



# **Text Analysis: The Basics**

# What is Text Analysis?

***“Systematic, objective, quantitative analysis of message characteristics”***

Kimberly A. Neuendorf, *The Content Analysis Guidebook*

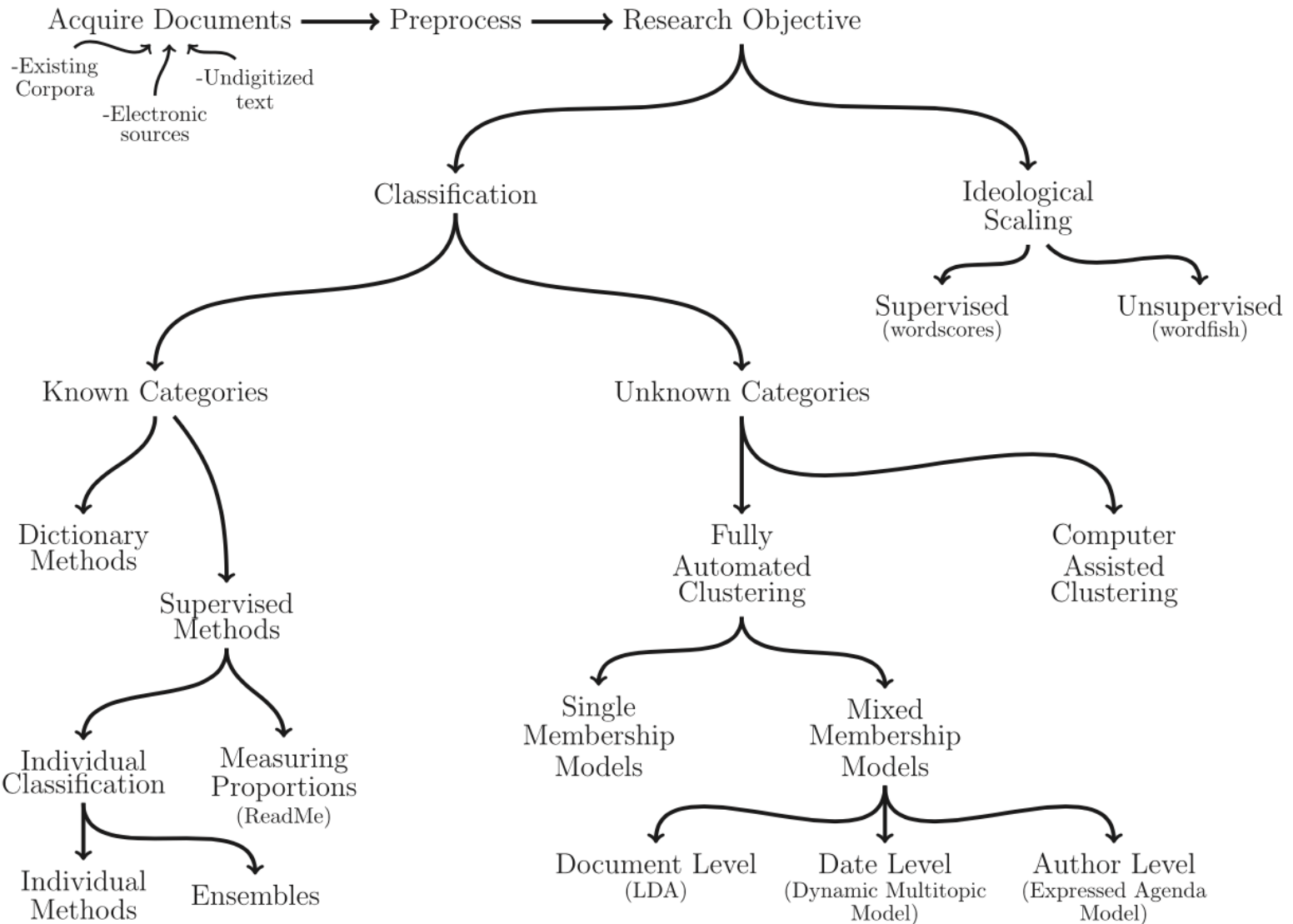
# Types of Text Analysis

## **Degree of human involvement:**

- Human coding (100%)
- Supervised
- Unsupervised (0%)

