

POIR 613: Measurement Models and Statistical Computing

Pablo Barberá

School of International Relations
University of Southern California
`pablobarbera.com`

Course website:

pablobarbera.com/POIR613/

Today

1. Solutions for last week's challenge
2. Reminder: project summary was due yesterday
 - ▶ Two-page detailed summary of project
 - ▶ Submit via Blackboard
 - ▶ Email to your peer (see next slide)
 - ▶ Feedback: 2-3 paragraphs with your reaction
 - ▶ Feedback due by October 9th
3. Other announcements:
 - ▶ No class on November 21st
 - ▶ Office hours back at regular time tomorrow
4. Topic models

Topic discovery in textual datasets

Overview of text as data methods

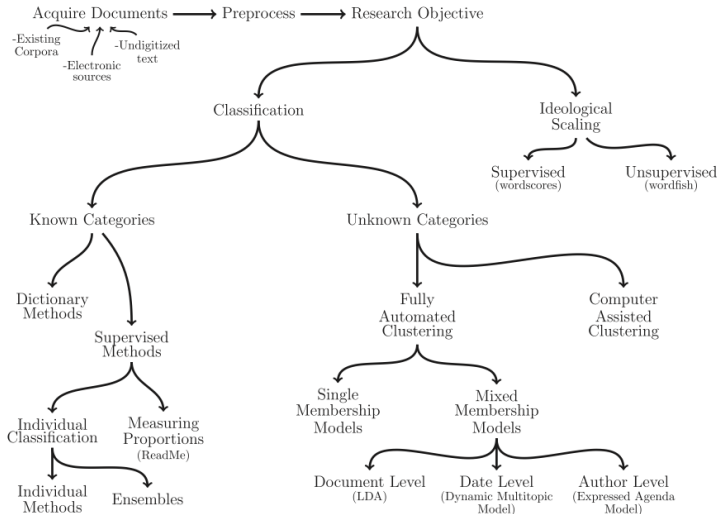


Fig. 1 in Grimmer and Stewart (2013)

Overview of techniques

▶ Descriptive analysis:

- ▶ What are the characteristics of this corpus? How do some documents compare to others?
- ▶ Keyness, collocation analysis, readability scores, Cosine/Jaccard similarity...

▶ Clustering and scaling:

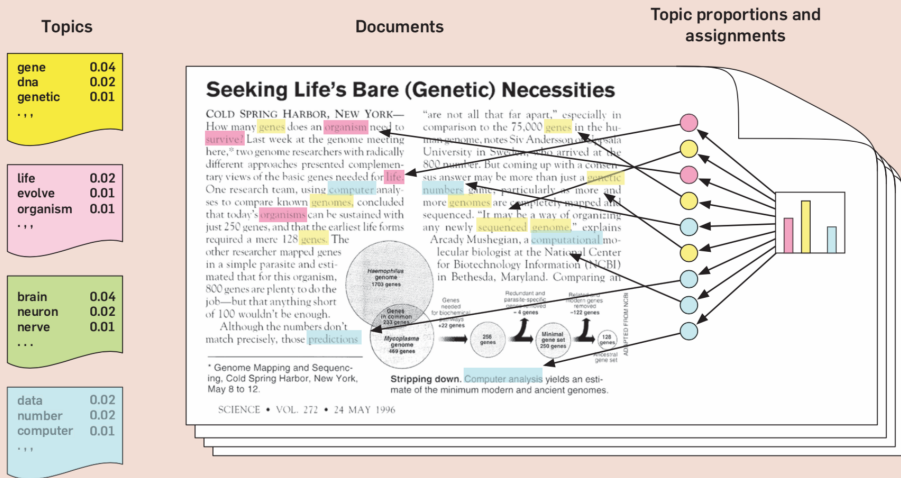
- ▶ What **groups** of documents are there in this corpus? Can documents be placed on a latent dimension?
- ▶ Cluster analysis, principal component analysis, wordfish..

