# LING 410X: Language as Data
## Semester: Spring '18

Instructor: Sowmya Vajjala

Iowa State University, USA

20 March 2018

# Class Outline

- Assignment 4 - brief discussion
- Introduction to topic modeling
- Assignment 5 description
- This week and next week: discussion about topic modeling, and how to make our own models

# Announcement

- Thursday: No class, as I am at a conference and our talk is scheduled at the same time. But, there is an optional exercise, based on Chapter 13 in the textbook.

# Rest of the semester

Macro Analysis of Texts:

1. Classification (Chapter 12) - This is what we discussed before the break
2. Topic modeling (Chapter 13) - this week and next
3. Clustering (Chapter 11) - When talking about visualization, after topic modeling

# Assignment 4 Grading

- First question: frequent words in the training and testing csv files - straight forward steps (everyone did this part well)

# Assignment 4 Grading

- First question: frequent words in the training and testing csv files - straight forward steps (everyone did this part well)
- Second question:
  - More like a exploratory assignment. The tutorial works, I wanted you to try changing somethings in that procedure and see what happens.
  - My take on your submissions: I liked the ones which documented their experiences -e.g., what worked, what failed? did they manage to get better results? etc.
  - The errors some of you saw when you tried to change the SVM algorithm to something else: Are they only in Windows? (some of you did not see this). One other possibility: memory limits.
  - I did not deduct any points for these errors - they are a part of the assignment!
- Where to go from there: Look for some publicly available classification dataset (there are many, I gave some links last time) and try to do classification.

Introduction to Topic Modeling

# What is topic modeling?

- Topic Models are a group of algorithms which attempt to discover latent themes in large collections of documents.
- They use statistical methods to analyze word usage in the texts to discover what "themes" run through them, how these themes connect to each other etc.

# What is topic modeling?

- Topic Models are a group of algorithms which attempt to discover latent themes in large collections of documents.
- They use statistical methods to analyze word usage in the texts to discover what "themes" run through them, how these themes connect to each other etc.
- Good thing about them: they do not expect us to provide any prior annotations/categories for texts. Topics will "emerge" from the analysis.
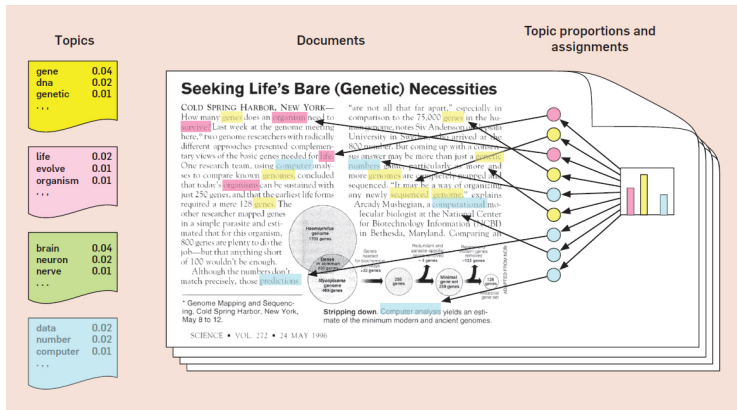- Bad thing: Lot of math behind it (but we do not have to understand the math to apply topic models)

# What is topic modeling?

- ▶ Topic Models are a group of algorithms which attempt to discover latent themes in large collections of documents.
- ▶ They use statistical methods to analyze word usage in the texts to discover what "themes" run through them, how these themes connect to each other etc.
- ▶ Good thing about them: they do not expect us to provide any prior annotations/categories for texts. Topics will "emerge" from the analysis.
- ▶ Bad thing: Lot of math behind it (but we do not have to understand the math to apply topic models)
- ▶ One of the most popular methods of analyzing unstructured text data.

