Unsupervised Learning and Topic Modeling

Yoav Goldberg

Bar Ilan University

(with slides from David Blei, Zornitsa Kozerava)

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● □ ● ● ● ●

1/1

Previously

Supervised Learning

- Get labeled training data
- Represent data as (features, label) pairs
- Train a classifier / model to predict labels based on features

Today

- What if we don't have training data?
- Can we still do something useful?

Unsupervised Learning

Things we can do without labeled data

Option 1: "Naturally occurring" labels / bootstrap

- Be creative and find data which can be used as labels.
 - e.g., we want to identify paragraphs. Maybe some website indicate this via their HTML tags?
- Automatically create your own training set
 - Write simple rule-based system to collect easy examples
 - high precision, low recall
 - Use the easy examples as training data
 - Hope it will generalize well.
 - Careful not to overlap your features with the rules too much!

Option 1: "Naturally occurring" labels / semi-supervised

- Be creative and find data which can be used as labels.
 - Want to identify sentiment? Look at tweets with happy and sad emojis.

Option 1: "Naturally occurring" labels / semi-supervised

- Be creative and find data which can be used as labels.
 - Want to identify sentiment? Look at tweets with happy and sad emojis.
 - what are the pros and cons here?

Option 1: "Naturally occurring" labels / semi-supervised

- Be creative and find data which can be used as labels.
 - Want to identify sentiment? Look at tweets with happy and sad emojis.
 - what are the pros and cons here?
- Can also use the *proxy* naturally occurring data for representation learning.
 - The Felbo et al 2017 paper. (next slides)

Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm

Bjarke Felbo¹, Alan Mislove², Anders Søgaard³, Iyad Rahwan¹, Sune Lehmann⁴

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Use Emoji Prediction to learn good representations for sentiment

Table 2: The number of tweets in the pretraining dataset associated with each emoji in millions.

😂 233.7	82.2	** 79.5	11 78.1	2 60.8	C) 54.7	5 4.6	51.7	d 50.5	44.0	9 .5	😌 39.1	55 34.8	34.4	😔 32.1	28.1
4 24.8	<u>©</u> 23.4	21.6	21.0	20.5	6 20.3	1 9.9	😌 19.6	18.9	😯 17.5	U 17.0	16 .9	Jjj 16.1	() 15.3	15.2	<mark>ප</mark> 15.0
••• 14.9	14.3	 14.2	<u>್</u> ಷ 14.2	😧 12.9	ര്ട്ടു 12.4	99 12.0	• <u>,</u> 12.0	8 11.7	e 11.7	W 11.3	~~ 11.2	11.1	••• 11.0	↓ 11.0	10.8
10.2	U)) 9.6	2 9.5	6 9.3	1 9.2	<mark>8.9</mark>	§ 8.7	0 8.6	6 8.1	6 .3	6 .0	5 .7	9 5.6	 5.5	6 .4	22 5.1

▲□▶ ▲□▶ ▲臣▶ ▲臣▶ 三臣 - のへで

Table 1: Example sentences scored by our model. For each text the top five most likely emojis are shown with the model's probability estimates.

I love mom's cooking

I love how you never reply back ..

I love cruising with my homies

I love messing with yo mind!!

I love you and now you're just gone ..

This is shit

This is the shit

2	U		6	\heartsuit	
49.1%	8.8%	3.1%	3.0%	2.9%	
0 0 0	-	25	•••	()	
14.0%	8.3%	6.3%	5.4%	5.1%	
.	4	¥.	00	100	
34.0%	6.6%	5.7%	4.1%	3.8%	
	T	53	13		
17.2%	11.8%	8.0%	6.4%	5.3%	
(?)	4	~	22	•••	
39.1%	11.0%	7.3%	5.3%	4.5%	
75	22	~	00	20	
7.0%	6.4%	6.0%	6.0%	5.8%	
	515	6		53	
10.9%	9.7%	6.5%	5.7%	4.8%	
< □		→ < ≥ >	∢ ≣ ≯	æ	9 Q (P

why is this useful?

