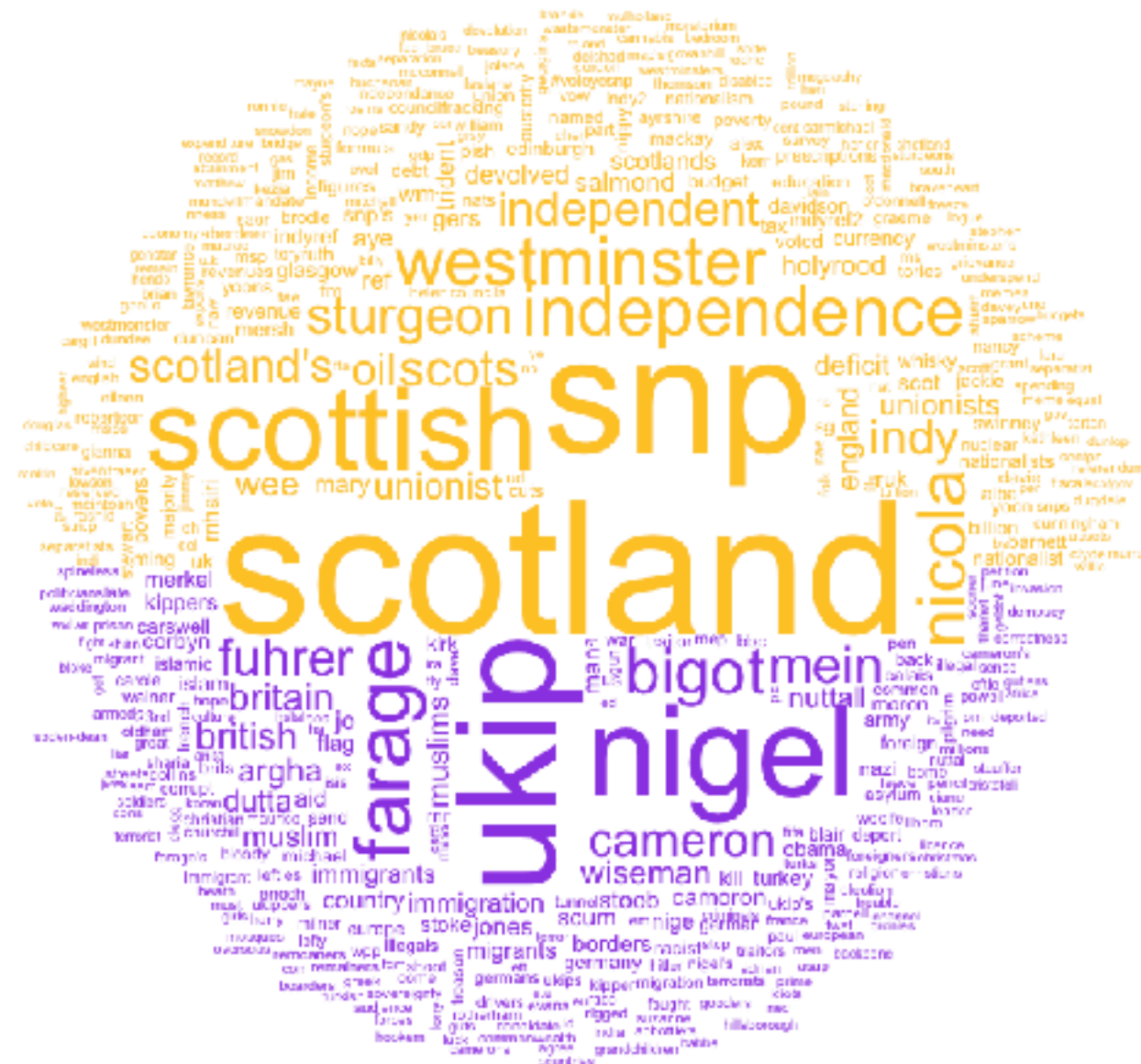


Basic Text Analysis using R

Justin Ho

What is Text Analysis?

SNP



UKIP

“Systematic, objective, quantitative analysis of message characteristics”

–Neuendorf, The Content Analysis Guidebook

Different Types of Text Analysis

- Degree of human involvement
 - Human coding (100%)
 - Supervised
 - Unsupervised (0%)
- Type of output
 - Scaling
 - Classification

4 Principles of Text Analysis

- Principle 1: All Quantitative Models of Language Are Wrong—But Some Are Useful
- Principle 2: Quantitative Methods Augment Humans, Not Replace Them
- Principle 3: There Is No Globally Best Method for Automated Text Analysis
- Principle 4: Validate, Validate, Validate

Read: Grimmer, J. and Steward, B. (2013). Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts. *Political Analysis*, 21(3), 267-297. doi:10.1093/pan/mps028

Bag of Words Assumption

- Word order doesn't matter
- The following three lines mean exactly the same:
 - I enjoy eating food and being with my family
 - I enjoy eating my family and being with food
 - and being eating enjoy family food I my with

End Product

Document-Feature Matrix

| | <i>word1</i> | <i>word2</i> | ... | <i>wordN</i> |
|-------------|--------------|--------------|-----|--------------|
| <i>doc1</i> | 0 | 0 | | 0 |
| <i>doc2</i> | 0 | 1 | | 0 |
| ... | | | | |
| <i>docN</i> | 2 | 0 | | 3 |