▶ Topic modeling:

- ▶ What are the main **themes** in this corpus? How do different documents relate to words differently?
- ▶ LDA, STM

Latent Dirichlet Allocation

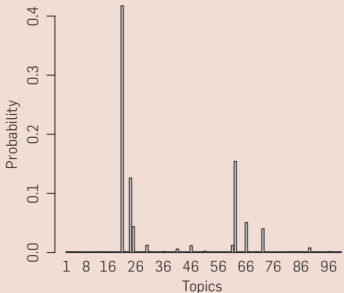
Figure 1. The intuitions behind latent Dirichlet allocation. We assume that some number of “topics,” which are distributions over words, exist for the whole collection (far left). Each document is assumed to be generated as follows. First choose a distribution over the topics (the histogram at right); then, for each word, choose a topic assignment (the colored coins) and choose the word from the corresponding topic. The topics and topic assignments in this figure are illustrative—they are not fit from real data. See Figure 2 for topics fit from data.



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

Latent Dirichlet Allocation

Figure 2. Real inference with LDA. We fit a 100-topic LDA model to 17,000 articles from the journal *Science*. At left are the inferred topic proportions for the example article in Figure 1. At right are the top 15 most frequent words from the most frequent topics found in this article.



“Genetics”

human
genome
dna
genetic
genes
sequence
gene
molecular
sequencing
map
information
genetics
mapping
project
sequences

“Evolution”

evolution
evolutionary
species
organisms
life
origin
biology
groups
phylogenetic
living
diversity
group
new
two
common

“Disease”

disease
host
bacteria
diseases
resistance
bacterial
new
strains
control
infectious
malaria
parasite
parasites
united
tuberculosis

“Computers”

computer
models
information
data
computers
system
network
systems
model
parallel
methods
networks
software
new
simulations

Latent Dirichlet Allocation

- ▶ Document = random mixture over latent topics
- ▶ Topic = distribution over n-grams

Probabilistic model with 3 steps:

1. Choose $\theta_i \sim \text{Dirichlet}(\alpha)$
2. Choose $\beta_k \sim \text{Dirichlet}(\delta)$
3. For each word in document i :
 - ▶ Choose a topic $z_m \sim \text{Multinomial}(\theta_i)$
 - ▶ Choose a word $w_{im} \sim \text{Multinomial}(\beta_{i,k=z_m})$

where:

α =parameter of Dirichlet prior on distribution of topics over docs.

θ_i =topic distribution for document i

δ =parameter of Dirichlet prior on distribution of words over topics

β_k =word distribution for topic k

Latent Dirichlet Allocation

Key parameters:

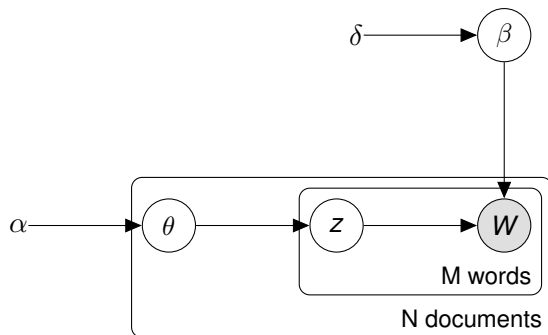
1. θ = matrix of dimensions N documents by K topics where θ_{ik} corresponds to the probability that document i belongs to topic k ; i.e. assuming $K = 5$:

	T1	T2	T3	T4	T5
Document 1	0.15	0.15	0.05	0.10	0.55
Document 2	0.80	0.02	0.02	0.10	0.06
...					
Document N	0.01	0.01	0.96	0.01	0.01

2. β = matrix of dimensions K topics by M words where β_{km} corresponds to the probability that word m belongs to topic k ; i.e. assuming $M = 6$:

	W1	W2	W3	W4	W5	W6
Topic 1	0.40	0.05	0.05	0.10	0.10	0.30
Topic 2	0.10	0.10	0.10	0.50	0.10	0.10
...						
Topic k	0.05	0.60	0.10	0.05	0.10	0.10