◆□▶ ◆□▶ ◆□▶ ◆□▶ ●□ ● ●



Figure 7: Correlation matrix of the model's predictions on the pretraining test set. $\langle \Box \rangle$ $\langle \Box \rangle$



Figure 6: Hierarchical clustering of the DeepMoji model's predictions across categories on the test set. The dendrogram shows how the model learns to group emojis into overall categories and subcategories based on emotional content. The y-axis is the distance on the correlation matrix of the model's predictions measured using average linkage.

イロト イヨト イヨト イヨト

크



Figure 6: Hierarchical clustering of the DeepMoji model's predictions across categories on the test set. The dendrogram shows how the model learns to group emojis into overall categories and subcategories based on emotional content. The y-axis is the distance on the correlation matrix of the model's predictions measured using average linkage.

Predicting a large set can be indicative of more coarse-grained trends.

(main idea. extra details in paper.)

- Train RNN (LSTM) to predict emojis based on a tweet.
 - Result: encoder that takes a tweet and returns a **vector** which is useful for predicting emojis.
- Take (smaller) sentiment dataset.
 - Encode sentences to vectors using above encoder.

▲ロト ▲団ト ▲ヨト ▲ヨト 三ヨー わらぐ

• Train to predict sentiment from vectors.

Option 2: Write and algorithm and hope it works

- Example: assignment 3.
 - Represent words by their contexts
 - Define the co-occurrence metric (PMI, word2vec)
 - Define similarity measure (cosine)
 - Use this to get a useful result lists of similar words
- Can be very effective
- But no "learning" involved.
- What to do when this doesn't work?

Option 3: Obtain Cheap / Easy Annotations

- Make easy annotation tasks for humans
 - Pose annotation as natural questions that are easy to answer.
- But how to come up with the right questions?

Option 3: Obtain Cheap / Easy Annotations

- Measure human behavior
 - Eye-tracking when reading
 - Mouse-movement when reading
 - Keyboard clicks when writing
 - ▶ ...
- How can these be leveraged to obtain useful data for learning?

Option 5: Latent-variable generative modeling

- > Define a "generative story" of how the data was generated
 - This story doesn't have to be very convincing or realistic
- The story can include "latent variables", stuff that you would like to see but you don't
 - For example: HMM POS-tagging, where we treat the tags as latent.
- Search for an assignment of latent variables such that the data has high probability under the model.
 - Usually, this search is hard.
 - Approximate!
 - ► EM
 - MCMC (Gibbs sampling)

Unsupervised Learning Example: HMM

Example: HMM

- We want to train a POS-tagger, but don't have labeled data.
- We do have a dictionary, associating some words with their possible POS tags, and also a lot of text.
- We will use the dictionary and the text to train a bigram HMM model.

Unsupervised Learning

Example: HMM

The Bigram-HMM generative story:

To generate a tagged sentence $(w, t) = (w_1, \dots, w_n, t_1, \dots, t_n)$:

- Start with tag $t_0 = START$.
- ► For *i* in 1,...,*n*:
 - Draw a random tag t_i from the transition distribution $P(t_i|t_{i-1})$
 - Draw a random word w_i from the tag distribution $P(w_i|t_i)$

Recall the supervised case

- We observe both the words and the tags.
- ► We estimate q = P(t_i|t_{i-1}) and e = P(w_i|t_i) based on our observations.
- Done

Notation – Discrete Distributions

we say that $X \sim Discrete(\theta, k)$ iff:

- X can get one of k values
- θ is a vector with k entries

•
$$\theta_i \ge 0$$

•
$$\sum_i \theta_i = 1$$

$$\blacktriangleright P(X=i) = \theta_i$$

Example

 $p(t_j|t_{j-1})$ is a discrete distribution.

