

# ConceptVector: Text Visual Analytics via Interactive Lexicon Building using Word Embedding

Deokgun Park<sup>\*</sup>, Seungyeon Kim<sup>\$</sup>, Jurim Lee<sup>&</sup>, Jaegul Choo<sup>&</sup> Nicholas Diakopoulos<sup>\*</sup>, and Niklas Elmqvist<sup>\*</sup>

\*University of Maryland USA {intuinno, nad, elm}@umd.edu <sup>\$</sup>Google USA sylund@gmail.com <sup>&</sup>Korea University South Korea {jurim0301, jchoo}@korea.ac.kr

# 











# Visualization





# Visualization







# NUMBERS











# **NUMBERS**







Lexicon = A set of keywords that is related to specific concept

Positive = { good, great, happy, ... }

Negative = {bad, worst, horrible, ... }



# NUMBERS

/, ... } e, ... }







# Positive = { good, great, happy, ... } Negative = {bad, worst, horrible, ... }



# **NUMBERS**







Positive = { good, great, happy, ... }

Negative = {bad, worst, horrible, ... }



# NUMBERS







Positive = { good, great, happy, ... }

Negative = {bad, worst, horrible, ... }



"The movie was great. I was happy."

"The movie opens today."

"The movie is *horrible*."



# **NUMBERS**







Positive = { good, great, happy, ... }
Negative = {bad, worst, horrible, ... }





# NUMBERS

ositive	Negative		
2	0		
0	0		
0	1		



# Linguistic Inquiry and Word Count (LIWC)

Г	nictionary / Thesaurus	
	Aaa Lorem Ipsum Dolor Sit A Etiam	met
C		

		83		84	85	
Power			Reward	Risk		
a-list*	forbids	obeyed	under	access*	abstain*	accepted
above	force*	obeying	underclass*	accrue*	alarm*	added
acclaimed	frail*	obeys	underdog*	accumul*	apprehens*	admitted
administr*	freshm*	officehold*	underling*	achievable	apprenhens*	affected
age	glori*	officer*	underprivileg*	achieve*	averse	ago
allow*	glory	official*	unimportant	achievi*	aversi*	already
amateur*	god	oppose*	unqualified	acquir*	avert*	appeared
ambition	god's	oppositi*	unwanted	add	avoid*	approached
ambitions	goddess*	order	up	added	bad	arrived
ambitious	good-for-nothing	orders	upper	adding	balk	asked
apolog*	govern*	outcast	upperclass*	adds	beware	ate
apprentic*	greatest	outrank*	useless	advanc*	careful*	attended
approv*	greatness	over	vanquish*	advantag*	caution*	attracted
armies	greed*	overpower*	verb	adventur*	cautious*	became
army	grown-up*	overrul*	veteran*	amass*	cease*	been
asham*	grownup*	owner*	victim*	approach	concern	began
assault*	help	pariah*	victor*	approached	consequen*	begged
assertive	helper*	parliament*	vip	approaches	crises	believed
attack*	helpless*	passiv*	vp*	approaching	crisis	born
attendant	hierarch*	pathetic	vulnerab*	award*	curb*	bought
authorit*	high	pathetically	war	benefit	danger	bounced
authorize	high-ranking	patriarch*	warfare*	benefits	dangerous	braved
battl*	higher	peasant*	warred	best	dangerously	broke
beat	highest	peon*	warring	bet	dangers	brought

90		
FocusPast		
mapped	we've	
mastered	weakened	
mated	weighed	
meant	weirded	
messaged	went	
met	wept	
might've	were	
mightve	weren't	
missed	werent	
mocked	weve	
mothered	what'd	
moved	whatd	
must've	where'd	
mustve	wished	
named	wobbled	
narrowed	woke	
neared	woken	
needed	won	
noticed	wondered	
obeyed	wore	
obtained	worked	
od'ed	worn	
okayed	worsen	
organized	would've	



# How to Build a Lexicon?



# How to build Lexicon

\_\_\_\_\_



### **Hand-picking**

- Linguistic Inquiry • and Word Count (LIWC)
- **General Inquirer** •

• High Quality, strong signal words

- Hard to build, scale
- Does not adapt to different domain (e.g. Twitter)

Manual

(GI)

### **Crowdsourcing**

Affective Norms for ٠

English Words

(ANEW)

Hedonometer ٠

 Scales up for single category Costly, limited category

### Manual







# How to build Lexicon

- Domain adaptable
- Scales to diverse topic
- Difficult to label each group



### **Topic Modeling**

#### Latent Semantic

٠

•

### Indexing (LSA)

#### Latent Dirichlet

#### Allocation (LDA)

# How to build Lexicon



- Scales to diverse topic
- Easy to interpret



Empath •

Manual



Pennington, Jeffrey, Richard Socher, and Christopher Manning. "Glove: Global vectors for word representation." Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). 2014.

