#### POIR 613: Computational Social Science

#### Pablo Barberá

University of Southern California pablobarbera.com

Course website: pablobarbera.com/POIR613/

# **Dictionary methods**

### Outline for today

- Dictionary methods: an overview
- Some well-known dictionaries
- Advantages and disadvantages
- Dictionary construction

### **Dictionary methods**

Classifying documents when categories are known:

- Lists of words that correspond to each category:
  - Positive or negative, for sentiment
  - Sad, happy, angry, anxious... for emotions
  - Insight, causation, discrepancy, tentative... for cognitive processes
  - Sexism, homophobia, xenophobia, racism... for hate speech

many others: see LIWC, VADER, SentiStrength, LexiCoder...

- Count number of times they appear in each document
- Normalize by document length (optional)
- ► Validate, validate, validate.
  - Check sensitivity of results to exclusion of specific words
  - Code a few documents manually and see if dictionary prediction aligns with human coding of document

Bridging qualitative and quantitative text analysis

- A hybrid procedure between qualitative and quantitative classification at the fully automated end of the text analysis spectrum
- "Qualitative" since it involves identification of the concepts and associated keys/categories, and the textual features associated with each key/category
- Dictionary construction involves a lot of contextual interpretation and qualitative judgment
- Perfect reliability because there is no human decision making as part of the text analysis procedure

#### Outline for today

- Dictionary methods: an overview
- Some well-known dictionaries
- Advantages and disadvantages
- Dictionary construction

#### Well-known dictionaries: General Inquirer

- General Inquirer (Stone et al 1966)
- Example: self = I, me, my, mine, myself selves = we, us, our, ours, ourselves
- Latest version contains 182 categories the "Harvard IV-4" dictionary, the "Lasswell" dictionary, and five categories based on the social cognition work of Semin and Fiedler
- Examples: "self references", containing mostly pronouns; "negatives", the largest category with 2291 entries

### Linquistic Inquiry and Word Count

- Created by Pennebaker et al see http://www.liwc.net
- Uses a dictionary to calculate the percentage of words in the text that match each of up to 82 language dimensions
- Consists of about 4,500 words and word stems, each defining one or more word categories or subdictionaries
- For example, the word *cried* is part of five word categories: sadness, negative emotion, overall affect, verb, and past tense verb. So observing the token *cried* causes each of these five subdictionary scale scores to be incremented
- Hierarchical: so "anger" words are part of an emotion category and a negative emotion subcategory
- You can buy it here:

http://www.liwc.net/descriptiontable1.php

## Example: Emotional Contagion on Facebook



#### VADER: an open-source alternative to LIWC

#### Valence Aware Dictionary and sEntiment Reasoner:

- Especially tuned for social media text
- Captures polarity and intensity of sentiments
- Includes emoticons, emoji, slang
- Feature-specific weights
- Python and R libraries:

https://github.com/cjhutto/vaderSentiment

Other open-source sentiment dictionaries: LexiCoder (media text), SentiStrength (social media text)

## Example: Laver and Garry (2000)

- A hierarchical set of categories to distinguish policy domains and policy positions – similar in spirit to the CMP
- Five domains at the top level of hierarchy
  - economy
  - political system
  - social system
  - external relations
  - a "general domain that has to do with the cut and thurst of specific party competition as well as uncodable pap and waffle"
- Looked for word occurrences within "word strings with an average length of ten words"
- Built the dictionary on a set of specific UK manifestos

