

Logistics
o

Unsupervised methods
oooooooooooooooooooo

Topic models
oooooooooooooooooooo

Topic model example
ooooooo

SOSC 4300/5500: Text Analysis; Unsupervised Methods

Han Zhang

Oct 20, 2020

Logistics
o

Unsupervised methods
oooooooooooooooooooo

Topic models
oooooooooooooooooooo

Topic model example
ooooooo

Outline

Logistics

Unsupervised methods

Topic models

Topic model example

Literature Review

- Due on **Nov 3** (four weeks from now)
- Ideally linking your literature review with your final project
- Chat with me or TA during office hours

Unsupervised vs supervised

- Supervised: you have a strong *a priori* set of known categories
 - e.g., sentiments, hate speech detection, fake news, election prediction
 - They serve as good **measurement** tools
 - But requires training data to start with for supervised machine learning
- Unsupervised: you do **not** have a strong *a priori* set of known categories
 - And want the machine to find the categories for you
 - Unsupervised methods are for exploration/discovery purposes
 - The categories can be thought as a group; observations with the same category belong to the same group

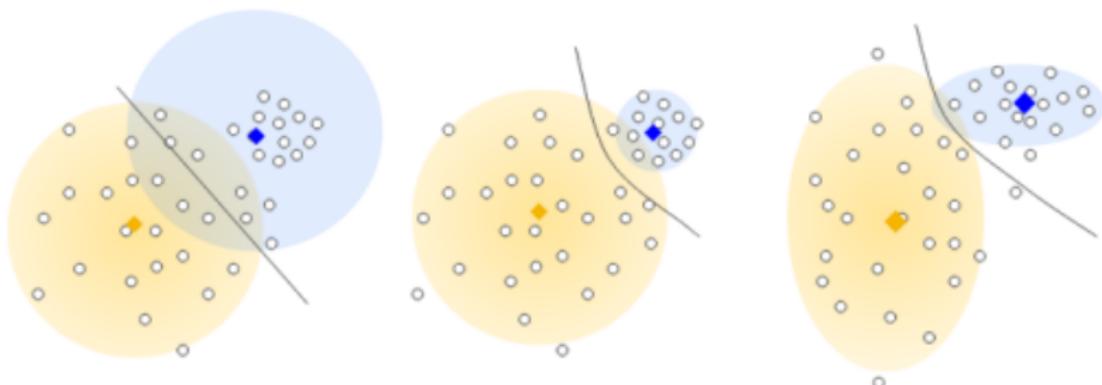
Unsupervised methods: K-means

- N observations into K clusters
- Each observation belongs to the cluster with the nearest mean (cluster centers, or **centroid**)
 - Usually the distance metric is **Euclidean distances**
 - For two observation $x = (x_1, x_2, \dots, x_m)$ and $y = (y_1, y_2, \dots, y_m)$
 - The euclidean distance is $\sqrt{\sum_{i=1}^n (x_i - y_i)^2}$
- The k-means algorithm
 - Step 0: select K initial “means” randomly
 - Step 1: associating every observation with the nearest mean
 - Step 2: the centroid of each of the K clusters becomes the new mean.
 - Repeat step 2 and step 3 until convergence
- Animation: https://commons.wikimedia.org/wiki/File:K-means_convergence.gif

K-means

- K-means is one of the simplest unsupervised methods
 - A good starting point for most data
 - But there are shortcomings:
 - Need to select K beforehand
 - This is a general problem for many other clustering methods
 - Sensitive to outliers
 - Sensitive to imbalanced data
 - Does not work very well if the number of dimensions is very large (**curse of dimensionality** again)

K-means and imbalanced data



Plain k-means

Varying widths across
clusters

Varying widths across
clusters & dimensions

- <https://developers.google.com/machine-learning/clustering/images/KmeansGeneralization.svg>
 - Most off-the-shelf packages do not allow you to vary the width of each cluster

K-means and curse of dimensionality

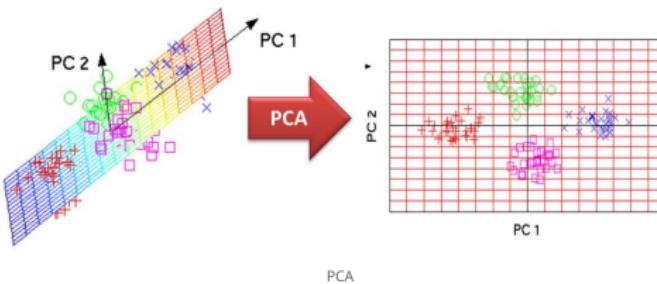
- K-means needs to calculate **distance** between data and K centroids to find the nearest cluster
- But when the dimension of data is high, the variance between distances decreases
- Then k-means becomes less effective at distinguishing between examples
- <https://developers.google.com/machine-learning/clustering/images/KmeansGeneralization.svg>

Dimensionality reduction

- If the dimension of data M on the same scale of N , or even larger than N
- It's better to perform **dimensionality reduction**, before running any clustering algorithms
- Unsupervised method is itself a dimension reduction
 - We reduce the observation to a single category