## **Type of Objective:**

- Scaling
- Classification



**Fig. 1** An overview of text as data methods.

(Source: Justin Grimmer and Brandon Stewart, 2013)

# Caution!

***"All Quantitative Models of Language Are Wrong — But Some Are Useful"***

Justin Grimmer and Brandon Stewart, *Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts*

# Bag of Words Assumption

- Word order doesn't matter
- The followings are 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



# The Pre-Processing

# Original Text

```
## <<PlainTextDocument>>  
## Metadata: 7  
## Content: chars: 181  
##  
## Article 1. All human beings are born free and equal in dignity
```

# Remove Punctuation

```
## <<PlainTextDocument>>  
## Metadata: 7  
## Content: chars: 178  
##  
## Article 1 All human beings are born free and equal in dignity
```

# To Lower Case

```
## <<PlainTextDocument>>  
## Metadata: 7  
## Content: chars: 178  
##  
## article 1 all human beings are born free and equal in dignity
```

# Remove Numbers

```
## <<PlainTextDocument>>  
## Metadata: 7  
## Content: chars: 177  
##  
## article all human beings are born free and equal in dignity a
```

# Remove Stopwords

```
## <<PlainTextDocument>>  
## Metadata: 7  
## Content: chars: 135  
##  
## article human beings born free equal dignity rights en
```

# Stemming

```
## <<PlainTextDocument>>  
## Metadata: 7  
## Content: chars: 108  
##  
## articl human be born free equal digniti right endow reason con
```

# Optional: Create N-Gram

```
##  
## article_1 1_all all_human human_beings beings_are are_born bor
```



# Create DFM

```
##      Terms
## Docs act  anoth  articl  born  brotherhood  conscienc  digniti  endow
##    1   1    1     1     1    1             1         1         1
##    2   0    0     1     0    0             0         0         0
```

# The Output

# 1. Document-Feature Matrix

	word1	word2	word3	...	wordN
doc1	0	2	0		3
doc2	1	0	0		0
doc3	0	0	2		1
...					
docN	0	1	0		0

## 2. Meta-data Martix

	var1	var2	var3	...	varN
doc1	Y	1	0		3
doc2	Y	0	0		2
doc3	N	0	0		1
...					
docN	Y	0	0		3

A photograph of a monkey sitting at a desk, typing on a silver laptop keyboard. The monkey is looking towards the camera with a neutral expression. The background is a warm, reddish-brown wall. A black horizontal bar is overlaid across the middle of the image, containing the text "Time for R" in white.

**Time for R**

# Keyword Analysis

# 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.”

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



# Keyword Analysis

## **What is a Chi-squared test?**

Comparison between the observed frequency and expected frequency.

	<i>love</i>	All other words	Total
Male	414	1714029	1714443
Female	1214	2592238	2593452
Total	1628	4306267	4307895

- Expected frequency: **Row total** times **Column total** divided by the **total number of words** in the corpus.
- Plug into this equation: 
$$X^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$
- Comparison between the **observed frequency** and **expected frequency** of the **word in question** and **all other words**.

