Week 3: Classifying Texts

POP77142 Quantitative Text Analysis for Social Scientists

Tom Paskhalis

Overview

- Bag-of-words
- Document-Term Matrix
- Dictionaries
- Text classification

Bag-of-words

Quantifying Texts

When I presented the supplementary budget to this House last April, I said we could work our way through this period of severe economic distress. Today, I can report that notwithstanding the difficulties of the past eight months, we are now on the road to economic recovery.

In this next phase of the Government's plan we must stabilise the deficit in a fair way, safeguard those worst hit by the recession, and stimulate crucial sectors of our economy to sustain and create jobs. The worst is over.

This Government has the moral authority and the well-grounded optimism rather than the cynicism of the Opposition. It has the imagination to create the new jobs in energy, agriculture, transport and construction that this green budget will

docs t06_kenny_fg t05_cowen_ff t14_ocaolain_sf t01_lenihan_ff	12 9 3 12	because 11 4 3	had 5 8 3	into 4 5 4 4	get 8 5 7 2	some 4 5 3 11	through 3 14 7 9	next 4 13 2 16	where 5 4 3 14	many 7 9 5	irish 10 8 6
t01_lenihan_ff t11_gormley_gree t04_morgan_sf t12_ryan_green t10_quinn_lab t07_odonnell_fg t09_higgins_lab t03_burton_lab t13_cuffe_green	n 0 11 2 1 5 2 4	1 0 8 2 4 4 2 8	5 0 7 3 4 2 5 12	4 3 15 7 2 1 4 10	2 0 8 0 8 5 0 5	11 2 19 3 4 0 1 5	9 0 6 0 1 1 0 4	16 3 5 1 0 1 0 5	14 1 3 6 1 0 2 8	6 1 6 0 2 3 0 15	9 2 6 0 0 0 8
t13_culle_green t08_gilmore_lab t02_bruton_fg	4 1	8 10	7 6	4 4	3	6	4 0	5 6	1 16	2 5	11 3

Descriptive statistics on words

Classifying documents

Extraction of topics

Vocabulary analysis

Sentiment analysis

Tokenization

• Previously, we *tokenized* a single text:

```
1 text <- "The quick brown fox jumps over the lazy dog."
2 tokens <- quanteda::tokens(text)
3 tokens

Tokens consisting of 1 document.
text1:
[1] "The" "quick" "brown" "fox" "jumps" "over" "the" "lazy" "dog"
[10] "."</pre>
```

• But what if we have multiple texts?

Document-Term Matrix (DTM)

• Imagine in addition to the text above, we have another one:

```
1 text2 <- "The quick brown fox jumps over the lazy cat."
```

• If we were to create a **document-term matrix (DTM)**, it would look like this:

Document	brown	cat	dog	fox	jumps	lazy	over	quick	the
1	1	0	1	1	1	1	1	1	2
2	1	1	0	1	1	1	1	1	2

• Where each row corresponds to a document and each column to a token.

DTM

- In real life the documents are unlikely to have as much overlap.
- A more realistic example of 2 texts would be something like this:

```
1 speech1 <- "I thank the Deputy."
2 speech2 <- "Deputy, please resume your seat."
```

• If we we were to create a document-term matrix, it would look like this:

Document	deputy	i	please	resume	seat	thank	the	your
1	1	1	0	0	0	1	1	0
2	1	0	1	1	1	0	0	1

• This is quite a few zeros!

Sparse Matrix

- In practice, DTMs are often stored as **sparse matrices**.
- Such matrix only stores non-zero values and their positions.

```
1 library("Matrix")
2 sm <- Matrix::sparseMatrix(
3    i = c(1, 1, 1, 1, 2, 2, 2, 2, 2),
4    j = c(1, 2, 6, 7, 1, 3, 4, 5, 8),
5    x = c(1, 1, 1, 1, 1, 1, 1, 1),
6    dims = c(2, 8),
7    dimnames = list(
8        c("Doc1", "Doc2"),
9        c("deputy", "i", "please", "resume", "seat", "thank", "the", "your"))
10   )
11   sm</pre>
```

• Of course, we will not be creating these matrices manually.

