

# Topic analysis using Mallet and network graphs

Rob McDaniel  
Founder  
[Linguistic.com](http://Linguistic.com)

# About me

[robmcdan@gmail.com](mailto:robmcdan@gmail.com)

<http://lingistic.com>

[github.com/robmcdan](https://github.com/robmcdan)

<https://github.com/Lingistic/>

<https://www.linkedin.com/in/robmcdan>

previous projects:

- Rakuten
- Payscale.com
- Microsoft

# About Linguistic

---

- Developed bias detection model for detection of political bias in webpages (currently in beta)
- using topics to categorize news articles and editorials
- entity and semantic extraction from articles
- topics (among other things) helps us disambiguate vocabulary

# Overview

---

what I'm going to cover:

Semantics vs. syntax

topic models and how they work

how to use mallet

measuring topic interaction

high-level, code available for off-line analysis

# Part 1: basics

---

What are topic models and what data does it produce

What is the vector space model

Training vs. inferring

# what are semantics

---

semantics: deals with meaning

the relationships of words together form the semantics

Syntax-based NLP tends to miss meaning; i.e. “Islamic Terrorism” and “Islamic Extremism” are syntactically dissimilar but semantically related.

Statistically prevalent syntax can improve semantic topic models (more later)

# what is a topic model

---

- unsupervised; discovers themes in unstructured text
- bag of words model
- generative model
- can be thought of as a clustering algorithm
- topics: distributions over words
- document: distribution of topics



## Topics

## Documents

## Topic proportions and assignments

gene 0.04  
dna 0.02  
genetic 0.01  
...

life 0.02  
evolve 0.01  
organism 0.01  
...

brain 0.04  
neuron 0.02  
nerve 0.01  
...

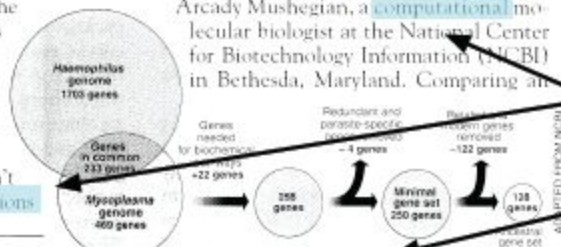
data 0.02  
number 0.02  
computer 0.01  
...

### Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson at Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic** **numbers** game, particularly as more and more **genomes** are completely mapped and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

**Stripping down.** **Computer analysis** yields an estimate of the minimum modern and ancient genomes.



music  
band  
songs  
rock  
album  
jazz  
pop  
song  
singer  
night

book  
life  
novel  
story  
books  
man  
stories  
love  
children  
family

art  
museum  
show  
exhibition  
artist  
artists  
paintings  
painting  
century  
works

game  
knicks  
nets  
points  
team  
season  
play  
games  
night  
coach

show  
film  
television  
movie  
series  
says  
life  
man  
character  
know

theater  
play  
production  
show  
stage  
street  
broadway  
director  
musical  
directed

clinton  
bush  
campaign  
gore  
political  
republican  
dole  
presidential  
senator  
house

stock  
market  
percent  
fund  
investors  
funds  
companies  
stocks  
investment  
trading

restaurant  
sauce  
menu  
food  
dishes  
street  
dining  
dinner  
chicken  
served

budget  
tax  
governor  
county  
mayor  
billion  
taxes  
plan  
legislature  
fiscal

# history of LDA

---

originally used for finding patterns in genetic data

highly useful in today's world of big data

many implementations available

# Steps of topic modelling

---

vectorize training documents

train

vectorize unseen documents

infer topics

# vectorizing

---

- every document is represented as a numerical vector
- large vocabulary = sparse matrix
- multi-dimensional vector space model
- Outputs word sequences for documents