 $t_j \sim Discrete(\theta, |T|)$

Where:

- ► |T| is the size of the tagset
- We can get a uniform distribution if we set:

 $\bullet \ \theta_i = 1/|T|$

We can also estimate θ from data using MLE:

$$\bullet \ \theta_{t_j} = \frac{count(t_{j-1}, t_j)}{count(t_{j-1})}$$

・ロ> < 部> < 言> < 言> こ
 ・11/1

Example HMM:

The UNsupervised case

- We don't get to see the tags. They are *latent*.
- But, for a given tag assignment, we can:
 - Estimate parameters
 - Calculate corpus probability
- Search for tag assignments such that if we estimate parameters from them, and then use the parameters to calculate the corpus probability, we will get high probability.
- This search looks hard!
- And it is.
- Two possible approximations:
 - EM algorithm
 - Gibbs sampling

- $w = w_1, \ldots, w_n$
- $t = t_1, \ldots, t_n$
 - We are interested in the tag assignment that will maximize P(w,t)
 - For a fixed w, $\arg \max_t P(w, t) = \arg \max_t P(t|w)$

- $w = w_1, \ldots, w_n$
- $t = t_1, \ldots, t_n$
 - We are interested in the tag assignment that will maximize P(w,t)
 - For a fixed w, $\arg \max_t P(w, t) = \arg \max_t P(t|w)$
 - ► If we could sample from P(t|w), we will, with high probability, get t such that P(t|w) is high.

- $w = w_1, \ldots, w_n$
- $t = t_1, \ldots, t_n$
 - We are interested in the tag assignment that will maximize P(w,t)
 - For a fixed w, $\arg \max_t P(w, t) = \arg \max_t P(t|w)$
 - ► If we could sample from P(t|w), we will, with high probability, get t such that P(t|w) is high.
 - Ok... but how do we sample from P(t|w)?

- $w = w_1, \ldots, w_n$
- $t = t_1, \ldots, t_n$
 - We are interested in the tag assignment that will maximize P(w,t)
 - For a fixed w, $\arg \max_t P(w, t) = \arg \max_t P(t|w)$
 - ► If we could sample from P(t|w), we will, with high probability, get t such that P(t|w) is high.
 - Ok... but how do we sample from P(t|w)?
 - Gibbs sampling is a "magical" way of doing that
 - To uncover the magic, see Graphical Models class

Main idea

- In order to sample $P(t|w) = P(t_1, t_2, \dots, t_n|w)$:
- Start with a random assignment of t_1, \ldots, t_n . Then:
 - sample t_1 based on t_2, \ldots, t_n, w
 - $\blacktriangleright P(t_1|t_2,t_3,\ldots,t_n,w)$
 - sample t_2 based on t_1, t_3, \ldots, t_n, w
 - ▶ ...
 - sample t_k based on $t_1, \ldots, t_{k-1}, t_{k+1}, \ldots, t_n, w$
 - ...and so on
- After many iterations, we will get samples from P(t|w)

Calculating $P(t_k|t_1,\ldots,t_{k-1},t_{k+1},\ldots,t_n,w)$

- Notation: $t^{-k} = t_1, ..., t_{k-1}, t_{k+1}, ..., t_n$.
- ► We can estimate q and e as previously, based on w and the assignments to t^{-k}.
- Now we get:

$$P(t_k|t^{-\mathbf{k}}) \propto q(t_k|t_{k-1})e(w_k|t_k)q(t_{k+1}|t_k)$$

- (why? and what does \propto means?)
- Calculate this for every possible value of t_k.
- Normalize

Draws from distributions

$X \sim Discrete(\theta, k)$

$$p = Math.random()$$

sum = 0.0

```
for i in 0...k-1 {
   sum += theta[i];
   if(sum >= p) return i
}
```



The Gibbs sampling algorithm

Sampling from P(t|w) for $t = t_1, \ldots, t_n$

Initialize *t* with random values Calculate parameters (collect counts) based on *t*,*w*. for many iterations **do** for $i \in 1, ..., n$ **do** "forget" value of t_i (decrease counts) Calculate $P(t_i|t^{-i})$ based on modified counts Sample new value for t_i from $P(t_i|t^{-i})$

Putting it all together

Training HMM from text and dictionary using Gibbs sampling

For each word, assign a random tag from the set allowed by the dictionary

Calculate q, e based on this tag assignment

for many iterations do

for every sentence do

for $i \in 1, \ldots, length$ do

"forget" value of t_i (decrease counts) Calculate $P(t_i|t^{-i})$ based on modified counts (Set prob of tags not in dictionary to 0. Normalize.) Sample new value for t_i from $P(t_i|t^{-i})$