DTM in R

- In R we can use the quanteda package to create DTMs.
- Note, it uses document-feature matrix (DFM) term.
- We start by tokenising our small corpus:

```
1 speeches <- c(speech1, speech2)
2 speeches_toks <- quanteda::tokens(tolower(speeches), remove_punct = TRUE)
3 speeches_toks

Tokens consisting of 2 documents.
text1 :
[1] "i" "thank" "the" "deputy"

text2 :
[1] "deputy" "please" "resume" "your" "seat"</pre>
```

• Then we create a DFM:

Bag of Words

- This representation of text is also known as **bag-of-words**.
- The key aspects of this approach are:
 - Single words are considered as relevant features of each document.
 - Documents are quantified by counting occurrences of words.
 - Word order is ignored.
 - Grammar and syntax are discarded.

Why Bag of Words?

- Simple.
- Efficient.
- Works well in many cases.
- Can be extended (co-occurrences aka "n-grams")
- Has been extensively used and validated in social sciences.
- Of course, there are cases when word order matters (e.g. text reuse).

Dictionaries

Word Meanings

- Words have meanings 🧐
- This allows us to take word usage as a proxy for the overall 'meaning' of a text.
- Certain kinds of words indicate certain kind of 'meanings'.
- Kinds of 'meanings':
 - Sentiment (e.g. positive, negative, etc.)
 - Emotions (e.g. anger, sad, happiness, etc.)
 - Topics (e.g. politics, sports, etc.)
 - Ideology (e.g. liberal, conservative, etc.)
 - Hate speech (e.g. sexism, homophobia, xenophobia, etc.)

Dictionaries

- Automated dictionary methods (ADM) exploit word usage to learn the 'meanings' of texts.
- Two steps:
 - 1. **Dictionary creation**: Define a list of words that represent a certain 'meaning'.
 - 2. **Dictionary application**: Count the number of words in a text that are in the dictionary.
- Dictionaries should be task-appropriate and validated.

Dictionary Structure

• We have seen dictionaries in the context of Python:

```
1 ideo_dict = {
2  "liberal": ["benefits", "worker", "trade union"],
3  "conservative": ["restriction", "immigration", "reduction"]
4  }
```

- Essentially, a dictionary is a set of **key-value pairs**.
- In the context of text analysis:
 - **Keys** labels for equivalence classes for the concept of interest.
 - Values terms or patterns that are declared equivalent occurrences of the key class

Dictionary vs Thesaurus

- A dictionary in a QTA sense is somewhat of a misnomer.
- Substantively, a dictionary is closer to a *thesaurus*.
- I.e. a list of canonical terms or concepts ('keys') associated with a list of synonyms.
- But unlike thesauruses, ADM dictionaries:
 - tend to be 'exclusive' (each value is associated with one key only)
 - do not always identify synonyms

Qualitative & Quantitative Text Analysis

- ADM dictionaries sit somewhere between more qualitative and fully automated approaches to text analysis.
- It is 'qualitative' in a sense that it requires identification of concepts and textual features associated with each of them.
- Dictionary construction involves a lot of contextual interpretation and qualitative judgment
- At the same time the application part is fully automated and perfectly reliable/replicable.

Some Famous Dictionaries

- General Inquirer (Stone et al. 1966): an early all-purpose dictionary (e.g. sentiment analysis) in general texts.
- Regressive Imagery Dictionary: designed to measure primordial vs. conceptual thinking.
- Linguistic Inquiry and Word Count (LIWC) (Pennebaker et al. 2001): large (paid) dictionary for many psychological and related concepts.

