

Quantitative text analysis: Topic Models

Blake Miller

MY 459: Quantitative Text Analysis

March 15, 2021

Course website: lse-my459.github.io

1. Overview and Fundamentals
2. Descriptive Statistical Methods for Text Analysis
3. Automated Dictionary Methods
4. Machine Learning for Texts
5. Supervised Scaling Models for Texts
6. *Reading Week*
7. Unsupervised Models for Scaling Texts
8. Similarity and Clustering Methods
9. Topic models
10. Word embeddings
11. Working with Social Media

Overview of text as data methods



Outline

- ▶ Overview of topic models
- ▶ Latent Dirichlet Allocation (LDA)
- ▶ Validating the output of topic models
- ▶ Examples
- ▶ Choosing the number of topics
- ▶ Extensions of LDA

Outline

- ▶ Overview of topic models
- ▶ Latent Dirichlet Allocation (LDA)
- ▶ Validating the output of topic models
- ▶ Examples
- ▶ Choosing the number of topics
- ▶ Extensions of LDA

Topic Models

- ▶ Topic models are algorithms for discovering the main “themes” in an unstructured corpus
- ▶ Can be used to organize the collection according to the discovered themes
- ▶ Requires no prior information, training set, or human annotation – only a decision on K (number of topics)
- ▶ Most common: Latent Dirichlet Allocation (LDA) – Bayesian mixture model for discrete data where topics are assumed to be uncorrelated
- ▶ LDA provides a generative model that describes how the documents in a dataset were created
 - ▶ Each of the K *topics* is a distribution over a fixed vocabulary
 - ▶ Each document is a collection of words, generated according to a multinomial distribution, one for each of K topics

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.

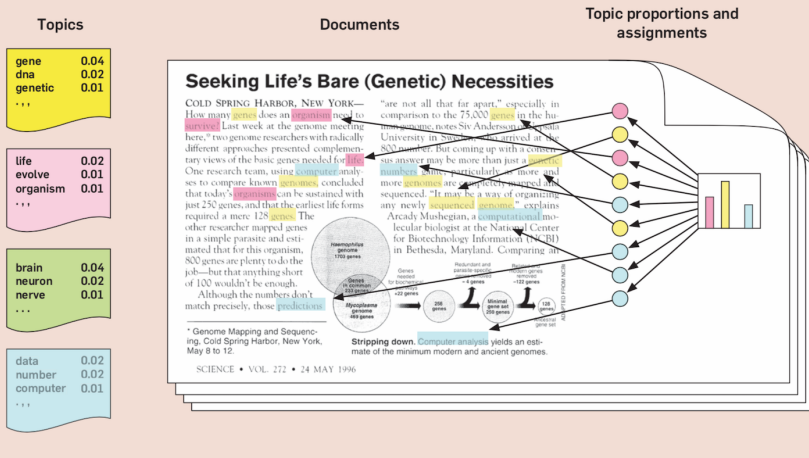


Illustration of the LDA generative process

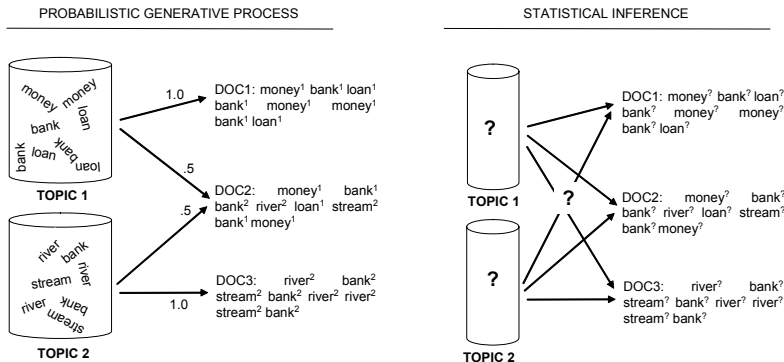


Figure 2. Illustration of the generative process and the problem of statistical inference underlying topic models

(from Steyvers and Griffiths 2007)

Topics example

Topic 247

word	prob.
DRUGS	.069
DRUG	.060
MEDICINE	.027
EFFECTS	.026
BODY	.023
MEDICINES	.019
PAIN	.016
PERSON	.016
MARIJUANA	.014
LABEL	.012
ALCOHOL	.012
DANGEROUS	.011
ABUSE	.009
EFFECT	.009
KNOWN	.008
PILLS	.008