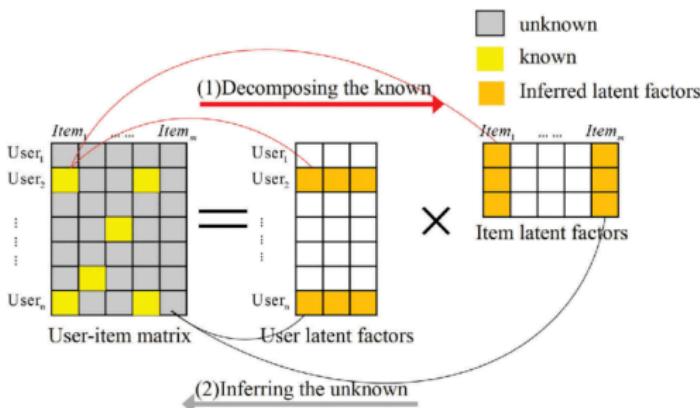
Dimensionality reduction: PCA

- PCA: principal component analysis
 - One of the simplest and widely used dimensionality reduction techniques
 - It seeks to project original data onto dimensions that best preserve the variances in the data



Dimensionality reduction: Nonnegative Matrix Factorization

- PCA only reduces the dimension of rows; not considering columns
- An two-dimensional extension of PCA, considering both features of rows and columns



Nonnegative matrix factorization in social sciences

- Matrix factorization is the idea behind item response theory (in educational testing)
 - Raw data:
 - Row are students
 - Columns are answers
 - Not every one answer the same set of questions, so there are missing entries
 - Nonnegative matrix factorization
 - Row: student ability (may be multi-dimensional)
 - Column: item difficulty

Nonnegative matrix factorization in social sciences

- Matrix factorization is also behind the **ideal point estimation** (in political science)
 - That is, infer ideological position of politicians and policies

Roll Call	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	ETC.	
Johnson (D-US)																		N	Y			
Sparkman (D-AL)	N	Y	N	N	Y	N	Y		Y	N	N	Y	N				N	N	Y	Y		
Hill (D-AL)	N	Y	N	N	Y	N	Y	Y	Y	Y	Y	Y	Y			Y	Y	Y	N	N	Y	Y
Gruening (D-AK)	N	Y	N		Y	N	Y	N	N	N	N	Y	Y			N		Y	Y			
Bartlett (D-AK)	N	Y	N	N	Y	Y	Y	N	N	N	Y	N	Y	Y	Y	Y	Y	Y	N	N	Y	
Hayden (D-AZ)	N		N	Y	Y	N	Y	Y	N	N	Y	N	Y	N	Y	Y	Y	N	N	Y	Y	
Fannin (R-AZ)	N	Y	N	Y	Y	N	Y	Y	N	Y	Y	N	Y	Y	Y	Y	N	N	N	N	N	
Fulbright (D-AR)	N	Y	N		Y	N	Y	N									N	Y				
McClellan (D-AR)	N	Y	N	N	Y	N	Y	Y	Y	N	Y	Y		Y	Y	Y	N	N	Y			
Kuchel (R-CA)	Y	N	Y	Y	Y	N	Y	N	N	N	N	N		Y	N	Y		N	Y	N		
Murphy (R-CA)	N	Y	N	Y	Y	N	Y	Y	N	N	Y	Y	Y	N	Y	Y	N	N	N	N		

- Fancy plots like the below actually use estimated latent points

Logistics

Unsupervised methods

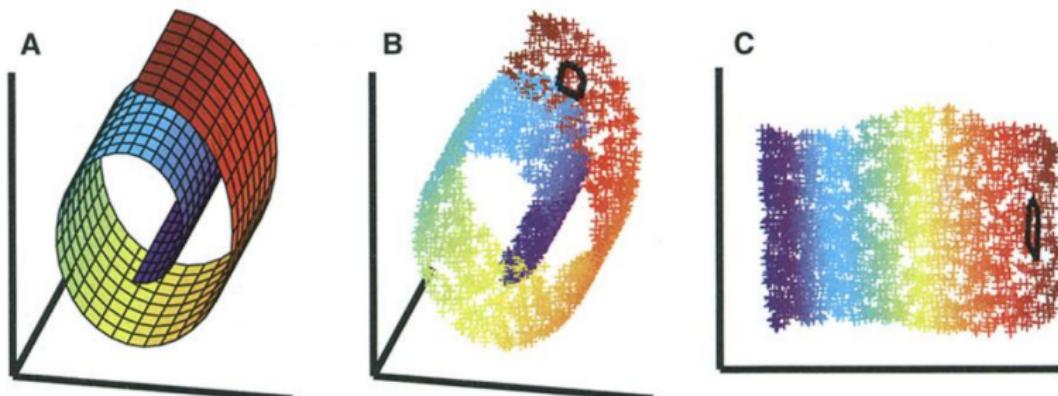
Topic models

Topic model example

Example

Dimensionality reduction: non-linear data

- PCA and nonnegative matrix factorization are linear methods
- They will not work if your data are nonlinear
-