Example: LexiCoder

• The LexiCoder Sentiment Dictionary (Young and Soroka 2012): a dictionary for sentiment analysis in political texts, validated with human-coded news content.

```
1 data("data dictionary LSD2015", package = "quanteda")
 1 str(data_dictionary_LSD2015)
Formal class 'dictionary2' [package "quanteda"] with 2 slots
  ..@ .Data:List of 4
  .. ..$ :List of 1
  ....$ : chr [1:2858] "a lie" "abandon*" "abas*" "abattoir*" ...
  ....$ :List of 1
  ....$ : chr [1:1709] "ability*" "abound*" "absolv*" "absorbent*" ...
  .. ..$ :List of 1
  ....$ : chr [1:1721] "best not" "better not" "no damag*" "no no" ...
  ....$ :List of 1
  .....$ : chr [1:2860] "not a lie" "not abandon*" "not abas*" "not abattoir*" ...
  ..@ meta :List of 3
  ....$ system:List of 5
  .....$ package-version:Classes 'package_version', 'numeric_version' hidden list of 1
  .. .. .. ..$ : int [1:3] 1 9 9009
                         :Classes 'R_system_version', 'package_version', 'numeric_version' hidden list of 1
  .. .. ..$ r-version
  .. .. ...$ : int [1:3] 3 6 2
  .....$ system : Named chr [1:3] "Darwin" "x86_64" "kbenoit"
  ..... attr(*, "names")= chr [1:3] "sysname" "machine" "user"
  .....$ directory : chr "/Users/kbenoit/Dropbox (Personal)/GitHub/quanteda/quanteda"
  ....$ created : Date[1:1], format: "2020-02-17"
```

Example: Laver and Garry (2000)

- A hierarchical set of categories to distinguish policy domains and policy positions.
- Derived from one of the longest content analysis exercises in political science Manifesto Project (previously known as CMP).
- Five domains at the top level of hierarchy:
 - economy
 - political system
 - social system
 - external relations
 - "general" domain
- The dictionary was developed on a set of specific UK manifestos.



Example: Laver and Garry (2000)

- An accompanying quanteda.dictionaries package contains a lot of mentioned dictionaries, including the Laver and Garry (2000) dictionary.
- Alternatively, you can download the 'raw' dictionary as text from https://provalisresearch.com/Download/LaverGarry.zip

```
1 # The package is not available on CRAN,
 2 # so need to install it from GitHub
 3 remotes::install_github("kbenoit/quanteda.dictionaries")
 1 data("data_dictionary_LaverGarry", package = "quanteda.dictionaries")
 1 str(data_dictionary_LaverGarry)
Formal class 'dictionary2' [package "quanteda"] with 2 slots
  ..@ concatenator: chr " "
  ..@ names : chr [1:9] "CULTURE" "ECONOMY" "ENVIRONMENT" "GROUPS" ...
  ..@ .Data :List of 9
  ....$ :List of 4
  .....$ CULTURE-HIGH :List of 1
  .....$ : chr [1:8] "art" "artistic" "dance" "galler*" ...
  .. .. .. $ CULTURE-POPULAR:List of 1
  .. .. ...$ : chr "media"
                         :List of 1
  .. .. ..$ SPORT
  .. .. ... * : chr "angler*"
  .....$ : chr [1:3] "people" "war_in_iraq" "civil_war"
  ....$ :List of 3
  .. .. ..$ +STATE+:List of 1
  .... $ : chr [1:50] "accommodation" "age" "ambulance" "assist" ...
```

```
.. .. ..$ =STATE=:List of 1
.. .. .. ..$ : chr [1:71] "accountant" "accounting" "accounts" "advert*" ...
.. .. ..$ -STATE-:List of 1
.. .. ..$ : chr [1:62] "assets" "autonomy" "barrier*" "bid" ...
.. ..$ :List of 2
.. ..$ CON ENVIRONMENT:List of 1
```