Topic 5

word	prob.
RED	.202
BLUE	.099
GREEN	.096
YELLOW	.073
WHITE	.048
COLOR	.048
BRIGHT	.030
COLORS	.029
ORANGE	.027
BROWN	.027
PINK	.017
LOOK	.017
BLACK	.016
PURPLE	.015
CROSS	.011
COLORED	.009

Topic 43

word	prob.
MIND	.081
THOUGHT	.066
REMEMBER	.064
MEMORY	.037
THINKING	.030
PROFESSOR	.028
FELT	.025
REMEMBERED	.022
THOUGHTS	.020
FORGOTTEN	.020
MOMENT	.020
THINK	.019
THING	.016
WONDER	.014
FORGET	.012
RECALL	.012

Topic 56

word	prob.
DOCTOR	.074
DR.	.063
PATIENT	.061
HOSPITAL	.049
CARE	.046
MEDICAL	.042
NURSE	.031
PATIENTS	.029
DOCTORS	.028
HEALTH	.025
MEDICINE	.017
NURSING	.017
DENTAL	.015
NURSES	.013
PHYSICIAN	.012
HOSPITALS	.011

Figure 1. An illustration of four (out of 300) topics extracted from the TASA corpus.

(from Steyvers and Griffiths 2007)

Often K is quite large!

Outline

- ▶ Overview of topic models
- ▶ Latent Dirichlet Allocation (LDA)
- ▶ Validating the output of topic models
- ▶ Examples
- ▶ Choosing the number of topics
- ▶ Extensions of LDA

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

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

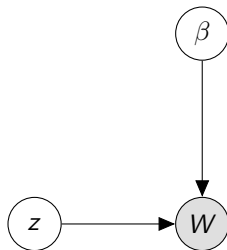
Plate notation



Plate notation

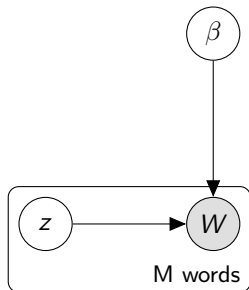


Plate notation



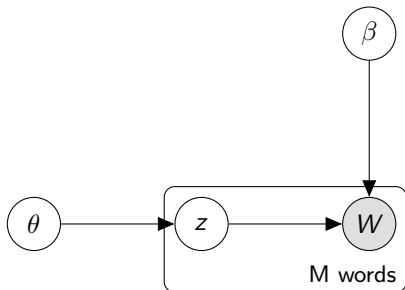
$\beta = M \times K$ matrix where β_{im} indicates $\text{prob}(\text{topic}=k)$ for word m

Plate notation



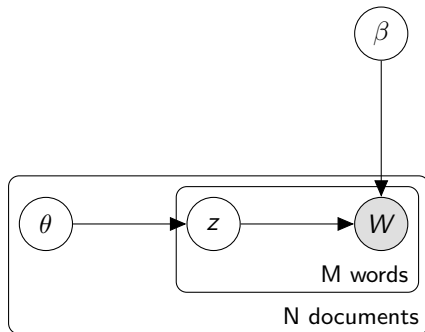
$\beta = M \times K$ matrix where β_{im} indicates $\text{prob}(\text{topic}=k)$ for word m

Plate notation



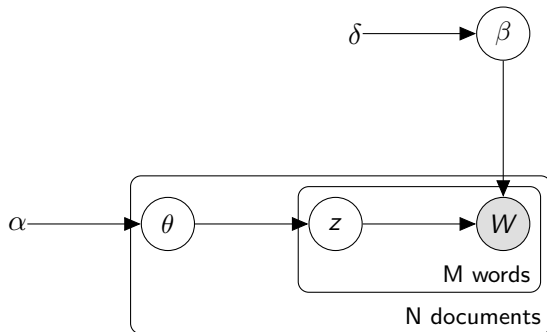
$\beta = M \times K$ matrix where β_{im} indicates $\text{prob}(\text{topic}=k)$ for word m
 $\theta = N \times K$ matrix where θ_{ik} indicates $\text{prob}(\text{topic}=k)$ for document
 i

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

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

Outline

- ▶ Overview of topic models
- ▶ Latent Dirichlet Allocation (LDA)
- ▶ Validating the output of topic models
- ▶ Examples
- ▶ Choosing the number of topics
- ▶ Extensions of LDA

Validation

From Quinn et al, AJPS, 2010:

1. Semantic validity

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?