Putting it all together

Training HMM from text and dictionary using Gibbs sampling

For each word, assign a random tag from the set allowed by the dictionary

Calculate q, e based on this tag assignment

for many iterations do

for every sentence do

for $i \in 1, \ldots, length$ do

"forget" value of t_i (decrease counts) Calculate $P(t_i|t^{-i})$ based on modified counts (Set prob of tags not in dictionary to 0. Normalize.) Sample new value for t_i from $P(t_i|t^{-i})$

Calculate final q and e based on the final state (can also average several states)

HMM - discussion

Why do you expect this to work?

Why do we need the tag dictionary?

Topic Modeling / LDA

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - のへで

The problem with information



www.betaversion.org/~stefano/linotype/news/26/

As more information becomes available, it becomes more difficult to access what we are looking for.

We need new tools to help us organize, search, and understand these vast amounts of information.
Topic modeling



Topic modeling provides methods for automatically organizing, understanding, searching, and summarizing large electronic archives.

- 1 Uncover the hidden topical patterns that pervade the collection.
- 2 Annotate the documents according to those topics.
- **3** Use the annotations to organize, summarize, and search the texts.

Discover topics from a corpus

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

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

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

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

Model the evolution of topics over time



D. Blei

Topic Models

Model connections between topics



Annotate images



SKY WATER TREE MOUNTAIN PEOPLE



SCOTLAND WATER FLOWER HILLS TREE



SKY WATER BUILDING PEOPLE WATER







FISH WATER OCEANPEOPLE MARKET PATTERNTREE CORALTEXTILE DISPLAY

BIRDS NEST TREE BRANCH LEAVES

Topic Models

ls 🗖

Discover influential articles



Derek E. Wildman et al., Implications of Natural Selection in Shaping 99.4% Nonsynonymous DNA Identity between Humans and Chimpanzees: Enlarging Genus Homo, PNAS (2003) 1178 citations!

イロト イヨト イヨト イヨト

æ

[3 citations]

Organize and browse large corpora



Stanley Kubrick



(film, series, show) (theory, work, human) (son, year, death) (black, white, people) (god, call, give) (math, energy, light)

Stanley Kuhrick (pi) 26, (128– Yuri) 7, (199) was an American IIIIn "discosen, writter products of an photographic who have in Edgind during most of the antipolicity of the standard during most of the structured and excessions, and the reclauses, has leave method of working, the variety of games ha worked in the structured applections, and the reclauses about 18 to confine at the Hollywood system, minimizing almost complexe atrids constraints, bus with the rate has a working and the constraints, bus with the rate workers in the structure of the subserver.

Kubrick's films are characterized by a formal visual spie and meticulous attendint to detail—his later films often have elements of surrealism and expressionism that eachews structured linear narratives. His films are repeatedly described as a low and methodical, and are often proceived as a reflection of his obsessive and perfectionist nazare.¹¹ A recurring theme in his films is much indumerito to man. While other weed as





words	related documents	related topics
theory	Meme	{work, book, publish}
work	Intelligent design	{law, state, case}
human	Immanuel Kant	(son, year, death)
idea	Philosophy of mathematics	{woman, child, man}
term	History of science	{god, call, give}
study	Free will	(black, white, people)
view	Truth	(film, series, show)
cience	Psychoanalysis	{war, force, army}
oncept	Charles Peirce	(language, word, form)
form	Existentialism	(@card@, make, design)
world	Deconstruction	(church, century, christian)
argue	Social sciences	{rate, high, increase}
social	Idealism	(company, market, business

イロト イロト イヨト イヨト 三日

Latent Dirichlet Allocation

Latent Dirichlet allocation (LDA)

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

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

SCIENCE • VOL. 272 • 24 MAY 1996

"are not all that far apart," especially in comparison to the 75.000 genes in the human genome, notes Siv Andersson of 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 nore 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



mate of the minimum modern and ancient genomes.

イロト イポト イヨト イヨト

Simple intuition: Documents exhibit multiple topics.

