types = FALSE` to quiet this message.

Rows: 367 Columns: 15

-- Column specification -----

Delimiter: ","

chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...

dbl (1): basepk

lgl (1): speakerriding

date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.

i Specify the column types or set `show_col_types = FALSE` to quiet this message.

Rows: 365 Columns: 15

-- Column specification -----

Delimiter: ","

chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...

dbl (1): basepk

lgl (1): speakerriding

date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.

i Specify the column types or set `show_col_types = FALSE` to quiet this message.

Rows: 425 Columns: 15

-- Column specification -----

Delimiter: ","

chr (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...

dbl (1): basepk

lgl (1): speakerriding

date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.

```

i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 297 Columns: 15
-- Column specification -----
Delimiter: ","
chr  (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl  (1): basepk
lgl  (1): speakerriding
date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 386 Columns: 15
-- Column specification -----
Delimiter: ","
chr  (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl  (1): basepk
lgl  (1): speakerriding
date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 444 Columns: 15
-- Column specification -----
Delimiter: ","
chr  (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl  (1): basepk
lgl  (1): speakerriding
date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 328 Columns: 15
-- Column specification -----
Delimiter: ","
chr  (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl  (1): basepk
lgl  (1): speakerriding
date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 485 Columns: 15
-- Column specification -----

```

```

Delimiter: ","
chr  (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl  (1): basepk
lgl  (1): speakerriding
date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
Rows: 467 Columns: 15
-- Column specification -----
Delimiter: ","
chr  (12): hid, pid, opid, speakeroldname, speakerposition, maintopic, subto...
dbl  (1): basepk
lgl  (1): speakerriding
date (1): speechdate

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.