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

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?

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?

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

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?

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

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?

Outline

- ▶ Overview of topic models
- ▶ Latent Dirichlet Allocation (LDA)
- ▶ Validating the output of topic models
- ▶ Examples
- ▶ Choosing the number of topics
- ▶ Extensions of LDA

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)

<p>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></p>
<p>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></p>
<p>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></p>
<p>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></p>
<p>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></p>
<p>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></p>
<p>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></p>
<p>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></p>

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)

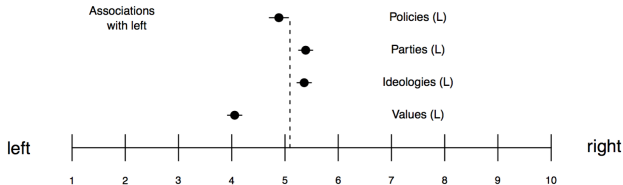
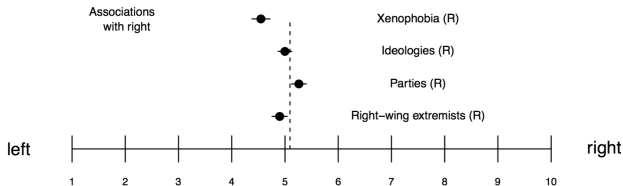
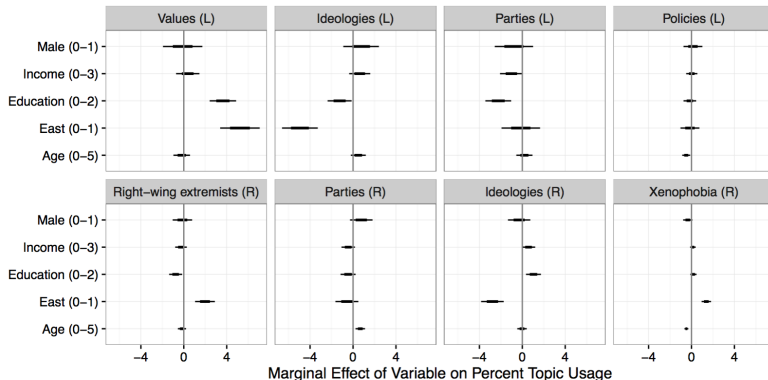


Fig. 7: Left-right scale means for different subsamples of associations with **right** (dashed = sample mean, bars = 95% CIs)



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>

Outline

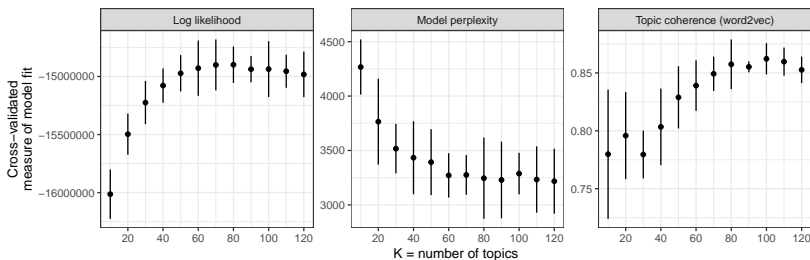
- ▶ Overview of topic models
- ▶ Latent Dirichlet Allocation (LDA)
- ▶ Validating the output of topic models
- ▶ Examples
- ▶ Choosing the number of topics
- ▶ Extensions of LDA

Choosing the number of topics

- ▶ Choosing K is “one of the most difficult questions in unsupervised learning” (Grimmer and Stewart, 2013, p.19)

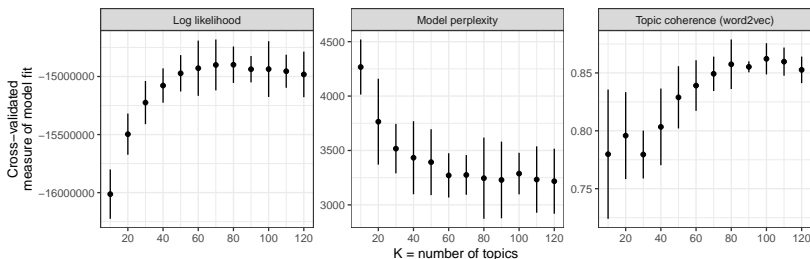
Choosing the number of topics

- ▶ Choosing K is “one of the most difficult questions in unsupervised learning” (Grimmer and Stewart, 2013, p.19)
- ▶ One approach is to decide based on cross-validated model fit



