

# Day 2: Descriptive statistical methods for textual analysis

Kenneth Benoit

TCD 2016: Quantitative Text Analysis

January 12, 2016

# Basic descriptive summaries of text

**Readability statistics** Use a combination of syllables and sentence length to indicate “readability” in terms of complexity

**Vocabulary diversity** (At its simplest) involves measuring a *type-to-token ratio* (TTR) where unique words are types and the total words are tokens

**Word (relative) frequency**

**Theme (relative) frequency**

**Length** in characters, words, lines, sentences, paragraphs, pages, sections, chapters, etc.

## Simple descriptive table about texts: Describe your data!

Speaker	Party	Tokens	Types
Brian Cowen	FF	5,842	1,466
Brian Lenihan	FF	7,737	1,644
Ciaran Cuffe	Green	1,141	421
John Gormley (Edited)	Green	919	361
John Gormley (Full)	Green	2,998	868
Eamon Ryan	Green	1,513	481
Richard Bruton	FG	4,043	947
Enda Kenny	FG	3,863	1,055
Kieran O'Donnell	FG	2,054	609
Joan Burton	LAB	5,728	1,471
Eamon Gilmore	LAB	3,780	1,082
Michael Higgins	LAB	1,139	437
Ruairi Quinn	LAB	1,182	413
Arthur Morgan	SF	6,448	1,452
Caoimhghin O'Caolain	SF	3,629	1,035
All Texts		49,019	4,840
<i>Min</i>		919	361
<i>Max</i>		7,737	1,644
<i>Median</i>		3,704	991
<i>Hapaxes with Gormley Edited</i>		67	
<i>Hapaxes with Gormley Full Speech</i>		69	

# Lexical Diversity

- ▶ Basic measure is the **TTR**: Type-to-Token ratio
- ▶ Problem: This is very sensitive to overall document length, as shorter texts may exhibit fewer word repetitions
- ▶ Special problem: length may relate to the introduction of additional subjects, which will also increase richness

# Lexical Diversity: Alternatives to TTRs

$$\text{TTR} \frac{\text{total types}}{\text{total tokens}}$$

$$\text{Guiraud} \frac{\text{total types}}{\sqrt{\text{total tokens}}}$$

**D** (Malvern et al 2004) Randomly sample a fixed number of tokens and count those

**MTLD** the mean length of sequential word strings in a text that maintain a given TTR value (McCarthy and Jarvis, 2010) – fixes the TTR at 0.72 and counts the length of the text required to achieve it

# Vocabulary diversity and corpus length

- ▶ In natural language text, the rate at which new types appear is very high at first, but diminishes with added tokens

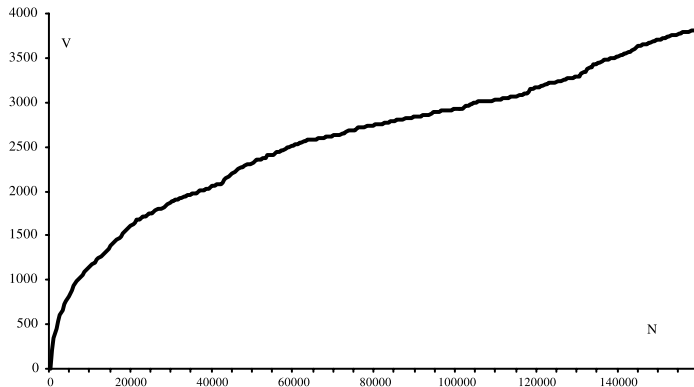


Fig. 1. Chart of vocabulary growth in the tragedies of Racine (chronological order, 500 token intervals).

# Vocabulary Diversity Example

- ▶ Variations use automated segmentation – here approximately 500 words in a corpus of serialized, concatenated weekly addresses by de Gaulle (from Labbé et. al. 2004)
- ▶ While most were written, during the period of December 1965 these were more spontaneous press conferences

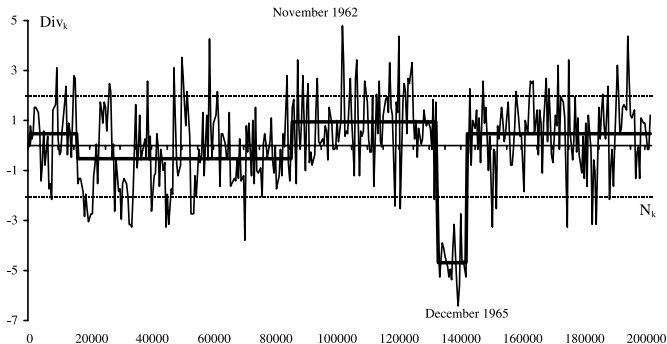


Fig. 8. Evolution of vocabulary diversity in General de Gaulle's broadcast speeches (June 1958–April 1969).