Example: Dictionary Application

- Imagine we want to know which of the parties discusses immigration the most in their electoral manifesto.
- We can start by creating a very simple dictionary to answer this question:

```
imm_dict <- quanteda::dictionary(list(
   immigration = c("asylum*", "border*", "immigra*", "migrant*", "refugee*")
))

manifestos <- readr::read_csv("../data/ireland_ge_2024_manifestos.csv")

manifestos_toks <- quanteda::tokens(
   manifestos$text,
   remove_punct = TRUE,
   remove_numbers = TRUE,
   remove_symbols = TRUE
</pre>
```

Example: Dictionary Application

• Now we can apply the dictionary to the manifestos:

```
1 manifestos imm <- quanteda::dfm(</pre>
      quanteda::tokens lookup(manifestos toks, dictionary = imm dict)
  1 manifestos_imm
Document-feature matrix of: 9 documents, 1 feature (0.00% sparse) and 0 docvars.
       features
        immigration
docs
  text1
                 53
  text2
                 24
  text3
                 32
  text4
                 24
  text5
                 31
  text6
                 31
[ reached max_ndoc ... 3 more documents ]
```

Calculating Quantities of Interest

- Of course, the absolute number of matched terms is not, necessarily, informative.
- In the immigration focus example we can use the total number of matched terms M_i divided by the total number of words in the document N_i :

immigration_focus_i =
$$\frac{M_i}{N_i}$$

• If we were to try to scale the manifestos as pro- or anti- immigration (assuming we had a relevant dictionary), we could then try something like:

immigration_position_i =
$$\frac{M_i^{anti} - M_i^{pro}}{N_i}$$

• In other words, we would calculate an absolute proportional difference.

Scaling

- The previously described approach was used extensively in Manifesto Project.
- The problems, however, are:
 - Addition of irrelevant content shifts the scale toward zero.
 - Assumes the additional mentions increase emphasis in a linear scale
- One alternative (Laver & Garry, 2000):

immigration_position_i =
$$\frac{M_i^{anti} - M_i^{pro}}{M_i^{anti} + M_i^{pro}}$$

• Another alternative (Lowe, Benoit, Mikhaylov & Laver, 2011):

immigration_position_i =
$$\log \frac{M_i^{ann}}{M_i^{pro}}$$

Example: Dictionary Application

```
immigration_focus <- cbind(
manifestos,
quanteda::convert(manifestos_imm, to = "data.frame")

| ' | ' |
| (\(df) transform(df, ntokens = quanteda::ntoken(manifestos_toks)))() | ' |
| (\(df) transform(df, rel_imm = immigration/ntokens))() | ' |
| -[, c("party", "immigration", "ntokens", "rel_imm")] | ' |
| (\(df) `[`(df, order(df$rel_imm, decreasing = TRUE),))()</pre>
immigration_focus
```

```
party immigration ntokens
                                     rel imm
text5
        ΙI
                          7295 0.0042494859
                     31
text7
       PBP
                         11976 0.0021710087
                         27749 0.0019099787
text1
        A0
                         29110 0.0008244589
text4
        GR
                     24
        FF
text2
                         33676 0.0007126737
                     24
text8
                         58281 0.0006176970
        SD
                     36
text3
        FG
                     32
                         52942 0.0006044350
text9
        SF
                     28 48813 0.0005736177
                         63107 0.0004912292
text6
        LAB
                     31
```