Goldberg, Yoav, and Omer Levy. "word2vec Explained: deriving Mikolov et al.'s negative-sampling word-embedding method." arXiv preprint arXiv:1402.3722 (2014).

# Nearest Neighbor





- 1. Frogs
- 2. Toad
- 3. Litoria
- 4. Leptodacitylidae
- 5. Rana
- 6. Lizard
- 7. eleutherodactylus







## Nearest Neighbor



Frog

- 1. Frogs
- 2. Toad
- 3. Litoria
- 4. Leptodacitylidae
- 5. Rana
- 6. Lizard
- 7. eleutherodactylus









Image from http://cs.stanford.edu/people/karpathy/tsnejs/wordvecs.html

## Previous Work

ŵ

# Empath

Empath can generate new lexical categories and analyze text over 200 built-in human-validated categories. This is a web demo; you can download our python tool to more easily run larger and or take a look at our recent CHI paper.

#### Generate category:

comma-separted terms, e.g., 'spoon, oven, counter

Generate

#### Analyze text over default categories

Fast, E., Chen, B., & Bernstein, M. S. (2016, May). Empath: Understanding topic signals in large-scale text. In Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems (pp. 4647-4657). ACM.







Compute Word Embedding

















Refinement of seed words Based on recommended words





Refinement of seed words Based on recommended words





Refinement of seed words Based on recommended words





#### Refinement of seed words Based on document analysis



Refinement of seed words Based on recommended words





X 3,000

Tweets from Trump and Hillary during 2016 election





{ pollen, cherry, underbrush, thicket, grape, sunflower, field, willow, rose, fenced, bouquet, flowered ... bush ... }



{ pollen, cherry, underbrush, thicket, grape, sunflower, field, willow, rose, fenced, bouquet, flowered ... bush ... }

Trump talked 3.31 times more about Plant than Hillary




## Trump talked 3.31 times more about Plant than Hill

## False positive errors



# Trump talked 3.31 times more about Plant than False positive errors

Polysemy bush -n. a low plant with many branches that arise from or near the ground. Bush – n. Jeff Bush



## Top 10 categories

## Trump vs Hillary











1st Word 2nd Word 3rd Word 4th Word etc

## Design Requirements



**D1: Diverse needs** 

**D2: Integrated loop** 



## **D3: In-context**

### Positive Words Input Please type a few words for the concepts you are looking for.



(d) Irrelevant words input

### Description: This is a concept about crude oil

This is a concept about crude oil

### Concept Type: Unipolar vs Bipolar

In unipolar type, the dictionary will be created around single concepts. For example, crude oil or car can be an unipolar type. In bipolar type, the dictionary will be created between two concepts. For example, happiness vs sadness or Democratic vs Republican can be an bipolar type.



### **Positive Words Input**

Please type a few words for the concepts you are looking for.



Click to add suggested words.



No Data Available.

Advanced Settings



**D1: Diverse needs** 

## Sophisticated concept modeling

**Irrelevant words** 

I am interest in "tidal flooding", not "storm flooding."



## Sophisticated concept modeling

**Irrelevant words** 

I am interest in "tidal flooding", not "storm flooding."

**Bipolar Conceps** 

Can I map the words continuously from "Democratic party" to "Republican party"?





## Sophisticated concept modeling

I am interest in "tidal flooding", not "storm flooding."

**Irrelevant words** 

**Bipolar Conceps** 

Can I map the words continuously from "Democratic party" to "Republican party"?



Kernel Density Estimation (KDE)

using Gaussian Kernel

where bandwidth represents selectivity of seed words



Cosine similarity for relevance score

## Previous approach



Cosine similarity for relevance score





Cosine similarity for relevance score





### Negative Seed







### Concept Name: Democrat vs Republic

Democrat vs Republic

### Description: This shows the difference between two parties

This shows the difference between two parties

### Concept Type: Unipolar vs Bipolar

In unipolar type, the dictionary will be created around single concepts. For example, crude oil or car can be an unipolar type. In bipolar type, the dictionary will be created between two concepts. For example, happiness vs sadness or Democratic vs Republican can be an bipolar type.