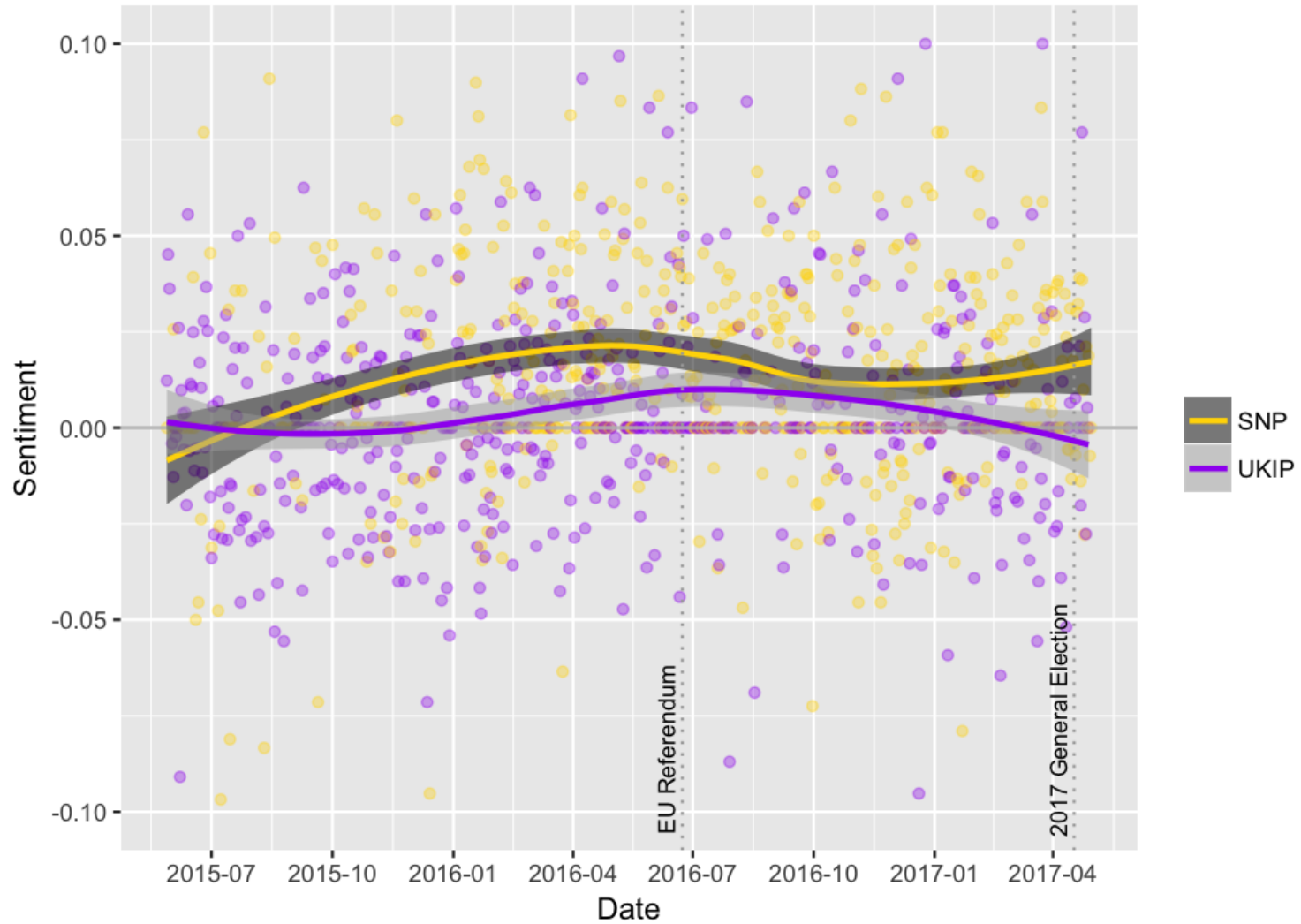
A photograph of a monkey sitting at a desk and typing on a laptop keyboard. The monkey is looking towards the camera with a neutral expression. The background is a warm, reddish-brown wall. A black horizontal bar is overlaid across the middle of the image, containing the text "Time for R" in white.

**Time for R**

# Other Applications

pre-cultural\_revolution  
international\_attention chinese\_government  
border\_regions organizational\_structure  
rely\_upon protest\_movement chiang\_kai-shek  
economic\_reforms economic\_development  
li\_hongzhi central\_committee chinese\_revolution  
chinese\_authorities  
red\_versus status\_quo chinese\_politics  
civil\_war  
general\_secretary great\_proletarian  
china\_prc soviet\_union hong\_kong took\_place three\_years  
civil\_society  
proletarian\_cultural chinese\_political zhao\_ziyang  
china\_today two\_lines  
political\_landscape party\_ccp people's\_republic china\_quarterly  
one\_hand student\_movement  
mao\_tse-tung falun\_gong united\_states  
mainland\_china people's\_liberation  
**cultural revolution**  
dalai\_lama communist\_party mao\_zedong  
religious\_movement red\_guard little\_impact  
inner\_mongolia chinese\_communist student\_protest  
chinese\_rule 20th\_century  
army\_pla human\_rights two\_decades  
red\_guards han\_chinese political\_system nobel\_prize  
social\_change deng\_xiaoping party\_congress  
political\_power tiananmen\_square taiwan's\_identity  
new\_york central\_asia public\_opinion state\_power  
contemporary\_chinese liberation\_army recent\_years  
jiang\_zemin autonomous\_region last\_decade  
political\_science important\_role turning\_point  
tiananmen\_massacre political\_scene versus\_expert  
article\_examined political\_institutions  
beijing\_government