Plate notation



$\beta = M \times K$ matrix where β_{im} indicates $\text{prob}(\text{topic}=k)$ for word m
 $\theta = N \times K$ matrix where θ_{ik} indicates $\text{prob}(\text{topic}=k)$ for document i

Validation

From Quinn et al, AJPS, 2010:

1. Semantic validity
 - ▶ Do the topics identify coherent groups of tweets that are internally homogenous, and are related to each other in a meaningful way?
2. Convergent/discriminant construct validity
 - ▶ Do the topics match existing measures where they should match?
 - ▶ Do they depart from existing measures where they should depart?
3. Predictive validity
 - ▶ Does variation in topic usage correspond with expected events?
4. Hypothesis validity
 - ▶ Can topic variation be used effectively to test substantive hypotheses?

Example: open-ended survey responses

Bauer, Barberá *et al*, *Political Behavior*, 2016.

- ▶ Data: General Social Survey (2008) in Germany
- ▶ Responses to questions: *Would you please tell me what you associate with the term “left”? and would you please tell me what you associate with the term “right”?*
- ▶ Open-ended questions minimize priming and potential interviewer effects
- ▶ Sparse Additive Generative model instead of LDA (more coherent topics for short text)
- ▶ $K = 4$ topics for each question

Example: open-ended survey responses

Table 1: Top scoring words associated with each topic, and English translations)

Left topic 1: Parties (proportion = .26, average lr-scale value = 5.38) linke, spd, partei, linken, pds, politik, kommunisten, parteien, grünen, punks <i>the left, spd, party, the left, pds, politics, communists, parties, greens, punks</i>
Left topic 2: Ideologies (proportion = .26, average lr-scale value = 5.36) kommunismus, links, sozialismus, lafontaine, rechts, aber, gysi, linkspartei, richtung, gleichmacherei <i>communism, left, socialism, lafontaine, right, but, gysi, left party, direction, levelling</i>
Left topic 3: Values (proportion = .24, average lr-scale value = 4.06) soziale, gerechtigkeit, demokratie, soziales, bürger, gleichheit, gleiche, freiheit, rechte, gleichberechtigung <i>social, justice, democracy, social, citizen, equality, equal, freedom, rights, equal rights</i>
Left topic 4: Policies (proportion = .24, average lr-scale value = 4.89) sozial, menschen, leute, ddr, verbinde, kleinen, einstellung, umverteilung, sozialen, vertreten <i>social, humans, people, ddr, associate, the little, attitude, redistribution, social, represent</i>
Right topic 1: Ideologies (proportion = .27, average lr-scale value = 5.00) konservativ, nationalsozialismus, rechtsradikal, radikal, ordnung, politik, nazi, recht, menschen, konservative <i>conservative, national socialism, right-wing radicalism, radical, order, politics, nazi, right, people, conservatives</i>
Right topic 2: Parties (proportion = .25, average lr-scale value = 5.26) npd, rechts, cdu, csu, rechten, parteien, leute, aber, verbinde, rechtsradikalen <i>npd, right, cdu, csu, the right, parties, people, but, associate, right-wing radicalists</i>
Right topic 3: Xenophobia (proportion = .25, average lr-scale value = 4.55) ausländerfeindlichkeit, gewalt, ausländer, demokratie, nationalismus, rechtsradikalismus, diktatur, national, intoleranz, faschismus <i>xenophobia, violence, foreigners, democracy, nationalism, right-wing radicalism, dictatorship, national, intolerance, fascism</i>
Right topic 4: Right-wing extremists (proportion = .23, average lr-scale value = 4.90) nazis, neonazis, rechtsradikale, rechte, radikale, radikalismus, partei, ausländerfeindlich, reich, nationale <i>nazis, neonazis, right-wing radicalists, rightists, radicals, radicalism, party, xenophobia, rich, national</i>

Note: “proportion” indicates the average estimated probability that any given response is assigned to a topic. “average lr-scale value” is the mean position on the left-right scale (from 0 to 10) of individuals whose highest probability belongs to that particular topic.