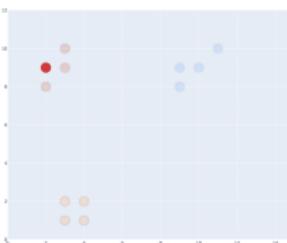


Dimensionality reduction: t-SNE

- t-SNE: t-distributed stochastic neighbor **embedding**
- Laurens van der Maaten and Geoffrey Hinton, *Visualizing Data using t-SNE*, Journal of Machine Learning Research **9** (2008), no. Nov, 2579–2605
- t-SNE is suited if you have nonlinear data
- And in particular, if you want to **visualize** them
- Note: this one can be even more slower than PCA

Intuition of t-SNE

- We can use the nearest data points to predict a single data X_i



- Similarly, if we map X_i to a lower-dimensional representation Z_i (called **embedding** here)
 - And do so for other data points X_j
- Then using Z_j to predict Z_i will also have a good prediction performance
- Later you will see that **word embedding** is essentially the same idea

Back to unsupervised clustering: Choosing K

- How do you choose K is one of the most challenging part of unsupervised methods
- Three general solutions:
- Data-driven method
- Theory-driven method
- Or combining both

Unsupervised: hierarchical agglomerative clustering

- **Bottom-up approach:**
 - Each observation starts in its own cluster
 - And pairs of clusters are merged as one moves up the hierarchy.
- Animation: <https://i.gifer.com/80Gy.mp4>
- Pros:
 - No need to specify K
 - Easy visualization
- Cons:
 - Slow!

Unsupervised methods

- There are many other clustering methods
- Their usage case can be viewed here
- <https://scikit-learn.org/stable/modules/clustering.html>

A complete workflow (Grimmer and Stewart, 2013)

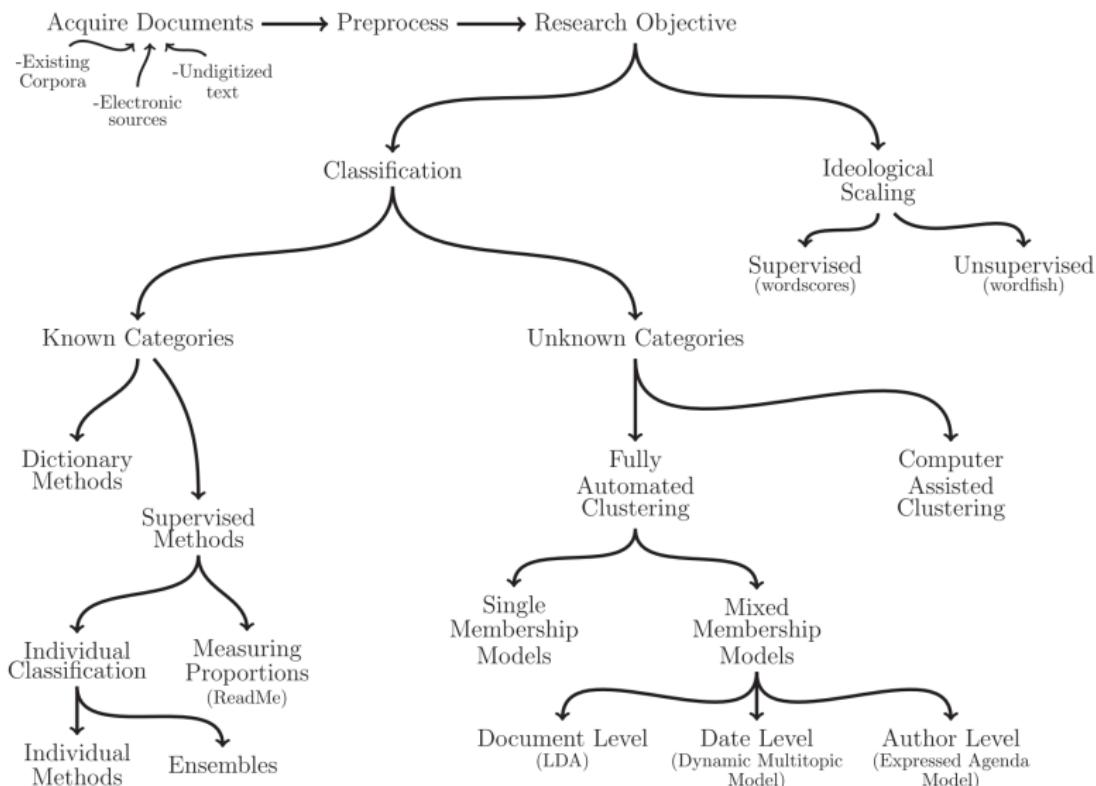


Fig. 1 An overview of text as data methods.