```

```
#fusion des bases de donnees
```

```

data_parl_merged <- bind_rows(
  data_parl_1, data_parl_2, data_parl_3, data_parl_4, data_parl_5, data_parl_6, data_parl_7, data_parl_8, data_parl_9, data_parl_10, data_parl_11, data_parl_12, data_parl_13, data_parl_14, data_parl_15, data_parl_16, data_parl_17, data_parl_18, data_parl_19, data_parl_20, data_parl_21, data_parl_22, data_parl_23, data_parl_24, data_parl_25, data_parl_26, data_parl_27, data_parl_28, data_parl_29, data_parl_30, data_parl_31, data_parl_32, data_parl_33, data_parl_34, data_parl_35, data_parl_36, data_parl_37, data_parl_38, data_parl_39, data_parl_40, data_parl_41, data_parl_42, data_parl_43, data_parl_44, data_parl_45, data_parl_46, data_parl_47, data_parl_48, data_parl_49, data_parl_50, data_parl_51, data_parl_52, data_parl_53, data_parl_54, data_parl_55, data_parl_56, data_parl_57, data_parl_58, data_parl_59, data_parl_60, data_parl_61, data_parl_62, data_parl_63, data_parl_64, data_parl_65, data_parl_66, data_parl_67, data_parl_68, data_parl_69, data_parl_70, data_parl_71, data_parl_72, data_parl_73, data_parl_74, data_parl_75, data_parl_76, data_parl_77, data_parl_78, data_parl_79, data_parl_80, data_parl_81, data_parl_82, data_parl_83, data_parl_84, data_parl_85, data_parl_86, data_parl_87, data_parl_88, data_parl_89, data_parl_90, data_parl_91, data_parl_92, data_parl_93, data_parl_94, data_parl_95, data_parl_96, data_parl_97, data_parl_98, data_parl_99, data_parl_100, data_parl_101, data_parl_102, data_parl_103, data_parl_104, data_parl_105, data_parl_106, data_parl_107, data_parl_108, data_parl_109, data_parl_110, data_parl_111, data_parl_112, data_parl_113, data_parl_114, data_parl_115, data_parl_116, data_parl_117, data_parl_118, data_parl_119, data_parl_120, data_parl_121, data_parl_122, data_parl_123, data_parl_124, data_parl_125, data_parl_126, data_parl_127, data_parl_128, data_parl_129, data_parl_130, data_parl_131, data_parl_132, data_parl_133, data_parl_134, data_parl_135, data_parl_136, data_parl_137, data_parl_138, data_parl_139, data_parl_140, data_parl_141, data_parl_142, data_parl_143, data_parl_144, data_parl_145, data_parl_146, data_parl_147, data_parl_148, data_parl_149, data_parl_150, data_parl_151, data_parl_152, data_parl_153, data_parl_154, data_parl_155, data_parl_156, data_parl_157, data_parl_158, data_parl_159, data_parl_160, data_parl_161, data_parl_162, data_parl_163, data_parl_164, data_parl_165, data_parl_166, data_parl_167, data_parl_168, data_parl_169, data_parl_170, data_parl_171, data_parl_172, data_parl_173, data_parl_174, data_parl_175, data_parl_176, data_parl_177, data_parl_178, data_parl_179, data_parl_180, data_parl_181, data_parl_182, data_parl_183, data_parl_184, data_parl_185, data_parl_186, data_parl_187, data_parl_188, data_parl_189, data_parl_190, data_parl_191, data_parl_192, data_parl_193, data_parl_194, data_parl_195, data_parl_196, data_parl_197, data_parl_198, data_parl_199, data_parl_200, data_parl_201, data_parl_202, data_parl_203, data_parl_204, data_parl_205, data_parl_206, data_parl_207, data_parl_208, data_parl_209, data_parl_210, data_parl_211, data_parl_212, data_parl_213, data_parl_214, data_parl_215, data_parl_216, data_parl_217, data_parl_218, data_parl_219, data_parl_220, data_parl_221, data_parl_222, data_parl_223, data_parl_224, data_parl_225, data_parl_226, data_parl_227, data_parl_228, data_parl_229, data_parl_230, data_parl_231, data_parl_232, data_parl_233, data_parl_234, data_parl_235, data_parl_236, data_parl_237, data_parl_238, data_parl_239, data_parl_240, data_parl_241, data_parl_242, data_parl_243, data_parl_244, data_parl_245, data_parl_246, data_parl_247, data_parl_248, data_parl_249, data_parl_250, data_parl_251, data_parl_252, data_parl_253, data_parl_254, data_parl_255, data_parl_256, data_parl_257, data_parl_258, data_parl_259, data_parl_260, data_parl_261, data_parl_262, data_parl_263, data_parl_264, data_parl_265, data_parl_266, data_parl_267, data_parl_268, data_parl_269, data_parl_270, data_parl_271, data_parl_272, data_parl_273, data_parl_274, data_parl_275, data_parl_276, data_parl_277, data_parl_278, data_parl_279, data_parl_280, data_parl_281, data_parl_282, data_parl_283, data_parl_284, data_parl_285, data_parl_286, data_parl_287, data_parl_288, data_parl_289, data_parl_290, data_parl_291, data_parl_292, data_parl_293, data_parl_294, data_parl_295, data_parl_296, data_parl_297, data_parl_298, data_parl_299, data_parl_300, data_parl_301, data_parl_302, data_parl_303, data_parl_304, data_parl_305, data_parl_306, data_parl_307, data_parl_308, data_parl_309, data_parl_310, data_parl_311, data_parl_312, data_parl_313, data_parl_314, data_parl_315, data_parl_316, data_parl_317, data_parl_318, data_parl_319, data_parl_320, data_parl_321, data_parl_322, data_parl_323, data_parl_324, data_parl_325, data_parl_326, data_parl_327, data_parl_328, data_parl_329, data_parl_330, data_parl_331, data_parl_332, data_parl_333, data_parl_334, data_parl_335, data_parl_336, data_parl_337, data_parl_338, data_parl_339, data_parl_340, data_parl_341, data_parl_342, data_parl_343, data_parl_344, data_parl_345, data_parl_346, data_parl_347, data_parl_348, data_parl_349, data_parl_350, data_parl_351, data_parl_352, data_parl_353, data_parl_354, data_parl_355, data_parl_356, data_parl_357, data_parl_358, data_parl_359, data_parl_360, data_parl_361, data_parl_362, data_parl_363, data_parl_364, data_parl_365, data_parl_366, data_parl_367, data_parl_368, data_parl_369, data_parl_370, data_parl_371, data_parl_372, data_parl_373, data_parl_374, data_parl_375, data_parl_376, data_parl_377, data_parl_378, data_parl_379, data_parl_380, data_parl_381, data_parl_382, data_parl_383, data_parl_384, data_parl_385, data_parl_386, data_parl_387, data_parl_388, data_parl_389, data_parl_390, data_parl_391, data_parl_392, data_parl_393, data_parl_394, data_parl_395, data_parl_396, data_parl_397, data_parl_398, data_parl_399, data_parl_400, data_parl_401, data_parl_402, data_parl_403, data_parl_404, data_parl_405, data_parl_406, data_parl_407, data_parl_408, data_parl_409, data_parl_410, data_parl_411, data_parl_412, data_parl_413, data_parl_414, data_parl_415, data_parl_416, data_parl_417, data_parl_418, data_parl_419, data_parl_420, data_parl_421, data_parl_422, data_parl_423, data_parl_424, data_parl_425, data_parl_426, data_parl_427, data_parl_428, data_parl_429, data_parl_430, data_parl_431, data_parl_432, data_parl_433, data_parl_434, data_parl_435, data_parl_436, data_parl_437, data_parl_438, data_parl_439, data_parl_440, data_parl_441, data_parl_442, data_parl_443, data_parl_444, data_parl_445, data_parl_446, data_parl_447, data_parl_448, data_parl_449, data_parl_450, data_parl_451, data_parl_452, data_parl_453, data_parl_454, data_parl_455, data_parl_456, data_parl_457, data_parl_458, data_parl_459, data_parl_460, data_parl_461, data_parl_462, data_parl_463, data_parl_464, data_parl_465, data_parl_466, data_parl_467, data_parl_468, data_parl_469, data_parl_470, data_parl_471, data_parl_472, data_parl_473, data_parl_474, data_parl_475, data_parl_476, data_parl_477, data_parl_478, data_parl_479, data_parl_480, data_parl_481, data_parl_482, data_parl_483, data_parl_484, data_parl_485, data_parl_486, data_parl_487, data_parl_488, data_parl_489, data_parl_490, data_parl_491, data_parl_492, data_parl_493, data_parl_494, data_parl_495, data_parl_496, data_parl_497, data_parl_498, data_parl_499, data_parl_500)