#### Example: Laver and Garry (2000): Economy

TABLE 1	Abridged Section of Revised Manifesto Coding Scheme					
1 ECONOMY Role of state	in econor	2014				
noie of state	In econor	ny				
Increa	ase role o	ate+ f state				
111	ECONOMY/+State+/Budget Budget					
	1111	ECONOMY/+State+/Budget/Spending Increase public spending				
		11111	ECONOMY/+State+/Budget/Spending/Health			
		11112	ECONOMY/+State+/Budget/Spending/Educ. and training			
		11113	ECONOMY/+State+/Budget/Spending/Housing			
		11114	ECONOMY/+State+/Budget/Spending/Transport			
		11115	ECONOMY/+State+/Budget/Spending/Infrastructure			
		11116	ECONOMY/+State+/Budget/Spending/Welfare			
		11117	ECONOMY/+State+/Budget/Spending/Police			
		11118	ECONOMY/+State+/Budget/Spending/Defense			
		11119	ECONOMY/+State+/Budget/Spending/Culture			
	1112	ECONOMY/+State+/Budget/Taxes Increase taxes				
		11121	ECONOMY/+State+/Budget/Taxes/Income			
		11122	ECONOMY/+State+/Budget/Taxes/Payroll			
		11123	ECONOMY/+State+/Budget/Taxes/Company			
		11124	ECONOMY/+State+/Budget/Taxes/Sales			
		11125	ECONOMY/+State+/Budget/Taxes/Capital			
		11126	ECONOMY/+State+/Budget/Taxes/Capital gains			
	1113	ECONOMY/+State+/Budget/Deficit Increase budget deficit				
		11131	ECONOMY/+State+/Budget/Deficit/Borrow			
		11132	ECONOMY/+State+/Budget/Deficit/Inflation			

# MFD (Graham and Haidt)

#### Moral Foundations dictionary:

- Moral foundations: dimensions of difference that explain human moral reasoning
- Measures the proportions of virtue and vice words for each foundation:
  - 1. Care/Harm
  - 2. Fairness/Cheating
  - 3. Loyalty/Betrayal
  - 4. Authority/Subversion
  - 5. Purity/Degradation

Link to dictionary. https:

//moralfoundations.org/other-materials/

#### Outline for today

- Dictionary methods: an overview
- Some well-known dictionaries
- Advantages and disadvantages
- Dictionary construction

### Potential advantage: Multi-lingual

	NL	UK	GE	IT
Core	elit*	elit*	elit*	elit*
	consensus*	consensus*	konsens*	consens*
	ondemocratisch* ondemokratisch*	undemocratic*	undemokratisch*	antidemocratic*
	referend*	referend*	referend*	referend*
	corrupt*	corrupt*	korrupt*	corrot*
	propagand*	propagand*	propagand*	propagand*
	politici*	politici*	politiker*	politici*
	*bedrog*	*deceit*	täusch*	ingann*
	*bedrieg*	*deceiv*	betrüg*	
			betrug*	
	*verraa*	*betray*	*verrat*	tradi*
	*verrad*			
	schaam*	shame*	scham* schäm*	vergogn*
	schand*	scandal*	skandal*	scandal*
	waarheid*	truth*	wahrheit*	verità
	oneerlijk*	dishonest*	unfair* unehrlich*	disonest*
Context	establishm* heersend*	establishm*	establishm* *herrsch*	partitocrazia
	capitul*	runng	nerrsen	
	kapitul*			
	kaste*			
	leugen* lieg*		lüge*	menzogn* mentir*

#### APPENDIX B DICTIONARY OF THE COMPUTER-BASED CONTENT ANALYSIS

(from Rooduijn and Pauwels 2011)

### Potential disadvantage: Context specific



Lexicons' Accuracy in Document Classification Compared to Machine-Learning Approach

Source: González-Bailón and Paltoglou (2015)

## Disadvantage: Highly specific to context

- Example: Loughran and McDonald used the Harvard-IV-4 TagNeg (H4N) file to classify sentiment for a corpus of 50,115 firm-year 10-K filings from 1994–2008
- found that almost three-fourths of the "negative" words of H4N were typically not negative in a financial context e.g. *mine* or *cancer*, or *tax*, *cost*, *capital*, *board*, *liability*, *foreign*, and *vice*
- Problem: polysemes words that have multiple meanings
- Another problem: dictionary lacked important negative financial words, such as *felony*, *litigation*, *restated*, *misstatement*, and *unanticipated*