Topic models

- A special type of unsupervised methods designed for discovering the main topics for text documents (especially in the format of document-term matrix)
- Each document can belong to multiple topics: **mixed membership model**
- Unsupervised; requires no prior information, training set, or human annotation
 - Need to make a decision on K (number of topics)
- The most basic and the common topic model:
- **Latent Dirichlet Allocation (LDA)**
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan, *Latent Dirichlet Allocation*, J. Mach. Learn. Res. **3** (2003), 993–1022

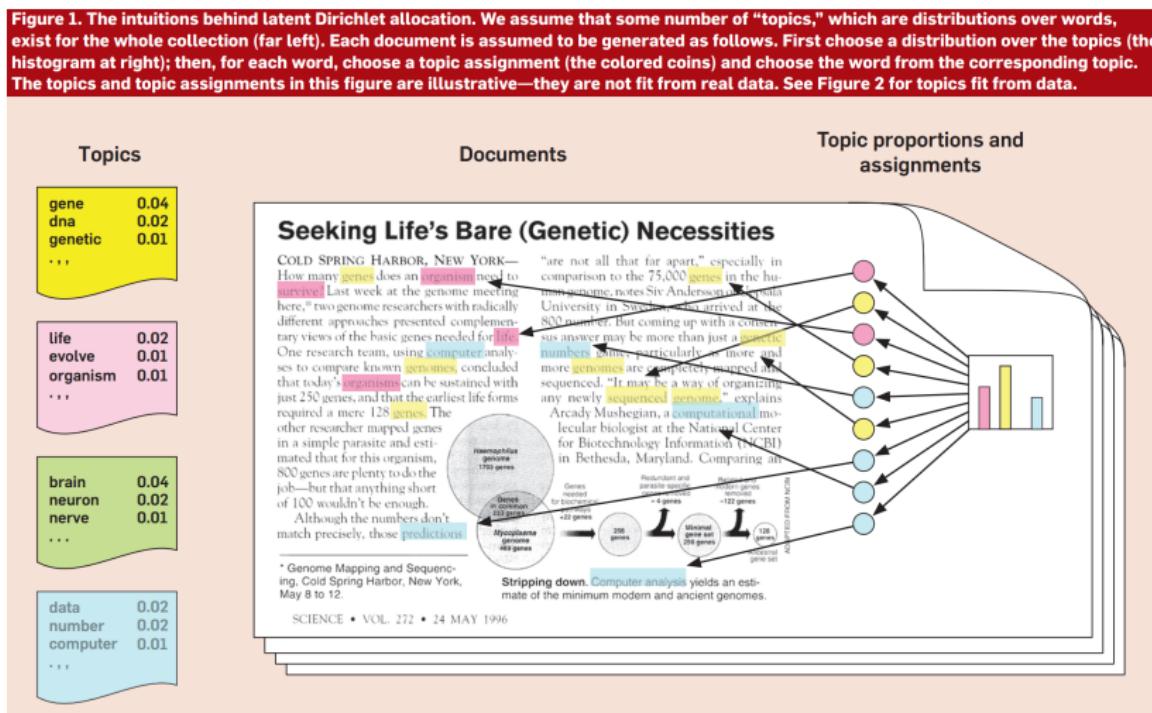
Why not simpler unsupervised methods?

- In principle, you can use K-means (or others) for clustering documents
- But their performance may be bad
 - K-means does not work very well with document-term matrix, which is high-dimensional by design.
 - The simplest K-means algorithm assumes single membership:
 - Each document belongs to a single category
 - Not realistic for text documents that often discuss several topics

LDA

David M. Blei, *Probabilistic Topic Models*, Commun. ACM 55 (2012), no. 4, 77–84

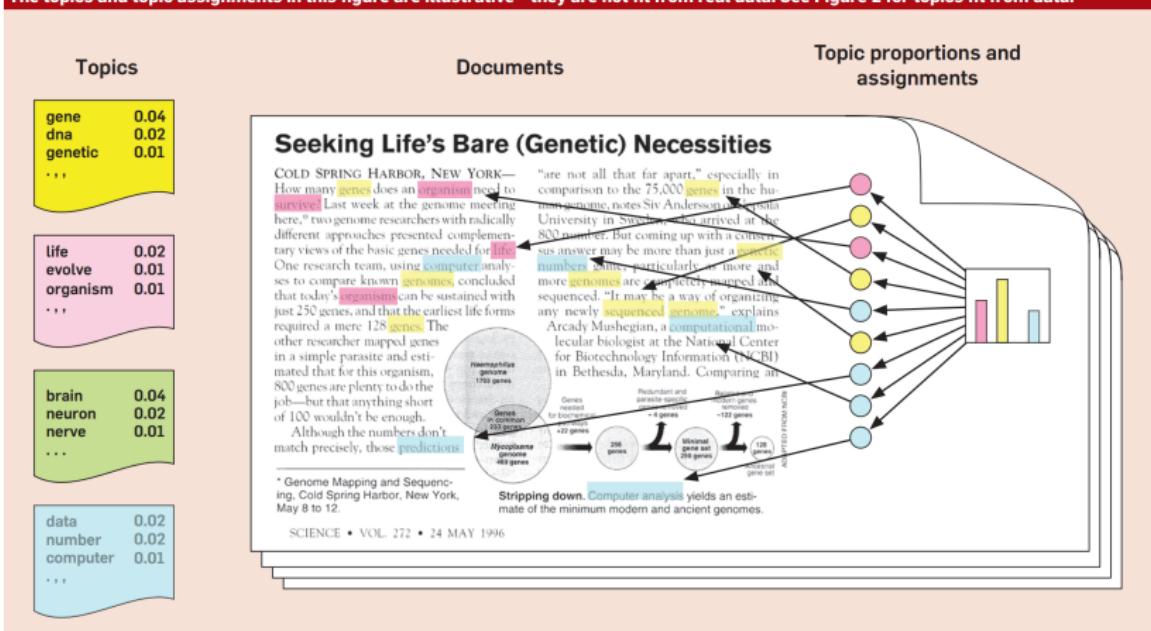
Figure 1. The intuitions behind latent Dirichlet allocation. We assume that some number of “topics,” which are distributions over words, exist for the whole collection (far left). Each document is assumed to be generated as follows. First choose a distribution over the topics (the histogram at right); 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.