Meta Data

| | <i>var1</i> | <i>var2</i> | ... | <i>varN</i> |
|-------------|-------------|-------------|-----|-------------|
| <i>doc1</i> | Y | 1 | | 0 |
| <i>doc2</i> | N | 1 | | 7 |
| ... | | | | |
| <i>docN</i> | Y | 0 | | 10 |

Workflow of Text Analysis

1. Loading the text into R
2. Preprocess the Text
3. Analysing
4. Back to Step 2

Loading the Text into R: Encoding

- Common Encoding Standards: ASCII, Unicode (UTF-8/16)

100 00000j 00>J+ 0r00n0r0, 0000us0053Ak[y000wG• 0T 00*0 '00 0a^0rt000!!0Y070[000St(p00S0 uG 0fA 0?x050Q0000j000•50R00CV8g0dG00\$000y0#080d0K0yF0, %~0?f0~004`0{ 200E000T0 (0o0=00L000B0 'V0 q0000b0Y0\$| @H0=f0 p0T0+v0F00+L00j70+U0000Y0000#X8O0'y → 09}0Ep 0 T0 00=000j 0!!0wG00 0]z0xJ00+10*00•Bm0000j0i00s0)=(10X00y[000 FLs00@0X0h 0R0+0 00~→ M0 \$ 0J| 00|F0.0E0^000000 0 00050<0h00Tbg00 0kxf_^0c• 0 00-0v0 000I0000b0000 _0Z005:o "3q0{0T H9Z000•gX00N0+?0Ie0w0000b0 092FFe00+00 0 0000000/0i0 \$| r0&0006'04 0000K001Y0Y00 0~00F00 j0060 0+00\$0R6I0kBR0}"MPTy4p0000-000T0<!! 5T0 ?00 00?!!000kO0q0 4V0↑ 000)00V0\=%#000Ê0 00jD| 000020 0t0U0 0+~+ 0&0 0*0W0F .0z=00 00*i0 q!! 000y000+^000(00%4?kb0|0•N 0000-000 0^000[x0h|0000c 000 0Cal0 0%000(VE0 0*050 0s00\$00+0aC 0000a*0|000>0e0&0f000 0 bd\$0%0C0n0-Hm0+0;0mS-00M0 0 00D0.Gly0 0X0*0/s+r 0+200q+jQ 000!!~ao 00`0| [zgmrZ%60S Km0+05→ 0+0b| 0s/ 0*0000S0 0!!0D0\$000 "0Mc 0~002NX00 0000-&-0081rW0000u0&00=0000bY(000400 0m 06m0+0004eN0;o000-00b0!!000k0\$00uk00 00P<00000k00kt0 刺0G0,0p040>000u0@.000-0000*03A5&00iO!!0

Loading the Text into R: Extracting the Text

- Extracting useful Text:
 - `<h1 class="chapter0" id="c01">BRAN</h1><p class="nonindent">The morning had dawned clear and cold, with a crispness that hinted at the end of summer. They set forth at daybreak to see a man beheaded, twenty in all, and Bran rode among them, nervous with excitement. This was the first time he had been deemed old enough to go with his lord father and his brothers to see the king's justice done. It was the ninth year of summer, and the seventh of Bran's life.</p>`

Preprocessing

- Remove capitalization, numbers, and punctuation
- Before:
Article 1. All human beings are born free and equal in dignity and rights. They are endowed with reason and conscience and should act towards one another in a spirit of brotherhood.
- After:
article all human beings are born free and equal in dignity and rights they are endowed with reason and conscience and should act towards one another in a spirit of brotherhood

Preprocessing

- Remove stop words
- Before:
article all human beings are born free and equal in dignity and rights they are endowed with reason and conscience and should act towards one another in a spirit of brotherhood
- After:
article all human beings born free equal dignity rights endowed reason conscience act one another spirit brotherhood

Preprocessing

- Apply Stemming Algorithm
- Before:
article^e all human be^{ings} born free equal dignity^y rights^s endowed^{ed}
reason conscience^e act one another^{er} spirit brotherhood
- After:
articl all human be born free equal digniti right endow reason
conscienc act one anoth spirit brotherhood

Your turn.

<https://github.com/justinchuntingho>

Keyword Analysis

- What is a Keyword?
 - “A keyword may be defined as a word which occurs with unusual frequency in a given text. This does not mean high frequency but unusual frequency, by comparison with a reference corpus of some kind”

(Scott, M. (1997). PC analysis of key words - and key key words. *System*, 25(2), 233-45.)

Keyword Analysis

- What is Keyness?
 - “The keyness of a keyword represents the value of **log-likelihood** or **Chi-square statistics**; in other words it provides **an indicator of a keyword’s importance** as a content descriptor for the appeal.”

(Biber, D., Connor, U. & Upton, A. with Anthony, M. & Gladkov, K. (2007). Rhetorical appeals in fundraising. In D. Biber, U. Connor & A. Upton. *Discourse on the Move: Using corpus analysis to describe discourse structure*. (121-151). Amsterdam: John Benjamin.)

- What is a Chi-squared test?
 - Comparison between the **observed frequency** and **expected frequency**.

Keyword Analysis

- What is p-value?
 - “The significance (p-value) represents the probability that this keyness is accidental”.
 - The odds of observing the same observation if in a purely random world.

Your turn.

Just a few more things...

- Email: Justin.Ho@ed.ac.uk
- Twitter: [@zjustin334](https://twitter.com/zjustin334)
- Github: <https://github.com/justinchuntingho>
- Edinburgh Text Analysis Research Group
 - <https://jiscmail.ac.uk/TEXTANALYSIS>
 - <https://edtextanalysis.wordpress.com/>