# Latent Dirichlet Allocation (LDA)

- LDA is the simplest topic modeling algorithm
- Intuitions:
  - each document is a mixture of multiple topics
  - each topic can be characterized by some set of keywords related to that topic.
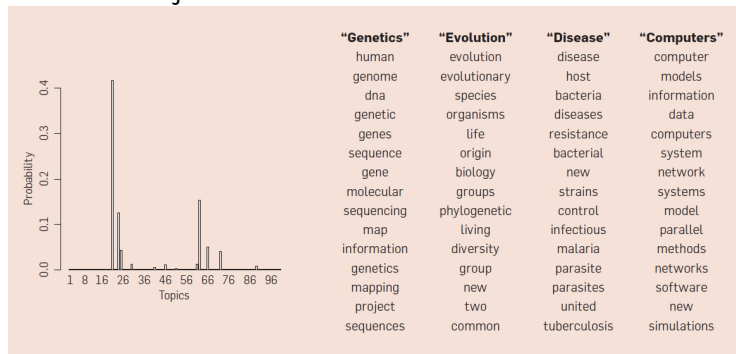  - a keyword can exist in multiple topics with different degrees of importance.

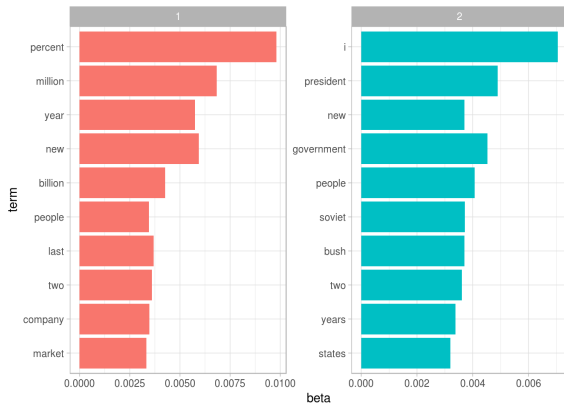# What does a Topic Model do?-1



source: `https://goo.gl/azc7Gc`

# What does a Topic Model do? -2

Real inference with LDA - topic model built using 17000 articles from Science journal.



source: https://goo.gl/azc7Gc

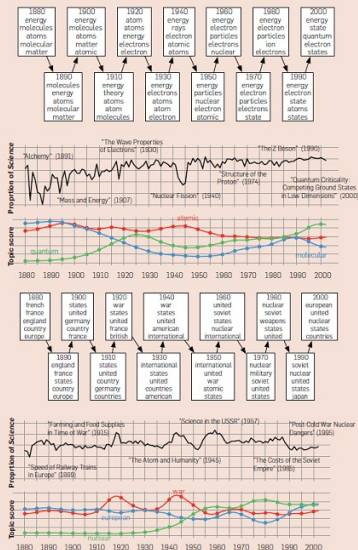# What does a Topic Model do? -3



source: `http://tidytextmining.com/topicmodeling.html`

# How are topic models useful? -1

- One application comes from the textbook author.
- He and his team analyzed best selling novels using topic models and concluded that best seller novels focus on a small number of topics instead of discussing 1000 things in one story :)
- Okay, it is more detailed than that, I am just telling a micro-summary of what they concluded.

# How are topic models useful? -2

# How are topic models useful? -3

Analyzing topics by author

| TOPIC 10 | |
|---|---|
| WORD | PROB. |
| SPEECH | 0.1134 |
| RECOGNITION | 0.0349 |
| WORD | 0.0295 |
| SPEAKER | 0.0227 |
| ACOUSTIC | 0.0205 |
| RATE | 0.0134 |
| SPOKEN | 0.0132 |
| SOUND | 0.0127 |
| TRAINING | 0.0104 |
| MUSIC | 0.0102 |
| **AUTHOR** | **PROB.** |
| Waibel_A | 0.0156 |
| Gauvain_J | 0.0133 |
| Lamel_L | 0.0128 |
| Woodland_P | 0.0124 |
| Ney_H | 0.0080 |
| Hansen_J | 0.0078 |
| Renals_S | 0.0072 |
| Noth_E | 0.0071 |
| Boves_L | 0.0070 |
| Young_S | 0.0069 |