## Complexity and Readability

- ▶ Use a combination of syllables and sentence length to indicate “readability” in terms of complexity
- ▶ Common in educational research, but could also be used to describe textual complexity
- ▶ Most use some sort of sample
- ▶ No natural scale, so most are calibrated in terms of some interpretable metric
- ▶ Not (yet) implemented in `quanteda`, but available from `koRpus` package



## Flesch-Kincaid readability index

- ▶ F-K is a modification of the original **Flesch Reading Ease Index**:

$$206.835 - 1.015 \left( \frac{\text{total words}}{\text{total sentences}} \right) - 84.6 \left( \frac{\text{total syllables}}{\text{total words}} \right)$$

**Interpretation:** 0-30: university level; 60-70: understandable by 13-15 year olds; and 90-100 easily understood by an 11-year old student.

- ▶ **Flesch-Kincaid** rescales to the US educational grade levels (1-12):

$$0.39 \left( \frac{\text{total words}}{\text{total sentences}} \right) + 11.8 \left( \frac{\text{total syllables}}{\text{total words}} \right) - 15.59$$

## Gunning fog index

- ▶ Measures the readability in terms of the years of formal education required for a person to easily understand the text on first reading
- ▶ Usually taken on a sample of around 100 words, not omitting any sentences or words
- ▶ Formula:

$$0.4 \left[ \left( \frac{\text{total words}}{\text{total sentences}} \right) + 100 \left( \frac{\text{complex words}}{\text{total words}} \right) \right]$$

where complex words are defined as those having three or more syllables, not including proper nouns (for example, Ljubljana), familiar jargon or compound words, or counting common suffixes such as -es, -ed, or -ing as a syllable

## Documents as vectors

- ▶ The idea is that (weighted) features form a vector for each document, and that these vectors can be judged using metrics of **similarity**
- ▶ A document's vector for us is simply (for us) the row of the document-feature matrix

## Characteristics of similarity measures

Let  $A$  and  $B$  be any two documents in a set and  $d(A, B)$  be the distance between  $A$  and  $B$ .

1.  $d(x, y) \geq 0$  (the distance between any two points must be non-negative)
2.  $d(A, B) = 0$  iff  $A = B$  (the distance between two documents must be zero if and only if the two objects are identical)
3.  $d(A, B) = d(B, A)$  (distance must be symmetric:  $A$  to  $B$  is the same distance as from  $B$  to  $A$ )
4.  $d(A, C) \leq d(A, B) + d(B, C)$  (the measure must satisfy the triangle inequality)

## Euclidean distance

Between document  $A$  and  $B$  where  $j$  indexes their features, where  $y_{ij}$  is the value for feature  $j$  of document  $i$

- ▶ Euclidean distance is based on the Pythagorean theorem
- ▶ Formula

$$\sqrt{\sum_{j=1}^j (y_{Aj} - y_{Bj})^2} \quad (1)$$

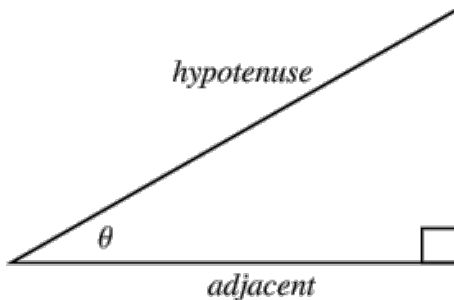
- ▶ In vector notation:

$$\|\mathbf{y}_A - \mathbf{y}_B\| \quad (2)$$

- ▶ Can be performed for any number of features  $J$  (or  $V$  as the vocabulary size is sometimes called – the number of columns in of the dfm, same as the number of feature types in the corpus)

## A geometric interpretation of “distance”

In a right angled triangle, the cosine of an angle  $\theta$  or  $\cos(\theta)$  is the **length of the adjacent side** divided by the **length of the hypotenuse**



We can use the vectors to represent the text location in a  $V$ -dimensional vector space and compute the angles between them

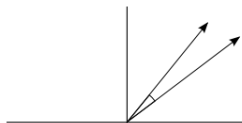
## Cosine similarity

- ▶ Cosine distance is based on the size of the angle between the vectors
- ▶ Formula