# Sentiment of Facebook Posts



Selected Topic:

Next Topic

Selected Topic:

Slide to adjust relevance metric:(2)

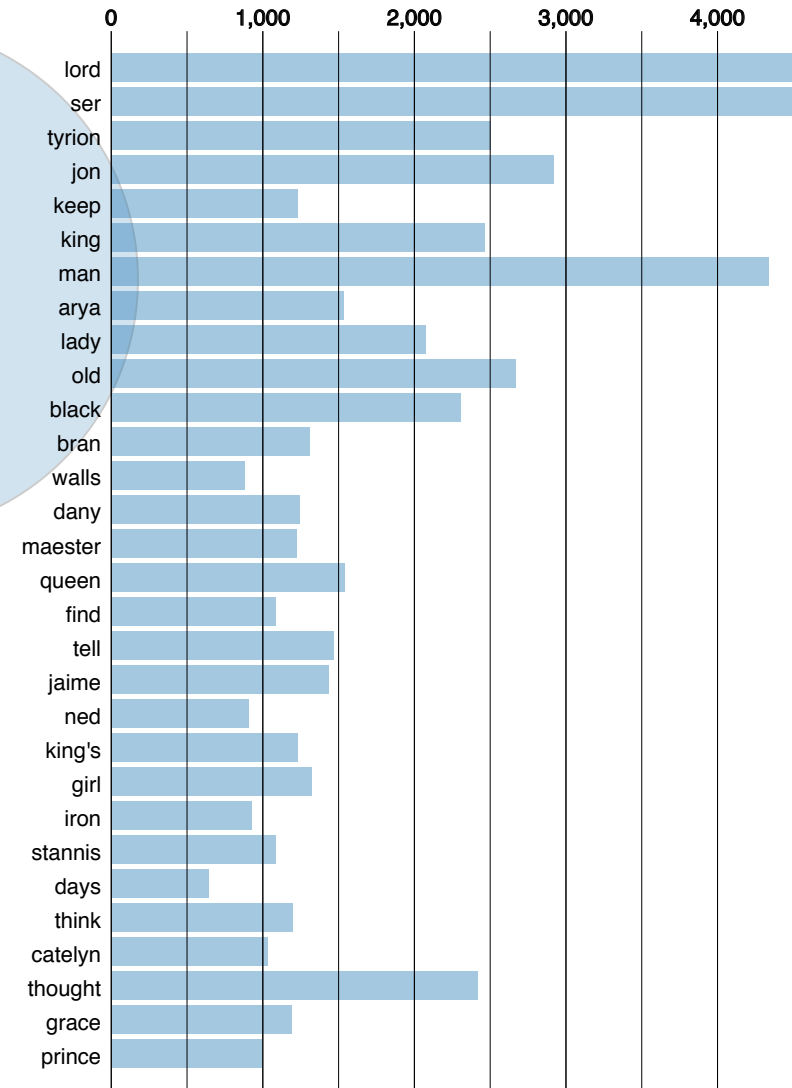
$$\lambda = 1$$

0.0

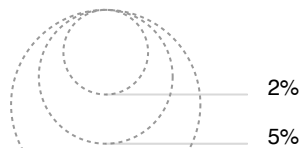
0.2

### Intertopic Distance Map (via multidimensional scaling)

### Top-30 Most Salient Term



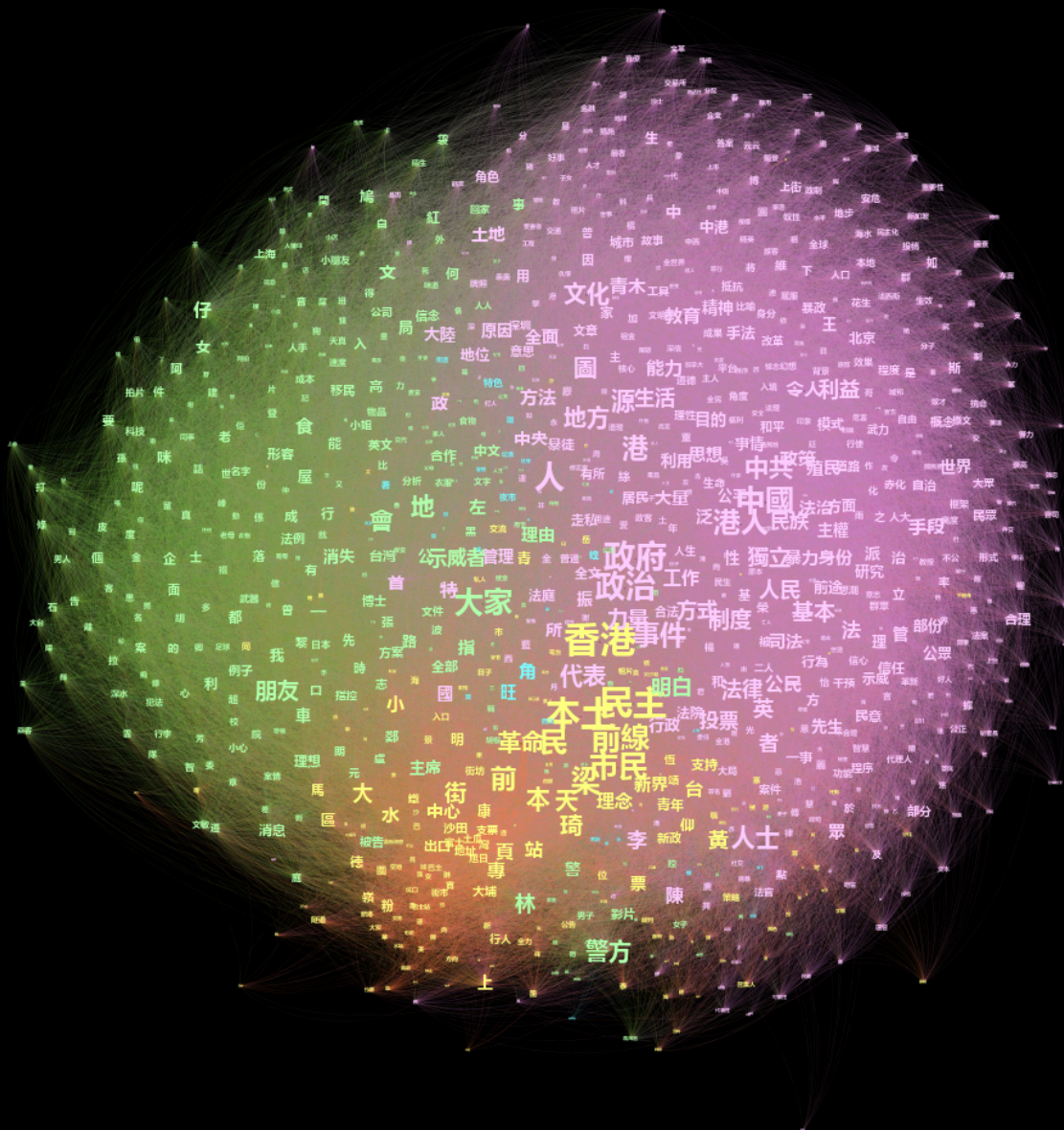
### Marginal topic distribution



Overall term frequency

Estimated term frequency within the selected topic

$$1. \text{saliency}(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$$





# Thank you very much!

Twitter: @justin\_ct\_ho

Github: justinchuntingho

Email: Jusitn.Ho@ed.ac.uk

Edinburgh Text Analysis Research Group:  
<https://jiscmail.ac.uk/TEXTANALYSIS>