Example: open-ended survey responses

Fig. 6: Left-right scale means for different subsamples of associations with **left** (dashed = sample mean, bars = 95% Cis)

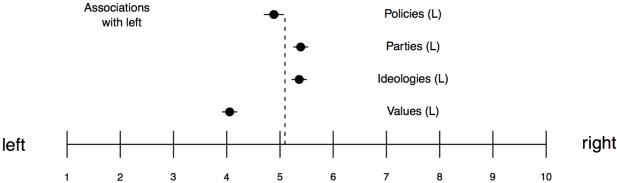
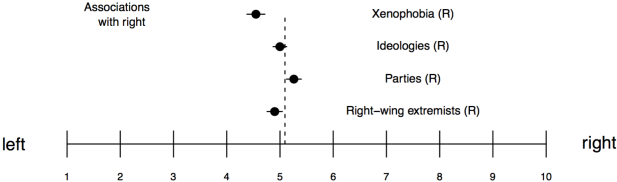


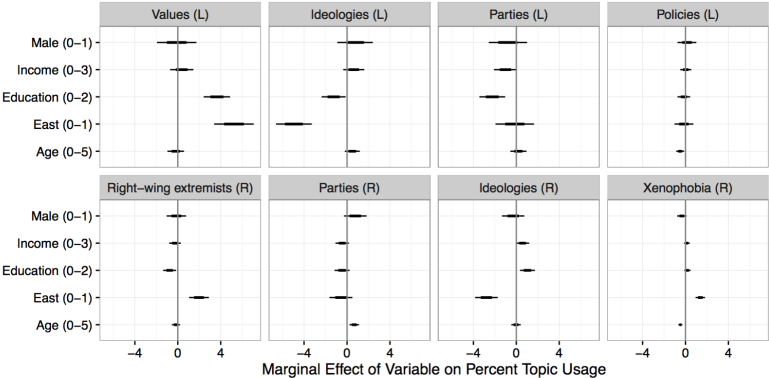
Fig. 7: Left-right scale means for different subsamples of associations with **right** (dashed = sample mean, bars = 95% Cis)



Bauer, Barberá *et al*, *Political Behavior*, 2016.

Example: open-ended survey responses

Fig. 9: Systematic relationship between associations with “left” and “right” and characteristics of respondents



Note: Each line indicates a 95% confidence interval (and 66% confidence interval in darker color) for the coefficient of eight different regressions of topic usage (in a scale from 0 to 100) at the respondent level on seven individual-level characteristics. The line on the bottom right corner (second row, second plot), for example, shows that individual a one-category change in age is associated with around one percentage point increase in the probability that the individual associated “right” with political parties.

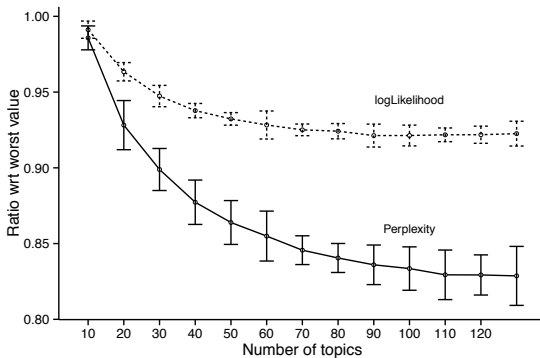
Bauer, Barberá *et al*, *Political Behavior*, 2016.

Example: topics in US legislators' tweets

- ▶ Data: 651,116 tweets sent by US legislators from January 2013 to December 2014.
- ▶ 2,920 documents = 730 days \times 2 chambers \times 2 parties
- ▶ Why aggregating? Applications that aggregate by author or day outperform tweet-level analyses (Hong and Davidson, 2010)
- ▶ $K = 100$ topics (more on this later)
- ▶ Validation: <http://j.mp/lda-congress-demo>

Choosing the number of topics

- ▶ Choosing K is “one of the most difficult questions in unsupervised learning” (Grimmer and Stewart, 2013, p.19)
- ▶ We chose $K = 100$ based on cross-validated model fit.