Unipolar Bipolar

from the Conceptvector Team

## Bipolar concepts



Advanced Settings 3



Refinement of seed words Based on recommended words



Donald J. Trump speaking in Milford, N.H., a week before the state's primary this month. Some establishment Republicans have been scrambling for a way to prevent him from becoming the party's presidential nominee. Damon Winter/The New York Times

### 2763 COMMENTS

Readers shared their thoughts on this article.

The comments section is closed. To submit a letter to the editor for publication, write to letters@nytimes.com.

All 2763	Readers' Picks 1721	NYT Picks 21

### David Weaver Orlando - February 27, 2016

The political industry does not realize that it is in the middle of a hostile takeover from people who are beyond sick and tired of insiders who use the aspirations of hard working people to line their pockets.

Is Trump perfect? Nope. But the leaders of this nations political industry have failed us miserably. Hopefully the "we know what's best for you" crowd has its' days numbered.

1100 Recommend

	Gfagan	PA -	February	27.	2016
27	Dr Frank	ensti	en frantic	ally	nailed

Dr. Frankenstien frantically nailed the planks into place, trying to seal his laboratory off from the rest of the world.

What would happen if the Monster got out? What if people realized he had created it?

The pounding from the inside increased in intensity. The Monster roared. It sounded like "Yuge!"

Dr. Frankenstein wondered what that might mean, then resumed hammering on the nails.

But the door gave way and the boards he was trying to nail in place began to bend under the Monster's fists.

Rage. That is all the doctor saw in the contorted, strangely orange face of his creation.

"Yuge!"

Splinters began flying from the boards.

Dr. Frankenstein turned his back and ran. He'd go to the House. There, at least, he'd be safe doing nothing.

When the questions came, the doctor would deny everything. He'd blame his neighbor, Mr. Obama.

Yes, that would work. That would deflect attention away from his own nefarious activities, spread over the preceding decades, that had led to this catastrophe.

It was now up to others to deal with his Creation. the substance of the same locale





## **Document Analysis**



### Home / Articles / CommentIQ

### Inside the Republican Party's Desperate Mission to Stop Donald Trump 🕑

Despite all the forces arrayed against Mr. Trump, a paralytic sense of indecision and despair has prevailed."

2016-02-28T00:00:00+00:00

Presidential Election of 2016; Trump, Donald J; Republican Party; Christie, Christopher J;Kasich, John R;Rubio, Marco;United States Politics and Government

By ALEXANDER BURNS, MAGGIE HABERMAN and JONATHAN MARTIN

U.S., Article





Rank comments by

This bar shows the weights for this ranking.

New (276	3)	Acc	epted (0)	Reje
Kevin O'Brier	1			
see no hop	be for a	GO	P evolution	in 8 ye
Accept	Rejec	ct	Pick	

n/a CParis

Search for ...

worse than a bunch of first graders.

Reject Accept

Pick

### PE

Cruz and Rubio tack so Tea Party right, so irrationally "conservative" they represent business as usual Republicans. Trump brings moderates in with hishedge on health care and vagueness on Planned Parenthood. Rove is wrong to think that a candidate



Seattle, WA , Feb 27, 2016 1:13:34 PM



## Remove Irrelevant words

conceptvector ora

### score = 0.0010129082575745594 Reject Pick Accept goeasyonus Trump exudes leadership, we havent had a leader as president since Reagan. The country is sick and tired of vanilla wrapper puppets in the WH. Trump is a real person, speaks his mind , has the country first attitude ..... and owes no one ...... score = 0.00034310946674149856 Reject Pick Accept goeasyonus being a democrat, arent you more knowledgeable of the dem party than the repub party ?? why dont u tell us all the great things the dem party is doing for the country. Why do those with the loudest complaints always ' belong ' to a party ? Too bad more voters dont think for themselves . score = 0.00030419896280519385 Accept Reject Pick

Ernesto

The first high profile defection to the Republicans was Strom Thurmond . Of course , by your estimation that had nothing to do with racial politics, right?

score = 0.00028805119766432566

Accept Reject Pick

2

great nw, Feb 29, 2016 2:29:17 AM

great nw, Feb 29, 2016 2:37:44 AM

NYC, Feb 29, 2016 1:12:22 AM

### CommentPlot remporal

and don't want to do anything about illegal immigration . Immigration effects everyth



Accep	ot	Re	ect	Pick
-------	----	----	-----	------

### Melvin

This could have been avoided by enforcing the immigration laws . It's a bipartisan failure .

score = 0.00565333256500768

Accept Reject Pick

### Bahtat

labor, they are not allowed to VOTE.

score = 0.005202291401691696

Accept Reject Pick

### Ali

Too bad they simply weren't willing to give up amnesty and more immigration, legal and illegal. They complain about gov't subsidies to business yet that's just what illeg aliens and immigrants are, subsidies. Immigration is the focus of much of what's wrong in this country, to those supporting Trump. from unemployment to under performing schools to high healthcare costs and water shortages .