| TOPIC 209 | |
|---|---|
| WORD | PROB. |
| PROBABILISTIC | 0.0778 |
| BAYESIAN | 0.0671 |
| PROBABILITY | 0.0532 |
| CARLO | 0.0309 |
| MONTE | 0.0308 |
| DISTRIBUTION | 0.0257 |
| INFERENCE | 0.0253 |
| PROBABILITIES | 0.0253 |
| CONDITIONAL | 0.0229 |
| PRIOR | 0.0219 |
| **AUTHOR** | **PROB.** |
| Friedman_N | 0.0094 |
| Heckerman_D | 0.0067 |
| Ghahramani_Z | 0.0062 |
| Koller_D | 0.0062 |
| Jordan_M | 0.0059 |
| Neal_R | 0.0055 |
| Raftery_A | 0.0054 |
| Lukasiewicz_T | 0.0053 |
| Halpern_J | 0.0052 |
| Muller_P | 0.0048 |

| TOPIC 87 | |
|---|---|
| WORD | PROB. |
| USER | 0.2541 |
| INTERFACE | 0.1080 |
| USERS | 0.0788 |
| INTERFACES | 0.0433 |
| GRAPHICAL | 0.0392 |
| INTERACTIVE | 0.0354 |
| INTERACTION | 0.0261 |
| VISUAL | 0.0203 |
| DISPLAY | 0.0128 |
| MANIPULATION | 0.0099 |
| **AUTHOR** | **PROB.** |
| Shneiderman_B | 0.0060 |
| Rauterberg_M | 0.0031 |
| Lavana_H | 0.0024 |
| Pentland_A | 0.0021 |
| Myers_B | 0.0021 |
| Minas_M | 0.0021 |
| Burnett_M | 0.0021 |
| Winiwarter_W | 0.0020 |
| Chang_S | 0.0019 |
| Korvemaker_B | 0.0019 |

| TOPIC 20 | |
|---|---|
| WORD | PROB. |
| STARS | 0.0164 |
| OBSERVATIONS | 0.0150 |
| SOLAR | 0.0150 |
| MAGNETIC | 0.0145 |
| RAY | 0.0144 |
| EMISSION | 0.0134 |
| GALAXIES | 0.0124 |
| OBSERVED | 0.0108 |
| SUBJECT | 0.0101 |
| STAR | 0.0087 |
| **AUTHOR** | **PROB.** |
| Linsky_J | 0.0143 |
| Falcke_H | 0.0131 |
| Mursula_K | 0.0089 |
| Butler_R | 0.0083 |
| Bjorkman_K | 0.0078 |
| Knapp_G | 0.0067 |
| Kundu_M | 0.0063 |
| Christensen-J | 0.0059 |
| Cranmer_S | 0.0055 |
| Nagar_N | 0.0050 |

Figure 3: An illustration of 4 topics from a 300-topic solution for the CiteSeer collection. Each topic is shown with the 10 words and authors that have the highest probability conditioned on that topic.
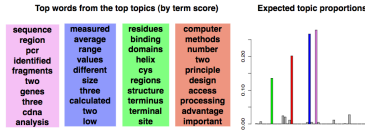
source:
https://mimno.infosci.cornell.edu/info6150/readings/398.pdf