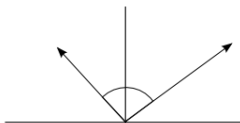
$$\frac{\mathbf{y}_A \cdot \mathbf{y}_B}{\|\mathbf{y}_A\| \|\mathbf{y}_B\|} \quad (3)$$

- ▶ The  $\cdot$  operator is the dot product, or  $\sum_j y_{Aj} y_{Bj}$
- ▶ The  $\|\mathbf{y}_A\|$  is the vector norm of the (vector of) features vector  $\mathbf{y}$  for document  $A$ , such that  $\|\mathbf{y}_A\| = \sqrt{\sum_j y_{Aj}^2}$
- ▶ Nice property: independent of document length, because it deals only with the angle of the vectors
- ▶ Ranges from -1.0 to 1.0 for term frequencies, or 0 to 1.0 for normalized term frequencies (or tf-idf)

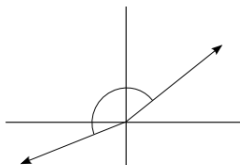
# Cosine similarity illustrated



Similar scores  
Score Vectors in same direction  
Angle between them is near 0 deg.  
Cosine of angle is near 1 i.e. 100%



Unrelated scores  
Score Vectors are nearly orthogonal  
Angle between them is near 90 deg.  
Cosine of angle is near 0 i.e. 0%



Opposite scores  
Score Vectors in opposite direction  
Angle between them is near 180 deg.  
Cosine of angle is near -1 i.e. -100%



## Example text

**Hurricane Gilbert** swept toward the Dominican Republic Sunday , and the Civil Defense alerted its heavily populated south coast to prepare for high **winds**, heavy **rains** and high seas.

The **storm** was approaching from the southeast with sustained **winds** of 75 mph gusting to 92 mph .

"There is no need for alarm," Civil Defense Director Eugenio Cabral said in a television alert shortly before midnight Saturday .

Cabral said residents of the province of Barahona should closely follow **Gilbert**'s movement .

An estimated 100,000 people live in the province, including 70,000 in the city of Barahona , about 125 miles west of Santo Domingo .

Tropical **Storm Gilbert** formed in the eastern Caribbean and strengthened into a **hurricane** Saturday night

The National **Hurricane** Center in Miami reported its position at 2a.m. Sunday at latitude 16.1 north , longitude 67.5 west, about 140 miles south of Ponce, Puerto Rico, and 200 miles southeast of Santo Domingo.

The National Weather Service in San Juan , Puerto Rico , said **Gilbert** was moving westward at 15 mph with a "broad area of cloudiness and heavy weather" rotating around the center of the **storm**.

The weather service issued a flash flood watch for Puerto Rico and the Virgin Islands until at least 6p.m. Sunday.

Strong **winds** associated with the **Gilbert** brought coastal flooding , strong southeast **winds** and up to 12 feet to Puerto Rico 's south coast.

## Example text: selected terms

- ▶ Document 1

Gilbert: 3, hurricane: 2, rains: 1, storm: 2, winds: 2

- ▶ Document 2

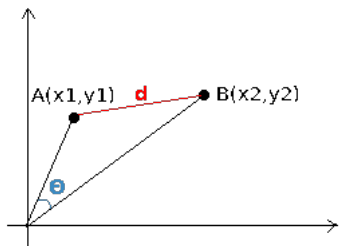
Gilbert: 2, hurricane: 1, rains: 0, storm: 1, winds: 2

## Example text: cosine similarity in R

```
> toyDfm <- matrix(c(3,2,1,2,2, 2,1,0,1,2), nrow=2, byrow=TRUE)
> colnames(toyDfm) <- c("Gilbert", "hurricane", "rain", "storm", "winds")
> rownames(toyDfm) <- c("doc1", "doc2")
> toyDfm
      Gilbert hurricane rain storm winds
doc1      3          2   1    2     2
doc2      2          1   0    1     2
> simil(toyDfm, "cosine")
      doc1
doc2 0.9438798
```

## Relationship to Euclidean distance

- ▶ Cosine similarity measures the similarity of vectors with respect to the origin
- ▶ Euclidean distance measures the distance between particular points of interest along the vector