### score = 0.004590839907424656

Reject Pick Accept

denver, co , Feb 27, 2016 3:15:43 PM jmf No its not . Who has stopped all ways to get rid of illegals ? Who has allowed massiv amounts of legal immigration ? Its both of them so that average people cant exist wh



You are viewing this concept as a guest. Please login and copy this concept to create a clone and edit.

## Temporal Trend

Immigration-related

SF, Feb 28, 2016 3:29:46 PM

San Diego , Feb 27, 2016 4:02:37 PM The GOP elites have another great love for Illegal Immigrants . As well as being chea

Michigan , Feb 27, 2016 7:10:03 PM



## Immigration



### Comment Proteration Test lype: unipolar

Positive inputs

αιιχεία κυτετιτ

want to end this rot in the system ? amend the constitution ! 1) Limit members of Congress to a 4yr single term . that gets them out of the business of endlessly seeking fat cat sponsors , and brings them to the business of true governance . the same applies to senators . let them all step aside after a term to give others also the chance to "serve the American people . " 2) Get the Citizens United judgement reversed . No more big money super-packs . 3) Apportion equal time to all candidates on the media . the media has overplayed the "elections are a spectator/gladiatorial sport" angle, by focusing on their ratings and the 'fun, and the bloodier' candidates . 4) Election day should be declared a holiday, so all can afford to go to vote. Even a far poorer country like India , has this . 5) make voting compulsory . We want a democratic country ? then let's get off our couches .

score = 0.0006589178749339568

Pick Accept Reject

great nw , Feb 29, 2016 2:29:17 AM goeasyonus Trump exudes leadership, we havent had a leader as president since Reagan. The country is sick and tired of vanilla wrapper puppets in the WH. Trump is a real person, speaks his mind , has the country first attitude ..... and owes no one .....

score = 0.00014680238075992264

Accept

Reject Pick

great nw , Feb 29, 2016 2:37:44 AM goeasyonus being a democrat, arent you more knowledgeable of the dem party than the repub party ?? why dont u tell us all the great things the dem party is doing for the country.

Why do those with the loudest complaints always ' belong ' to a party ? Too bad more voters dont think for themselves .

score = 0.00011201416478632322

Reject Accept

Pick





th	re	at	12.1	wrong problem
V	ali	d	at	e: observe and interview target users
F	t	hr	е	at: bad data/operation abstraction
Г	L		t	hreat: ineffective encoding/interaction t
L	L		۷	alidate: justify encoding/interaction des
L	L			threat: slow algorithm
	L			validate: analyze computational com
	L			implement system
	L			validate: measure system time/memo
	L	ŀ	v	alidate: qualitative/quantitative result in
	L	I	[	test on any users, informal usability stu
	L		V	alidate: lab study, measure human tim
		Va	ali	date: test on target users, collect anec
		Va	ali	date: field study, document human usa
V	ali	d	at	e: observe adoption rates
-				

Munzner, Tamara. "A nested model for visualization design and validation." IEEE transactions on visualization and computer graphics 15.6 (2009).



### Lab Experiment

- Lexicon-building interface
- Wordnet and

Thesarus.com

### Lab Experiment

- Lexicon-building interface
- Wordnet and

Thesarus.com

### <u>Quantitative</u> <u>Evaluation</u>

- Bipolar concept modeling
- Compared with
  crowdsourced data

### Lab Experiment

- Lexicon-building • interface
- Wordnet and •

Thesaurus.com

### **Quantitative Evaluation**

- Bipolar concept ٠ modeling
- Compared with • crowdsourced data



- - and usage

•

- expert and NLP
- expert

### **Expert Feedback**

- System limitation
- Visual analytics



## Conclusion



## Take-away Message

 You can build custom dictionary for your own domain and analyze text.

 Interactive refinement may improve quality of text analysis by reducing false positive errors.
Seungyeon Kim Asylund@gmail.com

Nicholas Diakopoulos @ndiakopoulos



Jaegul Choo jchoo@korea.ac.kr



Niklas Elmqvist @NElmqvist

1.17

# Meet the Team



#### Deokgun Park @intuinno



#### Jurim Lee jurim0301@korea. ac.kr









## Questions?

### Try demo at: Conceptvector.org intuinno@umd.edu

Positive Words Input Please type a few words for the concepts you are looking for.





(d) Irrelevant words input