Generative model for LDA



э

- Each topic is a distribution over words
- Each document is a mixture of corpus-wide topics
- Each word is drawn from one of those topics

The posterior distribution



- In reality, we only observe the documents
- The other structure are hidden variables

The posterior distribution



- Our goal is to **infer** the hidden variables
- I.e., compute their distribution conditioned on the documents
 p(topics, proportions, assignments | documents)

《口》 《卽》 《臣》 《臣

Graphical models (Aside)



- Nodes are random variables
- Edges denote possible dependence
- Observed variables are shaded
- Plates denote replicated structure

Graphical models (Aside)



- Structure of the graph defines the pattern of conditional dependence between the ensemble of random variables
- E.g., this graph corresponds to

$$p(y, x_1, ..., x_N) = p(y) \prod_{n=1}^N p(x_n | y)$$



<ロト <四ト <注入 <注下 <注下 <

- Nodes are random variables; edges indicate dependence.
- Shaded nodes are observed; unshaded nodes are hidden.
- Plates indicate replicated variables.

We have *K* topics, and a vocabulary *V* of |V| words. Each topic β^k is a distribution over words.

A document *d* is created by

- Sample length *n_d* from a Poisson distribution
 - (alternatively, assume n_d is given)
- Sample topic proportions θ^d from a Dirichlet distribution with parameter α.
- For each position $i \in 1, \ldots, n$:
 - Sample topic z_i from θ^d
 - Sample word w_i from the distribution β^{z_i}

We have *K* topics, and a vocabulary *V* of |V| words. Each topic β^k is a distribution over words.

A document *d* is created by

- Sample length *n_d* from a Poisson distribution
 - (alternatively, assume n_d is given)
- Sample topic proportions θ^d from a Dirichlet distribution with parameter α .
- For each position $i \in 1, \ldots, n$:
 - Sample topic z_i from θ^d
 - Sample word w_i from the distribution β^{z_i}

Assumptions

- We do not care about the word-order ("bag of words")
- Each word is independent of the other words given its topic

We have *K* topics, and a vocabulary *V* of |V| words. Each topic β^k is a distribution over words.

A document *d* is created by

- Sample length *n_d* from a Poisson distribution
 - (alternatively, assume n_d is given)
- Sample topic proportions θ^d from a Dirichlet distribution with parameter α.
- For each position $i \in 1, \ldots, n$:
 - Sample topic z_i from θ^d
 - Sample word w_i from the distribution β^{z_i}

Assumptions

- We do not care about the word-order ("bag of words")
- Each word is independent of the other words given its topic

- The Dirichlet distribution is a "distribution over distributions"
- When you sample $\theta \sim \text{DIRICHLET}(\alpha, K)$:
 - θ is a *K*-dim vector

•
$$\theta_i \ge 0$$

•
$$\sum_i \theta_i = 1$$

- The Dirichlet distribution is a "distribution over distributions"
- When you sample $\theta \sim \text{DIRICHLET}(\alpha, K)$:
 - θ is a *K*-dim vector

►
$$\theta_i \ge 0$$

•
$$\sum_i \theta_i = 1$$

The probability of seeing a particular vector θ is:

$$P_{\mathsf{Dirichlet}(lpha,K)}(heta) = rac{\prod_{i=1}^{K} heta_{i}^{lpha_{i}-1}}{B(lpha)}$$

$$B(\alpha) = \frac{\prod_{i=1}^{K} \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^{K} \alpha_i)}$$

- Γ is the gamma function, generalization of factorial.

- Generally, α is a *k*-dim vector, but we will assume "symmetric" dirichlet, in which α is a single scalar (and $\alpha_i = \alpha$ for all $i \in \{1, \dots, K\}$)

The Dirichlet distribution is a "distribution over distributions"

$$P_p(\theta|\alpha) = \frac{\prod_{i=1}^{K} \theta_i^{\alpha_i - 1}}{B(\alpha)}$$

• α controls the shape, mean and sparsity of θ



- + ロト + 御 ト + ヨト + ヨト - ヨ - のへで



- * ロ > * 御 > * 注 > * 注 > … 注 … の 9



(日) (日) (日) (日) (日)