Choosing the number of topics

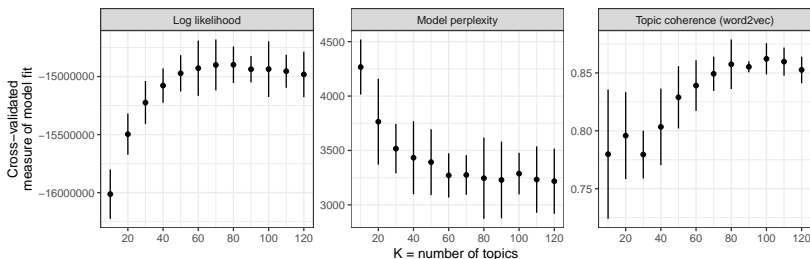
- ▶ Choosing K is “one of the most difficult questions in unsupervised learning” (Grimmer and Stewart, 2013, p.19)
- ▶ One approach is to decide based on cross-validated model fit



- ▶ **BUT:** “there is often a negative relationship between the best-fitting model and the substantive information provided”.

Choosing the number of topics

- ▶ Choosing K is “one of the most difficult questions in unsupervised learning” (Grimmer and Stewart, 2013, p.19)
- ▶ One approach is to decide 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.”

Model evaluation using “perplexity”

- ▶ can compute a likelihood for “held-out” data
- ▶ **perplexity**: can be computed as (using VEM):

$$\text{perplexity}(w) = \exp \left\{ - \frac{\sum_{d=1}^M \log p(w_d)}{\sum_{d=1}^M N_d} \right\}$$

- ▶ lower perplexity score indicates better performance

Evaluating model performance: human judgment

(Chang, Jonathan et al. 2009. “Reading Tea Leaves: How Humans Interpret Topic Models.” *Advances in neural information processing systems*.)

Uses human evaluation of:

- ▶ whether a topic has (human-identifiable) semantic coherence: **word intrusion**, asking subjects to identify a spurious word inserted into a topic
- ▶ whether the association between a document and a topic makes sense: **topic intrusion**, asking subjects to identify a topic that was not associated with the document by the model

Example

Word Intrusion

1 / 10

floppy alphabet computer processor memory disk

2 / 10

molecule education study university school student

3 / 10

linguistics actor film comedy director movie

4 / 10

islands island bird coast portuguese mainland

Topic Intrusion

6 / 10

DOUGLAS_HOFSTADTER

Douglas Richard Hofstadter (born February 15, 1945 in New York, New York) is an American academic whose research focuses on consciousness, thinking and creativity. He is best known for " , first published in

[Show entire excerpt](#)

student	school	study	education	research	university	science	learn
human	life	scientific	science	scientist	experiment	work	idea
play	role	good	actor	star	career	show	performance
write	work	book	publish	life	friend	influence	father

Example

Word Intrusion

1 / 10	floppy	alphabet	computer	processor	memory	disk
2 / 10	molecule	education	study	university	school	student
3 / 10	linguistics	actor	film	comedy	director	movie
4 / 10	islands	island	bird	coast	portuguese	mainland

Topic Intrusion

6 / 10

DOUGLAS_HOFSTADTER

Douglas Richard Hofstadter (born February 15, 1945 in New York, New York) is an American academic whose research focuses on consciousness, thinking and creativity. He is best known for " , first published in

Show entire excerpt

student	school	study	education	research	university	science	learn
human	life	scientific	science	scientist	experiment	work	idea
play	role	good	actor	star	career	show	performance
write	work	book	publish	life	friend	influence	father

- conclusions: the quality measures from human benchmarking were negatively correlated with traditional quantitative diagnostic measures!