Documents

- Each document is conceptualized as a probability distribution over topics

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.



Topics

- Each topic is defined as a probability distribution over words/n-grams
- Mark Steyvers and Tom Griffiths, *Probabilistic topic models*, Handbook of latent semantic analysis **427** (2007), no. 7, 424–440

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.

Word Polysemy

- Each word can belong to multiple topics
 - It's hard to achieve this with dictionary methods
- Mark Steyvers and Tom Griffiths, *Probabilistic topic models*, Handbook of latent semantic analysis **427** (2007), no. 7, 424–440

Topic 77

word	prob.
MUSIC	.090
DANCE	.034
SONG	.033
PLAY	.030
SING	.026
SINGING	.026
BAND	.026
PLAYED	.023
SANG	.022
SONGS	.021
DANCING	.020
PIANO	.017
PLAYING	.016
RHYTHM	.015
ALBERT	.013
MUSICAL	.013

Topic 82

word	prob.
LITERATURE	.031
POEM	.028
POETRY	.027
POET	.020
PLAYS	.019
POEMS	.019
PLAY	.015
LITERARY	.013
WRITERS	.013
DRAMA	.012
WROTE	.012
POETS	.011
WRITER	.011
SHAKESPEARE	.010
WRITTEN	.009
STAGE	.009

Topic 166

word	prob.
PLAY	.136
BALL	.129
GAME	.065
PLAYING	.042
HIT	.032
PLAYED	.031
BASEBALL	.027
GAMES	.025
BAT	.019
RUN	.019
THROW	.016
BALLS	.015
TENNIS	.011
HOME	.010
CATCH	.010
FIELD	.010

LDA algorithm

- Imagine how a computer writes a document with 5,000 words
 1. choose a topic according to the topic distribution (e.g., 0.8 prob of *economy* and 0.2 prob of *politics*)
 2. choose a word according to the topic's word distribution
 3. repeat Step 1 and 2 until you have selected 5,000 words

LDA algorithm

Mark Steyvers and Tom Griffiths, *Probabilistic topic models*,
Handbook of latent semantic analysis 427 (2007), no. 7, 424–440

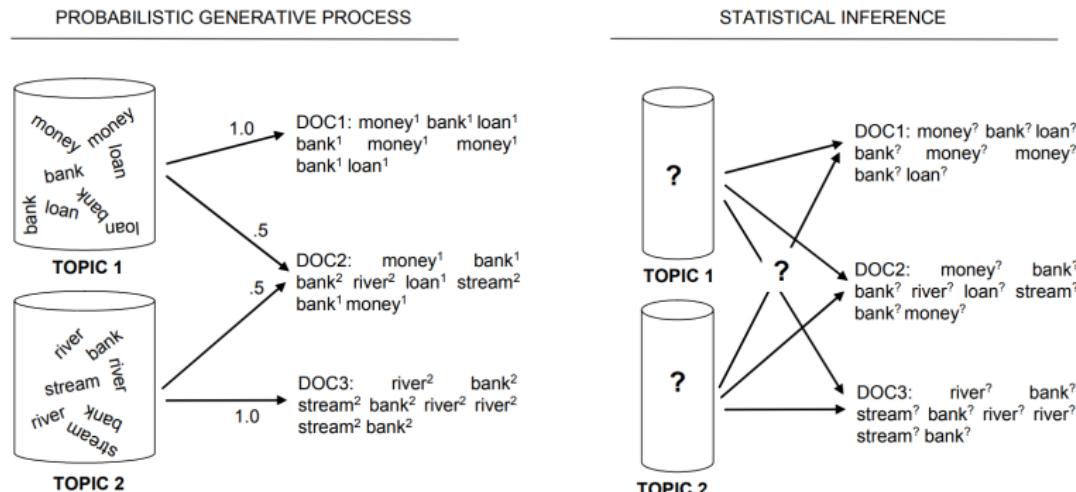


Figure 2. Illustration of the generative process and the problem of statistical inference underlying topic models