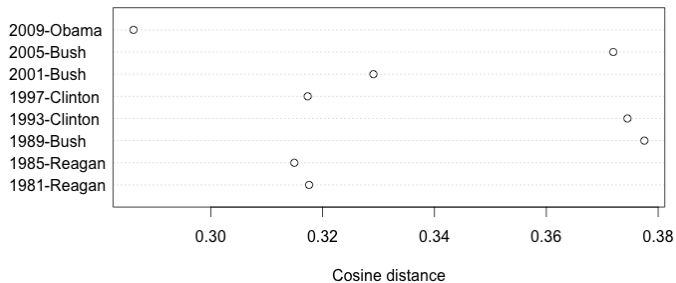
# Jacquard coefficient

- ▶ Similar to the Cosine similarity
- ▶ Formula

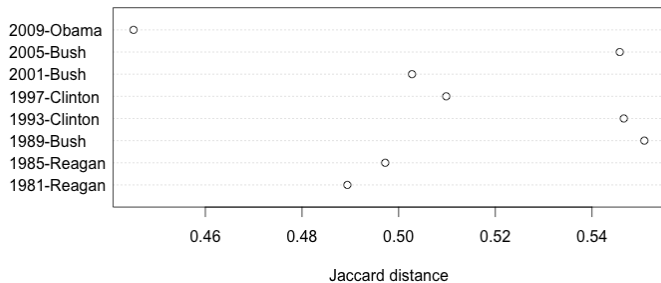
$$\frac{\mathbf{y}_A \cdot \mathbf{y}_B}{\|\mathbf{y}_A\| + \|\mathbf{y}_B\| - \mathbf{y}_A \cdot \mathbf{y}_B} \quad (4)$$

- ▶ Ranges from 0 to 1.0

## Example: Inaugural speeches



## Example: Inaugural speeches



## Can be made very general for binary features

Example: In the Choi et al paper, they compare vectors of features for (binary) absence or presence – called (“operational taxonomic

**Table 1** OTUs Expression of Binary Instances  $i$  and  $j$

$j \backslash i$	1 (Presence)	0 (Absence)	Sum
1 (Presence)	$a = i \cdot j$	$b = \bar{i} \cdot j$	$a+b$
0 (Absence)	$c = i \cdot \bar{j}$	$d = \bar{i} \cdot \bar{j}$	$c+d$
Sum	$a+c$	$b+d$	$n=a+b+c+d$

units”)

- ▶ Cosine similarity:

$$S_{\text{cosine}} = \frac{a}{\sqrt{(a+b)(a+c)}} \quad (5)$$

- ▶ Jaccard similarity:

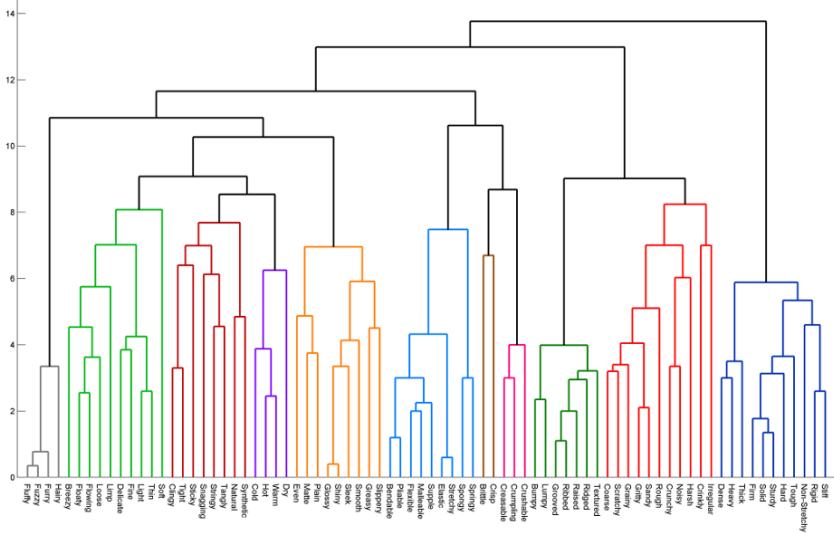
$$S_{\text{Jaccard}} = \frac{a}{\sqrt{(a+b+c)}} \quad (6)$$



## Typical features

- ▶ Normalized term frequency (almost certainly)
- ▶ Very common to use tf-idf – if not, similarity is boosted by common words (stop words)
- ▶ Not as common to use binary features