Outline

- ▶ Overview of topic models
- ▶ Latent Dirichlet Allocation (LDA)
- ▶ Validating the output of topic models
- ▶ Examples
- ▶ Choosing the number of topics
- ▶ Extensions of LDA

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)

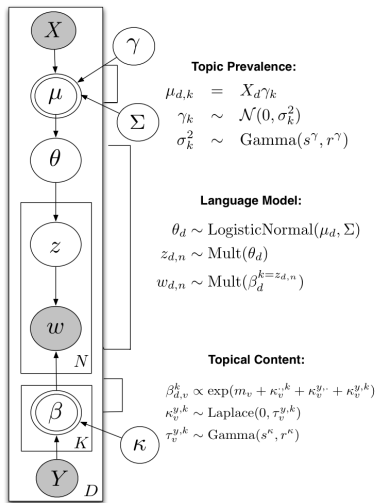
Why?

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

Structural topic model (STM)

- ▶ Basic idea: $STM = LDA + \text{Contextual Information}$
- ▶ STM provides two ways to include contextual information
 - ▶ Topic prevalence can vary by metadata (e.g. Democrats talk more about education than Republicans)
 - ▶ Topic content can vary by metadata (e.g. Democrats are less likely to use the word “life” when talking about abortion than Republicans)
- ▶ Including context improves the model:
 - ▶ more accurate estimation
 - ▶ better qualitative interpretability

Structural topic model (STM)

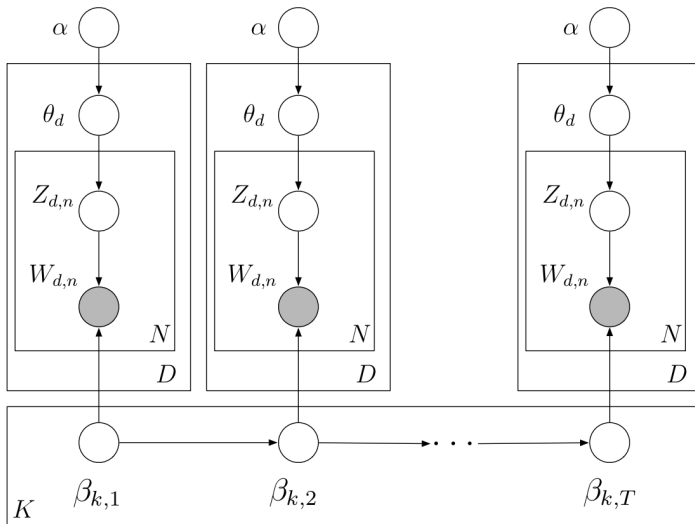


- **Prevalence:** Prior on the mixture over topics is now document-specific, and can be a function of covariates (documents with similar covariates will tend to be about the same topics)
- **Content:** distribution over words is now document-specific and can be a function of covariates (documents with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

Structural topic model (STM)

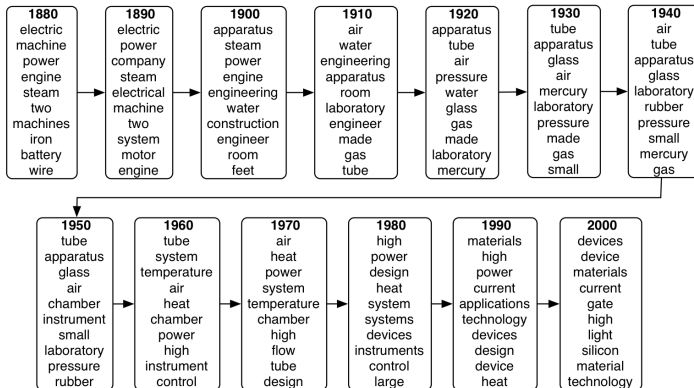
- ▶ User specifies the number of topics: K
- ▶ Observed data for standard topic models
 - ▶ Each document ($d \in 1 \dots D$) is a collection of N_d tokens
 - ▶ Each token is a particular word from a dictionary of V entries
 - ▶ Data summarized in a single matrix $D \times V$ matrix W
- ▶ Additional data for STM
 - ▶ Topic prevalence covariates: $D \times P$ matrix X
 - ▶ Topical content groups: D length vector Y
 - ▶ Latent variables
 - ▶ $D \times K$ matrix θ : proportion of document on each topic.
 - ▶ $K \times V$ matrix β : probability of drawing a word conditional on topic.
 - ▶ Low rank approximation to expected counts: $\tilde{W}_{D \times V} \sim \theta_{D \times K} \beta_{K \times V}$

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.