```

```

data_parl_clean <- data_parl_merged |>
  select(speechdate, speechtext, speakername) |>
  mutate(speechtext = tolower(speechtext)) |> #lowercase
  na.omit() # remove missing variables

```

```
#analyse du dictionnaire
```

```

nov_1984 <- run_dictionary(data = data_parl_clean,
  text = speechtext,
  dictionary = new_dictionary) |>
  bind_cols(data_parl_clean) |>
  select(-c(doc_id, speechtext)) |>
  pivot_longer(!c(speechdate, speakername),
    names_to = "category",
    values_to = "n") |>

```

```

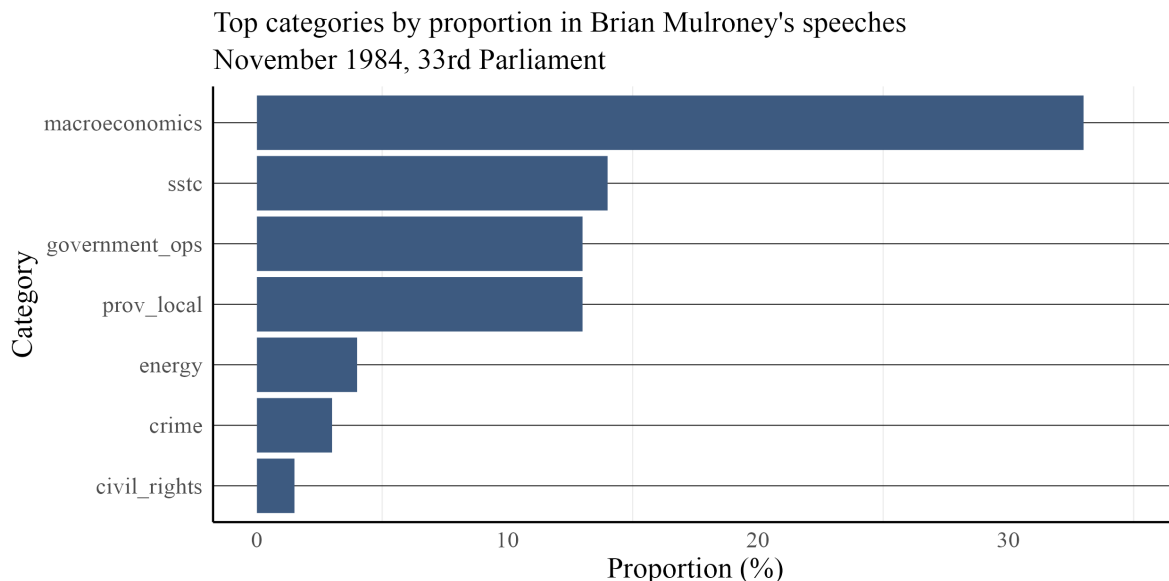
ungroup() |>
na.omit() |>
filter(speakername == "Martin Brian Mulroney") |>
group_by(category) |>
summarise(total = sum(n, na.rm = TRUE)) |>
ungroup() |>
mutate(prop = round(total/sum(total), 4)*100) |>
slice_max(total, n = 10)

```

100% expressions/words found

1.117 sec elapsed

3. Résultats



Dans nos résultats, nous mettons en évidence les principales catégories discutées par Brian Mulroney au cours du premier mois de novembre, qui marque le début de la 33e législature après qu'il a formé un gouvernement. Le sujet le plus fréquemment abordé est la macroéconomie, probablement lié aux discussions sur les accords de libre-échange avec les États-Unis. Cela conduit à des discussions plus approfondies et à des analyses détaillées dans la catégorie "macroéconomie". Cependant, la catégorie "Québec", un autre domaine d'intérêt, n'est pas apparue dans les discussions les plus importantes. Il est probable que les mots-clés que nous avons sélectionnés pour la catégorie "Québec" n'ont pas permis de capturer de manière adéquate les discussions relatives à l'accommodement du Québec. Par conséquent, il est nécessaire d'affiner la catégorie "Québec" dans notre dictionnaire.

Bibliographie

- Beelen, K., Thijm, T. A., Cochrane, C., Halvemaan, K., Hirst, G., Kimmins, M., Li-jbrink, S., Marx, M., Naderi, N., Rheault, L., Polyanovsky, R., and Whyte, T. (2017). “Digitization of the Canadian Parliamentary Debates.” *Canadian Journal of Political Science*, 50(3), 849–864. <https://doi.org/10.1017/S0008423916001165>.
- Albugh, Quinn, Julie Sevenans and Stuart Soroka. 2013. *Lexicoder Topic Dictionaries*, June 2013 versions, McGill University, Montreal, Canada.