## Documents



## Vector-space representation

However, complexity

We will see how small

Given a function based

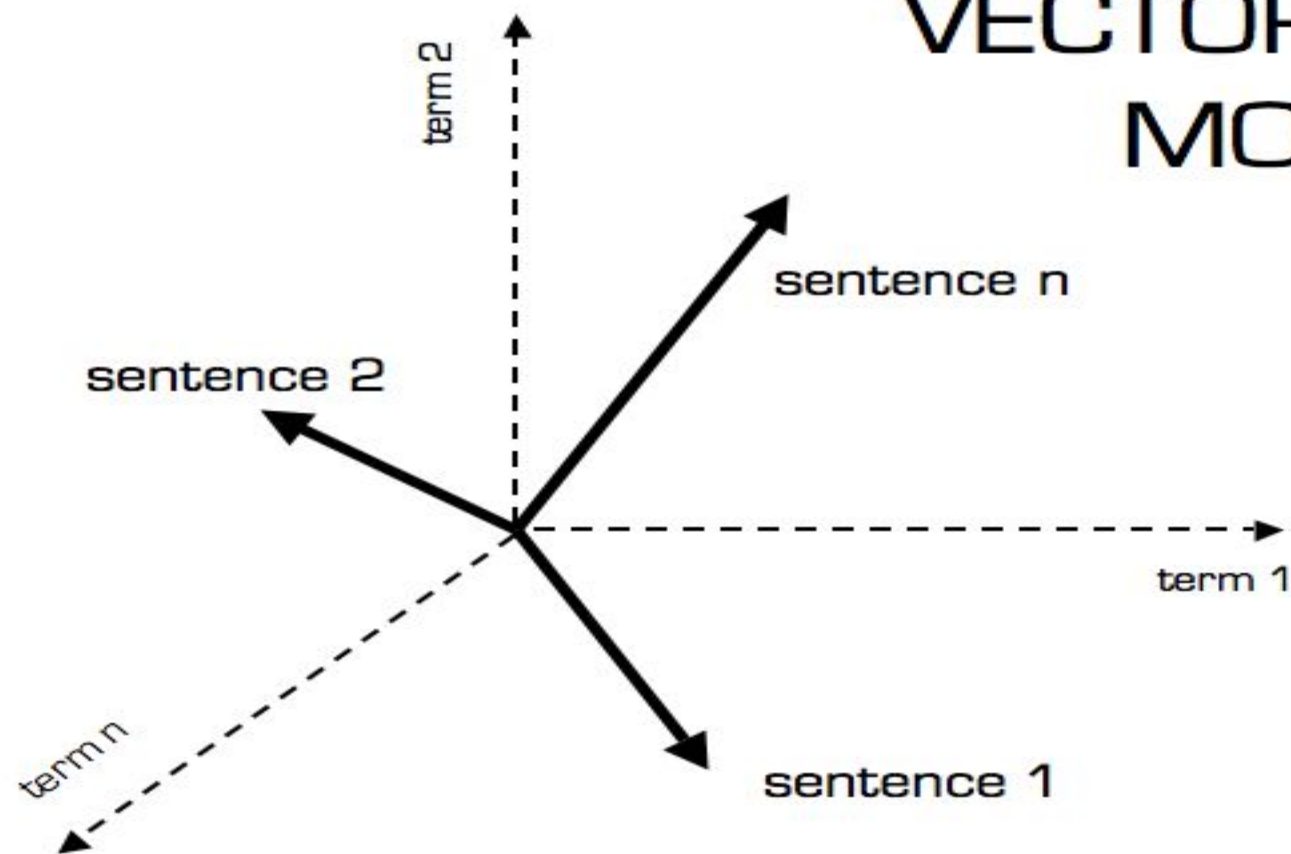
Using entropy of traffic

We study the complexity of influencing elections through bribery: How computationally complex is it for an external actor to determine whether by a certain amount of bribing voters a specified candidate can be made the election's winner? We study this problem for election systems as varied as scoring ...

|            | D1 | D2 | D3 | D4 | D5 |
|------------|----|----|----|----|----|
| complexity | 2  |    | 3  | 2  | 3  |
| algorithm  | 3  |    |    | 4  | 4  |
| entropy    | 1  |    |    | 2  |    |
| traffic    |    | 2  | 3  |    |    |
| network    |    | 1  | 4  |    |    |

Term-document matrix

# VECTOR SPACE MODEL



# training

---

- training on word sequences
- non-deterministic; set random seed for consistent re-modelling
- must predetermine number of topics

# inference

---

- infer topics from unseen documents
- deterministic
- must use original vocabulary



# Part 2: Preparing the data

---

- Data cleaning is everything
- Lots of tricks; I'll keep it simple

# The corpus

---

GOP and Democratic debates for 2016 election cycle

<http://www.presidency.ucsb.edu/debates.php>

# Why presidential debates?

---

- nicely chunked into little context chunklets
- wide variety of speakers and context
- limited vocabulary

# Prepare the corpus

---

- garbage in, garbage out
- stop word removal
- isolating key phrases
- parsing relevant items