# Uses for similarity measures: Clustering



## Other uses, extensions

- ▶ Used extensively in information retrieval
- ▶ Summary measures of how far apart two texts are – but be careful exactly how you define “features”
- ▶ Some but not many applications in social sciences to measure substantive similarity — scaling models are generally preferred
- ▶ Can be used to generalize or represent features in machine learning, by combining features using kernel methods to compute similarities between textual (sub)sequences without extracting the features explicitly (as we have done here)

## Edit distances

- ▶ Edit distance refers to the number of operations required to transform one string into another for strings of equal length
- ▶ Common edit distance: the **Levenshtein distance**
- ▶ Example: the Levenshtein distance between "kitten" and "sitting" is 3
  - ▶ kitten → sitten (substitution of "s" for "k")
  - ▶ sitten → sittin (substitution of "i" for "e")
  - ▶ sittin → sitting (insertion of "g" at the end).
- ▶ Hamming distance: for two strings of equal length, the Hamming distance is the number of positions at which the corresponding characters are different
- ▶ Not common, as at a textual level this is hard to implement and possibly meaningless

## Sampling issues in existing measures

- ▶ Lexical diversity measures may take sample frames, or moving windows, and average across the windows
- ▶ Readability may take a sample, or multiple samples, to compute readability measures
- ▶ But rather than simulating the “sampling distribution” of a statistic, these are more designed to:
  - ▶ get a representative value for the text as a whole
  - ▶ normalize the length of the text relative to other texts

## Sampling illustrated

- ▶ lexical diversity
- ▶ dictionaries (feature counts)  
Can construct multiple `dfm` objects, and count the “dictionary” features across texts. (Replicate the populism dictionary)

# Bootstrapping text-based statistics



# Simulation and bootstrapping

Used for:

- ▶ Gaining **intuition** about distributions and sampling
- ▶ Providing **distributional** information not distributions are not directly known, or cannot be assumed
- ▶ Acquiring **uncertainty** estimates

Both simulation and bootstrapping are **numerical approximations** of the quantities we are interested in. (Run the same code twice, and you get different answers)

Solution for replication: save the **seed**



# Quantifying Uncertainty

- ▶ Critical if we really want to compare texts
- ▶ Question: How?
  - ▶ Make parametric assumptions about the data-generating process. For instance, we could model feature counts according to a Poisson distribution.
  - ▶ Use a sampling procedure and obtain averages from the samples. For instance we could sample 100-word sequences, compute reliability, and look at the spread of the readability measures from the samples
  - ▶ Bootstrapping: a generalized resampling method

# Bootstrapping

- ▶ *Bootstrapping* refers to repeated resampling of data points **with replacement**
- ▶ Used to estimate the error variance (i.e. the **standard error**) of an estimate when the sampling distribution is unknown (or cannot be safely assumed)
- ▶ Robust in the absence of parametric assumptions
- ▶ Useful for some quantities for which there is no known sampling distribution, such as computing the standard error of a median

## Bootstrapping illustrated

```
> ## illustrate bootstrap sampling
> set.seed(30092014) # set the seed so that your results will match m
> # using sample to generate a permutation of the sequence 1:10
> sample(10)
[1] 4 2 1 9 8 5 7 3 6 10
> # bootstrap sample from the same sequence
> sample(10, replace=T)
[1] 8 6 6 2 5 8 4 8 4 9
> # bootstrap sample from the same sequence with probabilities that
> # favor the numbers 1-5
> prob1 <- c(rep(.15, 5), rep(.05, 5))
> prob1
[1] 0.15 0.15 0.15 0.15 0.15 0.05 0.05 0.05 0.05 0.05
> sample(10, replace=T, prob=prob1)
[1] 4 1 1 2 8 3 1 6 1 9
```

# Bootstrapping the standard error of the median

Using a user-defined function:

```
b.median <- function(data, n) {  
  resamples <- lapply(1:n, function(i) sample(data, replace=T))  
  sapply(resamples, median)  
  std.err <- sqrt(var(r.median))  
  list(std.err=std.err, resamples=resamples, medians=r.median)  
}  
summary(b.median(spending, 10))  
summary(b.median(spending, 100))  
summary(b.median(spending, 400))  
median(spending)
```

# Bootstrapping the standard error of the median

Using R's **boot** library:

```
library(boot)
samplemedian <- function(x, d) return(median(x[d]))
quantile(boot(spending, samplemedian, R=10)$t, c(.025, .5, .975))
quantile(boot(spending, samplemedian, R=100)$t, c(.025, .5, .975))
quantile(boot(spending, samplemedian, R=400)$t, c(.025, .5, .975))
```