æ



- + ロト + 御 ト + ヨト + ヨト - ヨ - のへで



- * ロ * * 母 * * き * モ * ヨ * うくで

 $\alpha = 0.01$



- イロト (個) (注) (注) モー のの(



- ・ロト ・回ト ・ヨト ・ヨー つく(

For draws θ from a symmetric dirichlet distribution:

- $\alpha = 1$ All θ are equally likely
- $\alpha > 1$ Uniform θ are more likely
- $\alpha < 1$ Spikey θ are more likely

We have *K* topics, and a vocabulary *V* of |V| words. Each topic β^k is a distribution over words. $\beta^k \sim Dirichlet(\eta, |V|)$

A document *d* is created by

- Sample length *n_d* from a Poisson distribution
 - (alternatively, assume n_d is given)
- Sample topic proportions θ^d from a Dirichlet distribution with parameter α .
- For each position $i \in 1, \ldots, n$:
 - Sample topic z_i from θ^d
 - Sample word w_i from the distribution β^{z_i}

We have *K* topics, and a vocabulary *V* of |V| words. Each topic β^k is a distribution over words. $\beta^k \sim Dirichlet(\eta, |V|)$

A document *d* is created by

- Sample length *n_d* from a Poisson distribution
 - (alternatively, assume n_d is given)
- Sample topic proportions θ^d from a Dirichlet distribution with parameter α .
- For each position $i \in 1, \ldots, n$:
 - Sample topic z_i from θ^d
 - Sample word w_i from the distribution β^{z_i}

$\boldsymbol{\alpha}$ controls how many topics we expect to see in our documents

Latent Dirichlet allocation (LDA)



- Our goal is to infer the hidden variables
- I.e., compute their distribution conditioned on the documents

p(topics, proportions, assignments|documents)

・ロト ・ 日ト ・ モト ・ モト



<ロト <四ト <注入 <注下 <注下 <

- Nodes are random variables; edges indicate dependence.
- Shaded nodes are observed; unshaded nodes are hidden.
- Plates indicate replicated variables.



・ロト ・日 ・ ・ ヨ ・ ・ ヨ ・ うへぐ



- This joint defines a posterior, $p(\theta, z, \beta | w)$.
- From a collection of documents, infer
 - Per-word topic assignment z_{d,n}
 - Per-document topic proportions θ_d
 - Per-corpus topic distributions β_k
- Then use posterior expectations to perform the task at hand: information retrieval, document similarity, exploration, and others.

Example inference

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 malically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compute Innowing genomes, concluded that today's organisms can be sastained with host 255 genes. and that the centiese 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

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

SCIENCE • VOL. 272 • 24 MAY 1996

"are not all that far apart," especially in comparison to the 75.002 eners in the human genome, notes Six Andersson of 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 Aready Mushegin, a comparational mo-

lecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an







Topic Models
Example inference

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

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

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

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

Example inference (II)

Chaotic Beetles

Charles Godfray and Michael Hassell

Ecologists have known since the pioneering work of May in the mid-1970s (1) that the population dynamics of animals and plants can be exceedingly complex. This complexity arises from two sources: The tangled web of interactions that constitute any natural community provide a myriad of different pathways for species to interact, both directly and indirectly. And even in isolated populations the nonlinear feedback processes present in all natural populations can result in complex dynamic behavior. Natural populations can show persistent oscillatory dynamics and chaos, the latter characterized by extreme sensitivity to initial conditions. If such chaotic dynamics were common in nature, then this would have important ramifications for the management and conservation of natural resources. On page 389 of this issue, Costantino et al. (2) provide the most

The authors are in the Department of Biology, Imperial College at Silvood Park, Ascot, Berks, SL5 7PZ UK. Email: m.hassell@ic.ac.uk convincing evidence to date of complex dynamics and chaos in a biological population—of the flour beetle, *Tribolium castaneum* (see figure).