# isolating key phrases

---

what are ngrams

finding likely ngrams

# nltk

---

- NLTK (natural language tool kit)
- built in collocation measures
- likelihood measure finds ngrams which go together often, based on prior occurrence

# ngram likelihood

---

```
##tokenize sentence into words
```

```
bigram_measures = nltk.collocations.BigramAssocMeasures()
```

```
bigram_finder = BigramCollocationFinder.from_words(words)
```

```
bigram_finder.score_ngrams(bigram_measures.likelihood_ratio)
```

# Example ngrams by likelihood measure

Secretary\_Clinton → 1928.6349782078692  
United\_States → 1582.65386804903  
Senator\_Sanders → 1430.490764995243  
Senator\_Rubio → 1109.1112295674834  
Wall\_Street → 1050.378694855774  
Senator\_Cruz → 1000.1855691365586  
New\_Hampshire → 908.0197777658559  
President\_Obama → 788.9489874026619  
Governor\_Christie → 732.7257542072496  
Governor\_Bush → 728.8350811850478  
North\_Korea → 597.3617719755027  
Governor\_Kasich → 589.1031551108515  
commercial\_break → 556.4690710489392  
Senator\_Paul → 547.8797358958528  
Hillary\_Clinton → 486.48080519364464  
health\_care → 460.21407947612664  
Donald\_Trump → 447.3320619557984  
bell\_rings → 439.30608437343244  
climate\_change → 436.45248108215134  
Barack\_Obama → 409.1356822232396  
foreign\_policy → 399.3244418178018  
White\_House → 386.00309121441626  
Des\_Moines → 380.3797449598837  
Dana\_Bash → 349.476587224584  
Ronald\_Reagan → 345.09198162490304  
Middle\_East → 325.27298204530655



# why not just nltk?

---

- ngrams by collocation are neat, but don't capture semantics
- when words are slightly different, they appear unrelated
- still useful for seeding LDA

# replacements and deletions

---

mallet allows for replacement and deletion of words

example:

New Hampshire -> New\_Hampshire

I'm -> <deleted>

# why replace?

---

- allows the inclusion of ngrams intoallet
- topic: “elections, new\_hampshire” instead of “elections, new, hampshire”
- remove words we don't care about

# process

---

Process ngrams (see code sample on github)

build replacement file (see code sample on github)

parse debates (see code sample on github)

vectorize and train lda

predict topics

build graph

# Part 3: usingallet

---

# generating topic models using `mallet`

---

installing `mallet` (fork available on my github account which ignores case sensitivity)

Java

`./bin/mallet` is a wrapper for accessing `mallet` features

# import data step

---

can import either a single file, one example per line (mallet import-file)

or can import a directory of files (mallet import-dir)

for an example, see Data/Debates/mallet\_files/input\_command.txt in code samples

# train model step

---

`./bin/mallet train-topics`

pass sequences file

specify outputs

for an example, see `Data/Debates/mallet_files/train_command.txt`



# interpreting the Mallet output files

---

doc\_topics -- the proportion of topics (columns) in each document (rows)

topic\_keys -- N words for each topic, to “describe” it

topic\_counts -- a count of each topic word and how many times it occurs in each topic

# import unseen document

---

import-file as before, but must use --use-pipe-from flag

vectorizes according to existing model vocab

# infer topics

---

infer-topics -- specify input doc, and inferencer file

# Putting topics together: interactions

---

measuring interactions using KL Divergence

- measures the differences in  $P(W | T)$  across documents
- captures how often topics occur with other topics
- topics that occur with others must be related
- threshold is important