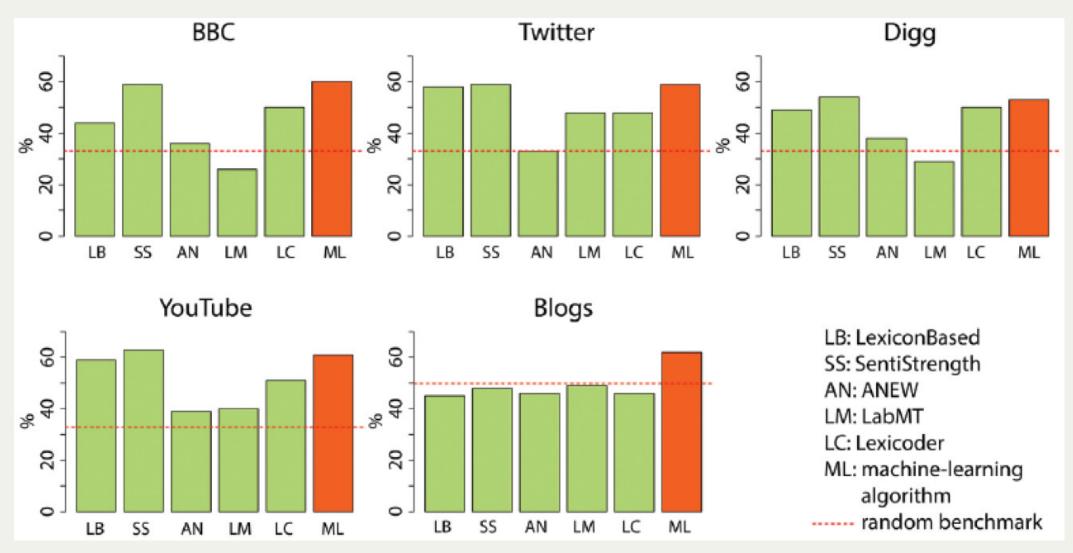
What to Do with Dictionary Results

- Describe the results.
- Scale the results (neg. vs pos., pro vs anti, left vs right, etc.).
- Could be used as features in downstream tasks:
 - Similarity measures (e.g. cosine)
 - ML-based classification
 - Topic modelling (seeded with keywords)
 - Prompt engineering for generative AI

How to Build a Dictionary

- 1. Identify "extreme" texts with known positions.
- E.g. opposition leader and PM, one-star and five-star reviews, etc.
- 2. Search for differentially occurring words using word frequencies.
- 3. Examine these words in context to assess their sensitivity and specificity.
- 4. Examine inflected forms to see whether stemming or wildcarding is required.
- 5. Use these words (or stems/lemmas) for categories.

Dictionary Performance



(González-Bailón & Paltoglou, 2015)

Dictionary vs Machine Learning

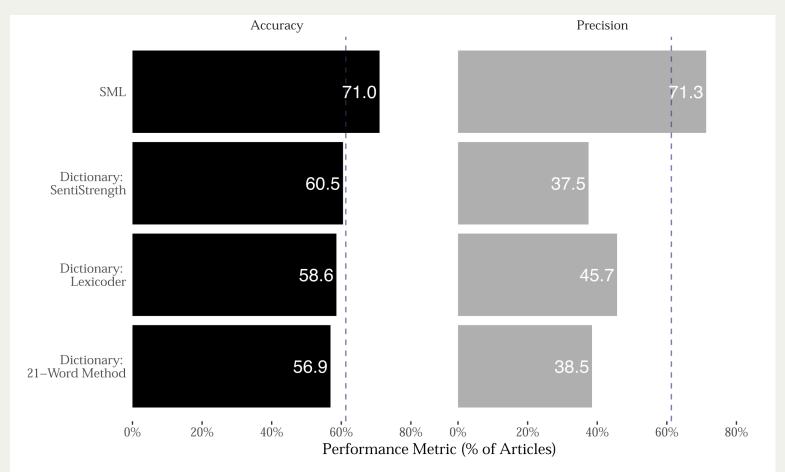


Figure 3. Performance of SML and Dictionary Classifiers—Accuracy and Precision. *Note:* Accuracy (percent correctly classified) and precision (percent of positive predictions that are correct) for the ground truth dataset coded by ten CrowdFlower coders. The dashed vertical lines indicate the baseline level of accuracy or precision (on any category) if the modal category is always predicted. The corpus used in the analysis is based on the keyword search of *The New York Times* 1980–2011 (see the text for details).

(Barberá, Boydstun, Linn, McMahon & Nagler, 2021)

Next

• Tutorial: Dictionaries and text classification

• Next week: Modelling texts