#### Potential disadvantage: sensitive to frequent words



Fig. 1. The timeline of sadness, anxiety, and anger on September 11 as expressed in messages sent to text pagers. Each data point represents the mean percentage of words related to the specific negative emotion, averaged across 30 min. The time slots start at 6:45 a.m. to 7:14 a.m. on September 11, 2001, and end at 12:15 a.m. to 12:44 a.m. on September 12, 2001. Exact times and brief descriptions of the most important events of September 11 are included above the timelines. WTC = World Trade Center

(from Back et al, Psychological Science, 2010)

#### Potential disadvantage: sensitive to frequent words

# Automation can lead to confounds in text analysis: Back, Küfner, and Egloff (2010) and the not-so-angry Americans.

🖙 EXPORT 🖈 Add To My List 🔛 🖶 <

Database: PsycINFO Comment/ Reply

#### Pury, Cynthia L. S.

Citation

Pury, C. L. S. (2011). Automation can lead to confounds in text analysis: Back, Küfner, and Egloff (2010) and the not-so-angry Americans. *Psychological Science*, 22(6), 835-836. <u>http://dx.doi.org/10.1177/056778711408735</u>

#### Abstract

Comments on an article by Milja D. Back et al. (see record 2010-25035-010). The authors used Linguistic Inquiry and Word Count (LWC) to analyze pager messages sent to more than 85,000 American pagers on September 11, 2001. They found that anger, as indexed by the words contained in those messages, rose steadily throughout the day. The data contained many technical codes; thus, Back et al. counted only words recognized by LIWC. However, this procedure did not exclude automatically generated messages. Consequently, LIWC words in such messages were counted, even if the words lacked emotional meaning in context. Furthermore, computers can send messages with superhuman frequency, turning an otherwise minor measurement error into a serious confound. This confound can be detected by treating individual text messages as primary units, reading samples of each key word in context, and looking for repeating false positives. Thus, it appears that much of the dramatic rise in anger reported by Back et al. was due to a repeated and emotionally neutral technical message associated with a single pager. Beccause today's e-mail, social media, and text messages can include automatically generated messages, future researchers of linguistic archives should consider ways to prevent similar confounds. (PsycINFO Database Record (c) 2016 APA, all rights reserved)

#### Potential disadvantage: sensitive to frequent words



Fig. 1. A revised timeline of anger as expressed in 37,606 social messages sent to text pagers on September 11, 2001. The graphs show (a) the mean percentage of words related to anger (as classified by Linguistic Inquiry and Word Count; Pennebaker, Francis, & Booth, 2001) and (b) the mean anger rating (0 = no onger, 1 = some onger, 2 = stong onger averaged across three raters for each message) across time slots staring at 6:45 a.m. to 7:14 a.m. on September 11, 2001, and ending at 12:15 a.m. to 12:44 a.m. on September 12, 2001.

(from Back et al, Psychological Science, 2011)

#### Outline for today

- Dictionary methods: an overview
- Some well-known dictionaries
- Advantages and disadvantages
- Dictionary construction

#### How to build a dictionary

- The ideal content analysis dictionary associates all and only the relevant words to each category in a perfectly valid scheme
- Three key issues:
  - ValidityIs the dictionary's category scheme valid?RecallDoes this dictionary identify *all* my content?PrecisionDoes it identify *only* my content?
- Imagine two logical extremes of including all words (too sensitive), or just one word (too specific)

#### How to build a dictionary

- 1. Identify "extreme texts" with "known" positions. Examples:
  - Tweets by populist vs mainstream parties (for populism dictionary)
  - Opposition leader and Prime Minister in a no-confidence debate (for opposition vs government dictionary)
  - Facebook comments to news about natural catastrophes vs football victories (for sentiment dictionary)
  - Subreddits for white nationalist groups vs regular politics (for racist rhetoric)
- 2. Search for differentially occurring words using word frequencies
- 3. Examine these words in context to check their precision and recall
- 4. Use regular expressions to see whether stemming or wildcarding is required