# Interactions as networks

---

Topic -> Node

Divergence -> weight

generates undirected network

networkX python package will output to graphml format

america, people, jobs

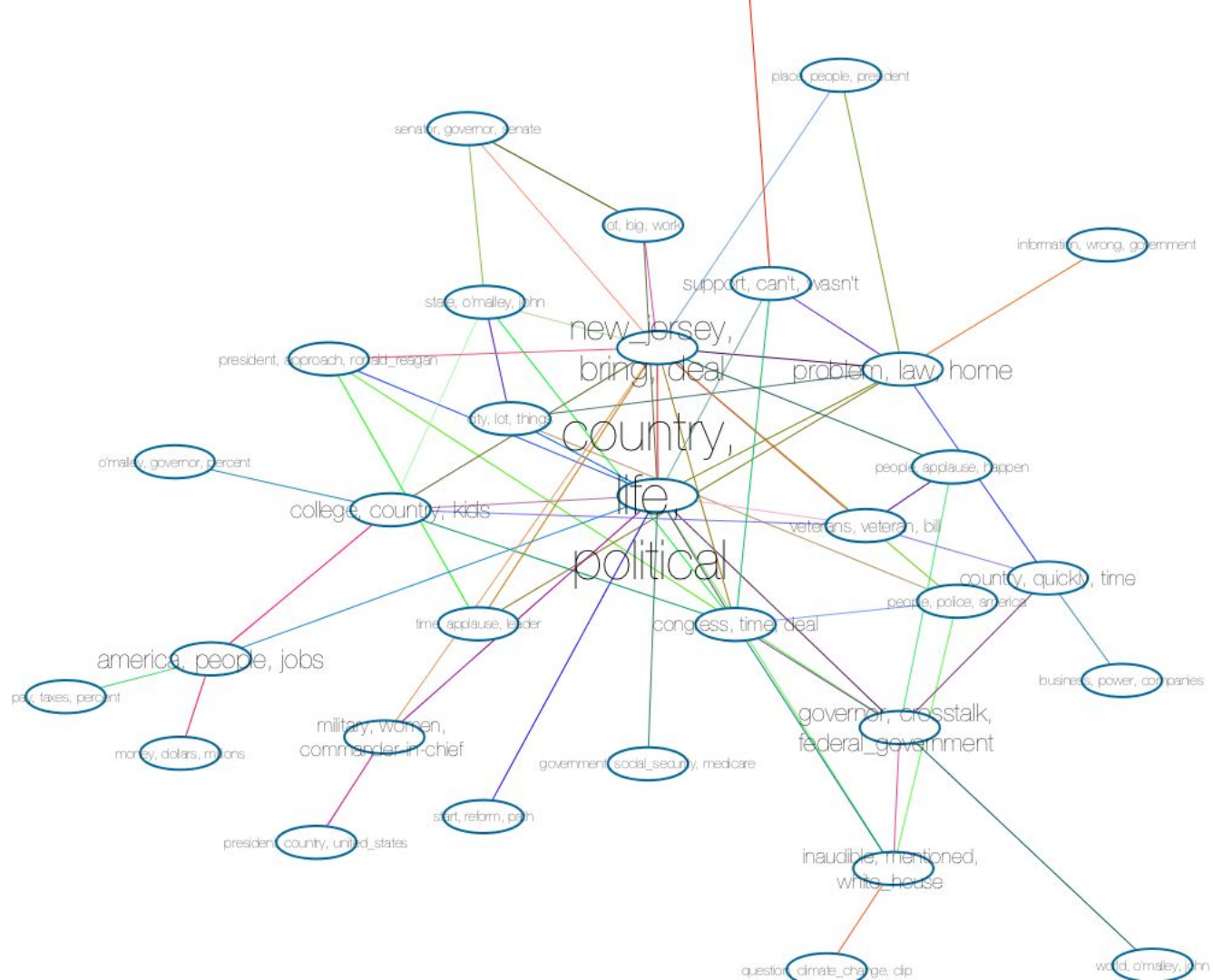
pay, taxes, percent

money, dollars, millions

# Graphing in Cytoscape

---

- cytoscape -- open source
- popular in bioinformatics
- complex networks
- <http://www.cytoscape.org/>
- <http://diging.github.io/tethne/api/tutorial.mallet.html>
-





# Possible improvements

---

remove noise / spelling correction

Train Sub-topic models

better data sampling -- some candidates speak more than others, which produces an imbalanced dataset

Model Topics by candidate and perform sentiment analysis/objectivity by speaker and topic

# Resources

---

<https://github.com/robmcdan/Mallet>

<https://networkx.github.io/>

<http://www.cytoscape.org/>

example code:

<https://github.com/Lingistic/DebateAnalysis>

# references & resources

---

<http://www.linguistic.com/blog/2016/2/1/5gtebcogi0xv2hiukgd1lvw9kg4322>

<http://diging.github.io/tethne/api/tutorial.mallet.html>

<https://www.cs.princeton.edu/~blei/papers/Blei2012.pdf>

<http://mimno.infosci.cornell.edu/papers/mimno-semantic-emnlp.pdf>

<http://mallet.cs.umass.edu/about.php>

<http://yosinski.com/mlss12/MLSS-2012-Blei-Probabilistic-Topic-Models/>