- ▶ **BUT:** “there is often a negative relationship between the best-fitting model and the substantive information provided”.
- ▶ GS propose to choose K based on “substantive fit.”

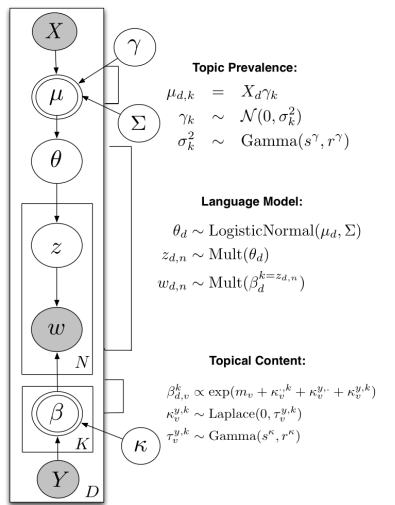
Extensions of LDA

1. Structural topic model (Roberts et al, 2014, AJPS)
2. Dynamic topic model (Blei and Lafferty, 2006, ICML; Quinn et al, 2010, AJPS)
3. Hierarchical topic model (Griffiths and Tenenbaum, 2004, NIPS; Grimmer, 2010, PA)

Why?

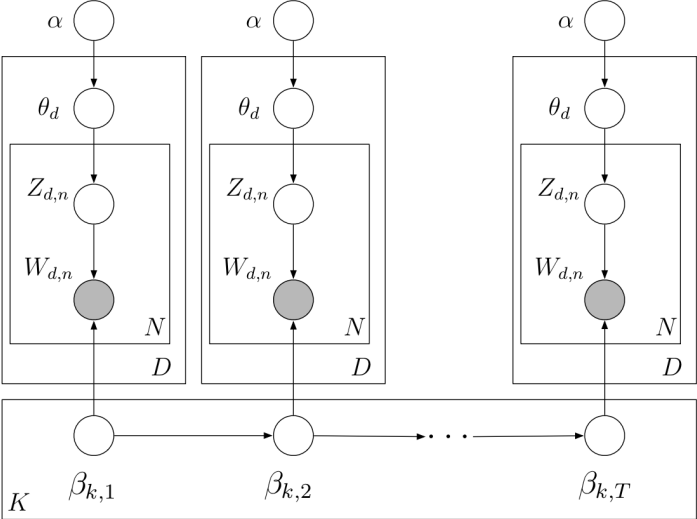
- ▶ Substantive reasons: incorporate specific elements of DGP into estimation
- ▶ Statistical reasons: structure can lead to better topics.

Structural topic model



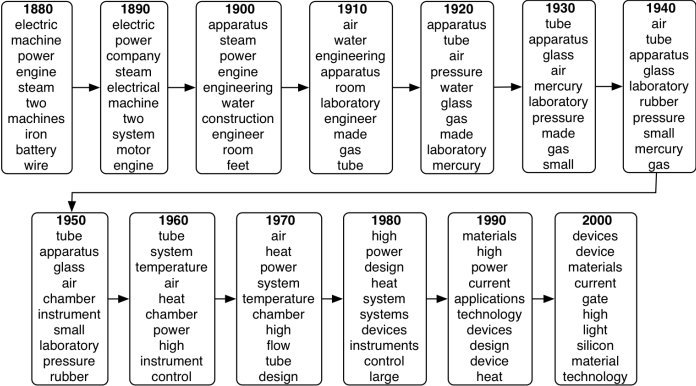
- ▶ **Prevalence:** Prior on the mixture over topics is now document-specific, and can be a function of covariates (documents with similar covariates will tend to be about the same topics)
- ▶ **Content:** distribution over words is now document-specific and can be a function of covariates (documents with similar covariates will tend to use similar words to refer to the same topic)

Dynamic topic model



Source: Blei, "Modeling Science"

Dynamic topic model



Source: Blei, "Modeling Science"

Figure 5. Two topics from a dynamic topic model. This model was fit to *Science* from 1880 to 2002. We have illustrated the top words at each decade.