LDA algorithm

- *Advanced topic*; it's okay if you do not fully get contents on this slide
- Nearly all others you have learned are **discriminate** statistical model
- LDA is a **generative** statistical model
- Differences in math:
 - Discriminative model directly maps features to outcome:
 $P(Y|X)$
 - E.g., linear regression
 - $p(Y | X) = N(\beta X, \sigma^2)$, where N is a standard normal distribution
 - But we do know know Y yet!
 - Generative model: model $P(X|Y)$ instead, using Bayes' rule
 - Assumes that we know $P(Y)$
 - And it is easy to calculate $P(X|Y)$

Math of LDA

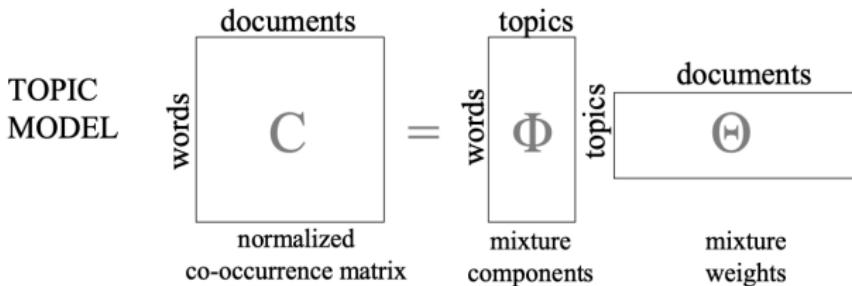
- Advanced topic; it's okay if you do not fully get contents on this slide
- Here, topics are latent outcome Y ; and X is the document-term matrix
- A discriminative model will directly model $P(Z|w)$
- What LDA does:
-

$$p(w_{id}) = \sum_{j=1}^K P(w_i|Z_i=j)P(z_i=j)$$

- $p(w_{id})$ is the probability of observing word i in document d
- $P(w_{id}|Z_d=j) = \phi^j$: word probability in topic j
- $P(Z_d=j) = \theta^d$: topic j 's probability of document d

Matrix version of LDA

- Advanced topic; it's okay if you do not fully get contents on this slide
- A different way to look at LDA is that it decompose the document-term matrix into the produce of the following two:
 - term-topic matrix Φ
 - each element is the ϕ in the previous slide)
 - topic-document matrix Θ
 - each element is the θ in the previous slide)



Choosing the number of topics

- Topic models: easy to run because no labels needed, but requires significant care in **validation**
- Choosing K is “one of the most difficult questions in unsupervised learning” (Grimmer and Stewart, 2013, p.19)
- In supervised methods, we choose parameter values based on whether they improve the prediction performance
- In unsupervised methods, there is no such luck
 - How do I know whether I should choose 5 or 10 topics?

Choosing the number of topics

- Two general approaches of choosing parameter values in LDA
- Data-driven method:
 - **perplexity**; how well the LDA model predicts a new document in the validation set
 - The original metric used in Blei et al., 2013
 - It's the inverse probability of the validation set, normalized by the number of words
 - $$PP(W) = p(w_1 w_2 \cdots w_V)^{1/n} = \sqrt[n]{\frac{1}{p(w_1 \cdots w_V)}}$$
 - The **lower** the perplexity, the **better** the model
- Cross-validation again:
 - choose K that minimizes the perplexity on validation set
 - or other metrics; more later in this class

Choosing the number of topics

- Substantive fit: human in the loop; requires domain-knowledge
- Jonathan Chang, Sean Gerrish, Chong Wang, Jordan Boyd-Graber, and David M. Blei, *Reading Tea Leaves: How Humans Interpret Topic Models*, NIPS 2009
 - Human often disagree with the model chosen by reducing perplexity

Human-Validations

- Grimmer and Stewart, 2013
- Semantic validity:
 - Do the topics identify coherent groups of documents?
- Convergent/discriminant validity
 - Do the topics match existing measures where they should match?
 - Do they depart from existing measures where they should depart?
- Predictive validity
 - Does variation in topic usage correspond with expected events?

Semantic validity

- Chang et al., 2009
- Word intrusion:
 1. select 5 words with the highest probabilities in a topic
 2. select another word that has a low probability in the topic, but high prob in other topics. This is an *intruder* word.
 3. present the 6 words to a human coder, and see if he/she can easily picks up the intruder word.
- e.g., {dog, cat, horse, apple, pig, cow}
 - Easily see that *apple* is an intruder
 - because {dog, cat,horse, pig, cow} make sense together as an animal topic
- You can compare two models on their word intrusion scores

Logistics
o

Unsupervised methods
oooooooooooooooooooo

Topic models
oooooooooooooooooooo●oo

Topic model example
oooooooo

Convergent validity

- Give each topic a label
- Ask human coders to read a sample of document and assign them a label
 - but do **not** show them words in a topic, of course
- And see if the human coding agrees with topic modeling results

Table 4 An example of topic labeling

<i>Description</i>	<i>Discriminating words</i>
Iraq War	Iraq, iraqi, troop, war, sectarian
Honorary	Honor, prayer, remember, fund, tribute
Fire Department Grants	Firefight, homeland, afgp, award, equipment

Predictive validity

- Does variation in topic usage correspond with expected events?