**Note:** There is a good reference on using `boot()` from <http://www.mayin.org/ajayshah/KB/R/documents/boot.html>

# Bootstrapping methods for textual data

- ▶ Question: what is the "sampling distribution" of a text-based statistic? Examples:
  - ▶ a term's (relative) frequency
  - ▶ lexical diversity
  - ▶ complexity
- ▶ Could use to compare subgroups, analagous to ANOVA or t-tests, for statistics such as similarity (below)

## Guidelines for bootstrapping text

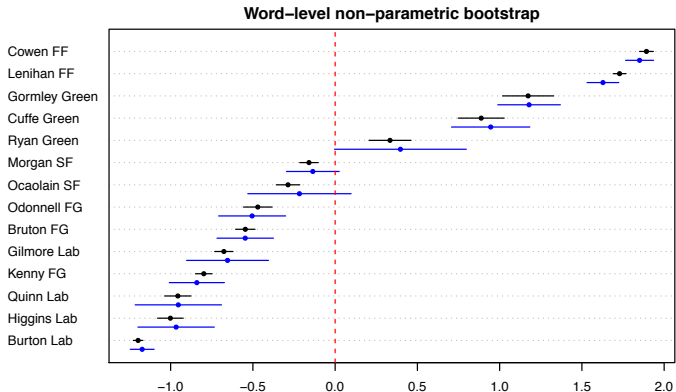
- ▶ Bootstrap by resampling tokens.  
Advantage: This is easily done from the document-feature matrix.  
Disadvantage: Ignores the natural units into which text is grouped, such as sentences
- ▶ Bootstrap by resampling sentences.  
Advantage: Produces more meaningful (potentially readable) texts, more faithful to data-generating process.  
Disadvantage: More complicated, cannot be done from dfm, must segment the text into sentences and construct a new dfm for each resample.

## Guidelines for bootstrapping text (cont)

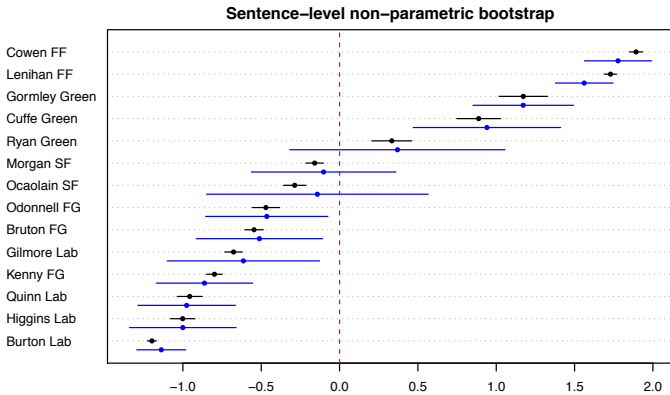
- ▶ Other options for bootstrapping text: generalize the notion of the “block” bootstrap. In block bootstrap, consecutive blocks of observations of length  $K$  are resampled from the original time series, either in fixed blocks (Carlstein, 1986) or overlapping blocks (Künsch, 1989)
  - ▶ paragraphs
  - ▶ pages
  - ▶ chapters
  - ▶ stratified: words within sentences or paragraphs



# Different bootstrapping methods: example



# Different bootstrapping methods: example



# Different bootstrapping methods: example

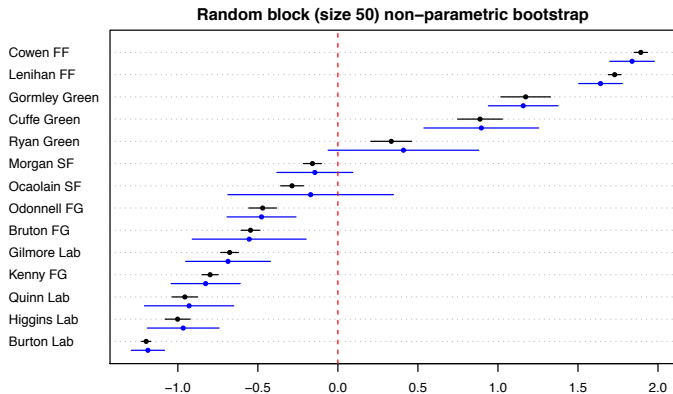


Figure 3: *Estimates of  $\theta_i$  and 95% confidence intervals from the 2009 Irish Budget Debates using Block Bootstrapping.* Black points and lines are analytical SEs; blue point estimates and 95% CIs correspond to the labelled method.