# How are topic models useful? -4
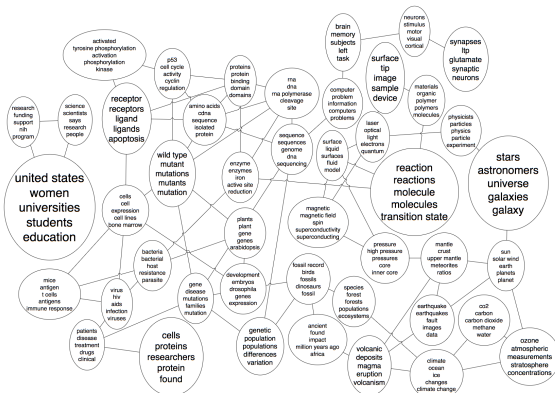## Picking up similar documents



FIGURE 4. The analysis of a document from *Science*. Document similarity was computed using Eq. (4); topic words were computed using Eq. (3).

source:
http://www.cs.columbia.edu/~blei/papers/BleiLafferty2009.pdf

# How are topic models useful? -5

Topic Graphs



source:
http://www.cs.columbia.edu/~blei/papers/BleiLafferty2009.pdf

# Fun with topic models

```
https://gist.github.com/inkhorn/9044779#
file-recipe-analysis-r
```

# A small exercise

- Think for 5-10 minutes and try to list some (up to 5) applications of topic modeling in your discipline/topics of your interest
- Now, discuss with our neighbor and compare your and their ideas
- After about 15 minutes, let us discuss what you think are potential applications of topic models.

# Another small exercise

What do you think of these topics (and their 5 most frequent keywords)? If you are asked to evaluate this topic model now, what will you look for? Think and discuss with your classmates for 5 minutes and we will later collect all answers.

- ▶ Topic 1 : Onion, Cream, Black pepper, Milk, Cinnamon

- ▶ Topic 2: Cumin, Coriander, Turmeric, Fenugreek, Lemongrass

- ▶ Topic 3: Vanilla, Cream, Almond, Coconut, Oat

- ▶ Topic 4: Olive oil, tomato, parmesan cheese, lemon juice, garlic

- ▶ Topic 5: soy sauce, scallion, sesame oil, cane molasses, roasted sesame seed

- ▶ Topic 6: Milk, pepper, yeast, potato, lemon juice

- ▶ Topic 7: Scallion, garlic, ginger, soy bean, pepper

- ▶ Topic 8: Pepper, vinegar, onion, tomato, milk

# Some questions to ponder on:

- Coherence among the keywords for a topic (Is some word looking out of place?)
- Are there two topics that perhaps should be one?
- Can we name the topics with what we think is the group?
- Do you think the topic model learnt something about ingredients in this example?

# Topic Models in R

- different libraries: mallet, topicmodels, LDA etc.
- Textbook follows mallet
- Your Assignment 5 will use tm and topicmodels.

# Assignment 5 description

- ▶ Deadline: 31st March
- ▶ grade: 15%
- ▶ Num. questions: 1
- ▶ What to do?: Build a topic model with the given data, following given instructions, and answer questions about what you did.
- ▶ Difficulty level: moderate, but the program takes a few minutes to complete running.
- ▶ R libraries needed: tm, topicmodels

# Topic Models - The Textbook Way

- The author used mallet library to develop topic models for a corpus of novels and authors (same one he used in Chapters 11-12).
- I will follow a different method (using tm, which we used before), but I recommend you to also go through this.
- Note: You won't be an expert in text mining with one undergrad course. It is okay if you don't have a 100% understanding of this.
- The goal is to introduce you to different possible ways, give some ideas, and make you think.

# For Thursday

- I uploaded a Zip file, start with that. It contains all you need to follow Chapter 13's example.
- Attendance question: Try to follow the textbook example (materials provided in a zip file), and write your notes in the forum for 22nd March.

# Next Week

- Read this before coming to class: `https://goo.gl/L8MFfG`
- Or this: `http://www.scottbot.net/HIAL/index.html@p=19113.html`
- Optional, additional reading (for next week): `https://goo.gl/azc7Gc` (Has some math)
- Other references: `http://tidytextmining.com/topicmodeling.html`
- Attendance question for today: What are four things you need to build a topic model? - Answer can be found by reading first url.