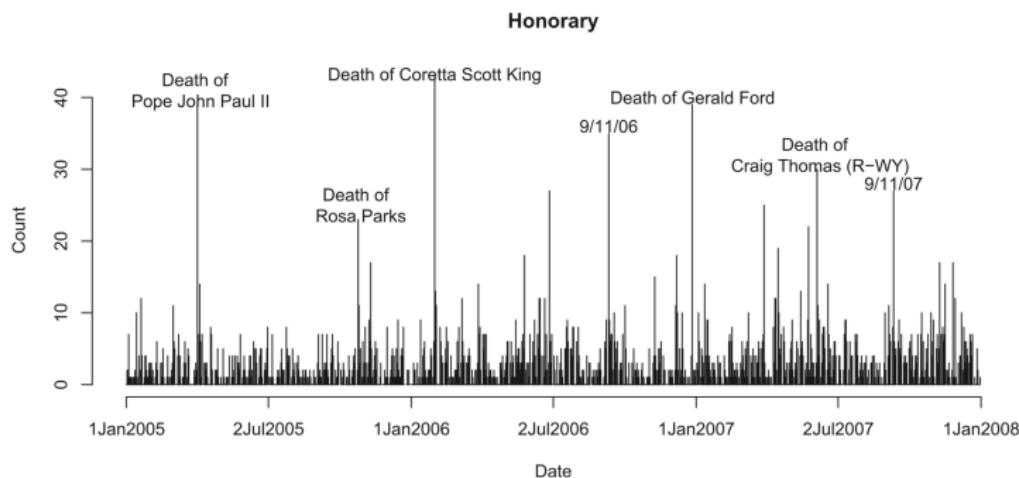
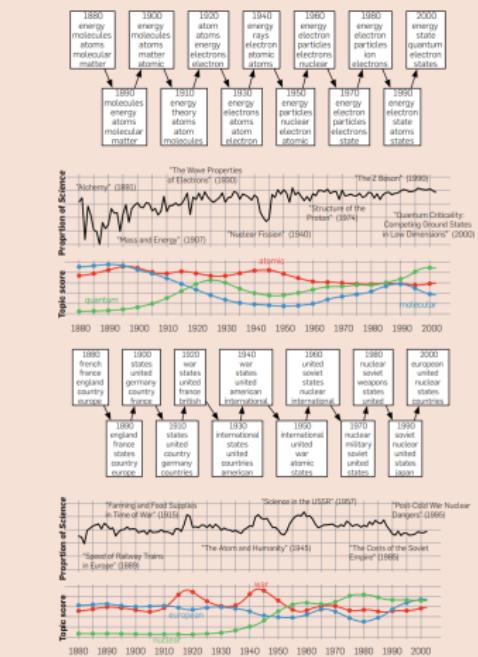


Fig. 4 Predictive validity of topics.

Topic modeling for discovery (Blei et al., 2012)

- Supervised methods categories *a priori*; it will miss many information

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.



Set up

- Pablo Barberá, Andreu Casas, Jonathan Nagler, Patrick J. Egan, Richard Bonneau, John T. Jost, and Joshua A. Tucker, *Who Leads? Who Follows? Measuring Issue Attention and Agenda Setting by Legislators and the Mass Public Using Social Media Data*, American Political Science Review **113** (2019), no. 4, 883–901
- Questions:
 - Do Congressmen follow issues raised by the public
 - Or, public follows issues raised by Congressmen?
- To answer these questions, we need to measure **attention to issues** of political discussions

Data

TABLE 1. Description of the Tweets in the Dataset

Group	N	Avg	Min	Max	Tweets
House Republicans	238	1,215	70	8,857	267,311
House Democrats	207	1,177	113	5,993	222,491
Senate Republicans	46	1,532	73	6,627	67,412
Senate Democrats	56	1,616	150	10,736	87,307
Random sample	25k	465	1	8,926	11,316,396
Informed public	10k	948	100	5,861	9,487,382
Republican supporters	10k	1,091	100	8,804	10,911,813
Democratic supporters	10k	1,306	100	5,122	13,058,947
Media outlets	36	7,803	8	15,858	273,121

Note: Period of analysis: January 1, 2013, to December 31, 2014. *N* corresponds to the number of Twitter accounts in each sample. *Avg*, *Min*, and *Max* correspond to the average, minimum, and maximum number of tweets, respectively, sent by individual users in each group during the whole period of analysis. *Tweets* corresponds to the total number of tweets sent by all users in each group during the period of analysis.