It has proven extremely difficult to demonstrate complex dynamics in populations in the field. By its very nature, a chaotically fluctuating population will superficially resemble a stable or cyclic population buffeted by the normal random perturbations experienced by all species. Given a long enough time series, diagnostic tools from nonlinear mathematics can be used to identify the telltale signatures of chaos. In phase space, chaotic trajectories come to lie on "strange attractors," curious geometric objects with fractal structure and hence noninteger dimension. As they move over the surface of the attractor, sets of adjacent trajectories are pulled apart, then stretched and folded, so that it becomes impossible to predict exact population densities into the future. The strength of the mixing that gives rise to the extreme sensitivity to initial conditions can be measured mathematically estimating the Liapunov expo-



Cannibalism and chaos. The flour beetle, *Tribolium castaneum*, exhibits chaotic population dynamics when the amount of cannibalism is altered in a mathematical model

Topic Models

somating the traditions exploment of the stand normaltic dynamics and normaltic dynamics and normaltic observations of the standard norway attempts to estimate attractor dimension and Liapunov exponents from time series data, and some candidate chaotic population have been identified (some insects, nodents, and most convincingly, human childhood diseases), but the statistical difficculties preclude any broad generalization (3).

An alternative approach is to parameterize population models with data from natural populations and then compare their predictions with the dynamics in the field. This technique has been gaining populatity in recent years, helped by statistical advances in parameter estimation. Good ex-

SCIENCE . VOL. 275 . 17 JANUARY 1997

D Blei

Example inference (II)

problem problems mathematical number new mathematics university two first numbers work time mathematicians chaos chaotic

model rate constant distribution time number size values value average rates data density measured models

selection male males females sex species female evolution populations population sexual behavior evolutionary genetic reproductive

species forest ecology fish ecological conservation diversity population natural ecosystems populations endangered tropical forests ecosystem

Why does LDA work?

- LDA trades off two goals.
 - For each document, allocate its words to as few topics as possible.
 For each topic, assign high probability to as few terms as possible.
- These goals are at odds.
 - Putting a document in a single topic makes #2 hard: All of its words must have probability under that topic.
 - Putting very few words in each topic makes #1 hard: To cover a document's words, it must assign many topics to it.

• Trading off these goals finds groups of tightly co-occurring words.

What do we get out of LDA?

- Topic assignments z
- Topic proportions (how strong is topic k in document j?)
- Topics distributions (how strong is word i in topic k?)
 - Also: which topics are related to word i?

What do we get out of LDA?

- Topic assignments z
- Topic proportions (how strong is topic k in document j?)
- Topics distributions (how strong is word i in topic k?)
 - Also: which topics are related to word i?

So?

- Which topics are in our corpus?
- Find similar docs (by comparing "topic vectors" of docs)
- Find related words (by comparing "topic vectors" of words)
- Query expansion: find documents related to words X,Y,Z, even if all or some of these words did not appear in the document

▶ ...

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services." Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

Topic Model

- P(tlk) for all t and k, is a term by topic matrix (gives which terms make up a topic)
- P(kldoc) for all k and doc, is a topic by document matrix (gives which topics are in a document)

<ロト <四ト <注入 <注下 <注下 <



EXAMPLE

Analysis of TASA Corpus

- Given a text collection written by first grade to college students
- Data has following characteristics
 - 26,000+ word types (stop words removed)

▲ロト ▲団ト ▲ヨト ▲ヨト 三ヨー わらぐ

- 37,000+ documents
- 6,000,000+ word tokens
- Find topics in the data

Topics in the Educational Corpus (TASA)

- 37K docs, 26K words
- 1700 topics, e.g.:

Polysemy



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで

Three Documents with the word "play"

(numbers & colors → topic assignments)

A Play⁰⁸² is written⁰⁸² to be performed⁰⁸² on a stage⁰⁸² before a live⁰⁹³ audience⁰⁸² or before motion²⁷⁰ picture⁰⁰⁴ or television⁰⁰⁴ cameras⁰⁰⁴ (for later⁰⁵⁴ viewing⁰⁰⁴ by large²⁰² audiences⁰⁸²). A Play⁰⁸² is written⁰⁸² because playwrights⁰⁸² have something