Topic modeling

- Use LDA to identify issues: each issue is a topic
- Choose 100 topics; based on minimizing perplexity (data-driven approach)
- Why they do not choose supervised methods? What are their arguments?
 1. too many categories; requires too many labeled documents
 - Say 500 documents per each category; that is 50,000; not a small number
 2. they do passed the convergent validity check

Topics

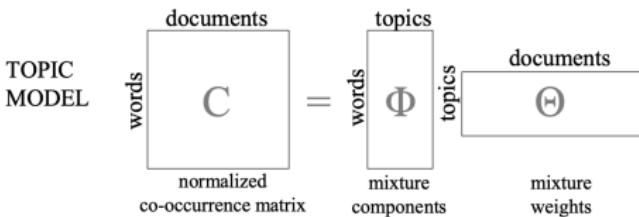
TABLE 2 List of Political Issues

Topic Number	Label	Topic number	Label
3	Investigation of Benghazi attack	50	Climate change
7	100 days of #BringBackOurGirls campaign	51	Lame duck congress
9	Gender wage gap	53	Minimum wage
12	Republican issues Spring 2013	58	Affordable Care Act
14	Marriage equality	62	Border crisis in Texas
15	Gun violence	63	Obamacare (employer mandate)
16	Abortion (pro-life)	64	FAA furloughs cause flight delays
18	Veteran affairs delays scandal	66	Malaysia Airlines crash in Ukraine
20	NSA surveillance scandal	67	Comprehensive immigration reform
23	#BringBackOurGirls campaign	70	#MiddleClassFirst campaign
28	Employment Non-Discrimination Act	75	Military Justice Improvement Act
32	Islamic state	81	Poverty (SNAP program)
33	Use of military force in Syria	83	Twenty-first century cures initiative
36	Ebola	85	Unemployment insurance
37	Social security	88	IRS scandal
39	Keystone XL pipeline	89	Obamacare (website and implementation)
41	Immigration (border security)	93	Jobs bills omnibus
43	Executive action on immigration	96	Violence Against Women Act
46	Unemployment numbers reports	97	Protests in Ukraine and Venezuela
47	Paul Ryan budget proposal	99	CIA detentions and interrogations report
48	Black history month	100	#ObamacareInThreeWords campaign
(101)	Student debt	(102)	Hobby lobby supreme court decision
(103)	Budget discussion	(104)	2013 government shutdown

Note: The topic number in parentheses indicate issues that have been created *ad hoc* by merging very similar topics from the topic model.

Measure attention to issues

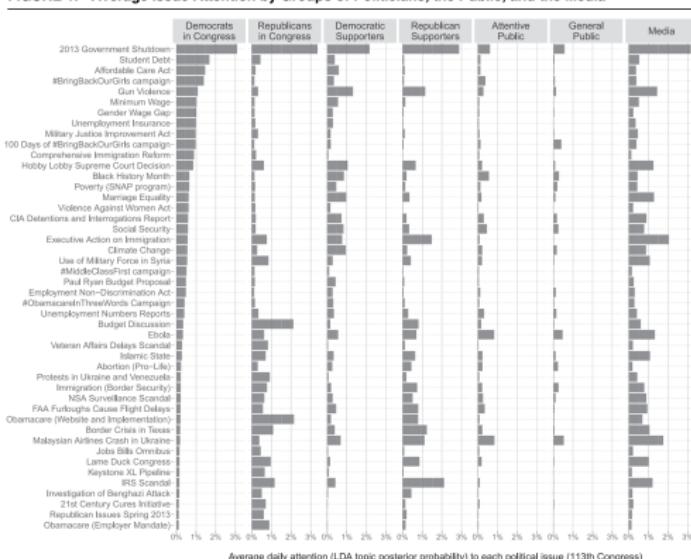
- After getting a list of issues
- They measure **attention** to issues as the daily posterior LDA probabilities for each group
 - They basically mean $P(Z_d = j) = \theta^d$: topic j 's probability of document d
- Or the column means of the Θ matrix (topic-document matrix)



Measure attention to issues

- Result: Y_{ijt} ; proportion of topic i for group j at time t

FIGURE 1. Average Issue Attention by Groups of Politicians, the Public, and the Media



Statistical analysis

- Topic modeling help the authors to obtain the key independent and dependent variables Y_{ijt}
- To answer their questions:
 - Do Congressmen follow issues raised by the public
 - Or, public follows issues raised by Congressmen?
- A series of vector auto regression regressions (a standard model dealing with time-series data)
 - Basically, can Y_{ijt} predict $Y_{i,j',t+1}$?
 - Or in plain language, can public's issue attention predict Congressmen's future issue attention
 - Or vice versa

Revisiting unsupervised vs supervised methods

- For analyzing text data, unsupervised methods (especially topic modeling) are perhaps used by frequently in social sciences
 - Because they allows discovery of new topics over time (theoretical reason)
 - And there is no need gather thousands of training documents (practical consideration)
- But remember topic models are for exploration; they are not designed to predict Y given X
- If your goal is prediction, and the outcome is designed very clearly
 - and you know the outcome do not change that much
- Then using topic models for predictions is nearly always going be performing worse than supervised methods
- Choose supervised/unsupervised methods depending on your tasks

Logistics
o

Unsupervised methods
oooooooooooooooooooo

Topic models
oooooooooooooooooooo

Topic model example
oooooooo●

Some good resources for more

- Text as Data conference
<https://www.textasdata2019.net/>