He was listening⁰⁷⁷ to music⁰⁷⁷ coming⁰⁰⁹ from a passing⁰⁴³ riverboat. The music⁰⁷⁷ had already captured⁰⁰⁶ his heart¹⁵⁷ as well as his ear¹¹⁹. It was jazz⁰⁷⁷. Bix beiderbecke had already had music⁰⁷⁷ lessons⁰⁷⁷. He wanted²⁶⁸ to play⁰⁷⁷ the cornet. And he wanted²⁶⁸ to play⁰⁷⁷ jazz⁰⁷⁷

Jim²⁹⁶ plays¹⁶⁶ the game¹⁶⁶. Jim²⁹⁶ likes⁰⁸¹ the game¹⁶⁶ for one. The game¹⁶⁶ book²⁵⁴ helps⁰⁸¹ jim²⁹⁶. Don¹⁸⁰ comes⁰⁴⁰ into the house⁰³⁸. Don¹⁸⁰ and jim²⁹⁶ read²⁵⁴ the game¹⁶⁶ book²⁵⁴. The boys⁰²⁰ see a game¹⁶⁶ for two. The two boys⁰²⁰ play¹⁶⁶ the game¹⁶⁶.

LDA Inference

How do we fit an LDA model to the data?

Use an existing tool!

- Mallet (java)
- gensym (python)
- Many other tools available
 - (see David Blei's website)

But how are the tools implemented? And what if we want a slightly different story?

But how are the tools implemented?

And what if we want a slightly different story?

- Exact inference is intractable.
- Use an approximate algorithm.

But how are the tools implemented?

And what if we want a slightly different story?

- Exact inference is intractable.
- Use an approximate algorithm.
- Current tools use modern complex algorithms:
 - Fast
 - Scale well to huge number of topics and documents
 - Beyond the scope of this course

But how are the tools implemented?

And what if we want a slightly different story?

- Exact inference is intractable.
- Use an approximate algorithm.
- Current tools use modern complex algorithms:
 - Fast
 - Scale well to huge number of topics and documents
 - Beyond the scope of this course
- but for fitting a small to medium data, we can use Gibbs sampling.
 - (Gibbs is also our best bet for implementing modifications of LDA)

Recall:

- Inputs: α, η, Κ
- Obeserved variables: words, $W = w_{d,n}$
- Unobserved: $\theta = \theta^1, \dots, \theta^D, \beta = \beta^1, \dots, \beta^K, Z = z_{d,n}$

We need to sample from $p(Z, \theta, \beta | W, \alpha, \eta)$

In Gibbs:

Initialize random z

Then, repeatedly:

- For each $k \in \{1, \dots, K\}$, sample β^k based on Z, W, η
- Sample θ^d based on Z, W, α
- Sample $z_{d,1}$ based on $Z^{-\mathbf{d},\mathbf{1}}$, θ^d , β , W
- Sample $z_{d,2}$ based on $Z^{-\mathbf{d},\mathbf{2}}$, θ^d , β , W

- For each $k \in \{1, \dots, K\}$, sample β^k based on Z, W, η
- Sample θ^d based on Z, W, α
- Sample $z_{d,1}$ based on $Z^{-\mathbf{d},\mathbf{1}}$, θ^d , β , W
- Sample $z_{d,2}$ based on $Z^{-\mathbf{d},\mathbf{2}}$, θ^d , β , W

▶ ...

- For each $k \in \{1, \dots, K\}$, sample β^k based on Z, W, η
- Sample θ^d based on Z, W, α
- Sample $z_{d,1}$ based on $Z^{-d,1}$, θ^d , β , W
- Sample $z_{d,2}$ based on $Z^{-d,2}$, θ^d , β , W

▶ ...

These lines are easy:

$$p(z_{d,i} = k \mid Z^{-\mathbf{d},\mathbf{i}}, \theta^d, \beta, W) = \theta_k^d \cdot \beta_{W_{d,i}}^k$$

 θ_{k}^{d} probability of generating topic *k* in doc *d* $\beta_{W_{d,i}}^{k}$ probability of generating word $W_{d,i}$ from topic *k*

- For each $k \in \{1, \dots, K\}$, sample β^k based on Z, W, η
- Sample θ^d based on Z, W, α
- Sample $z_{d,1}$ based on $Z^{-\mathbf{d},\mathbf{1}}$, θ^d , β , W
- Sample $z_{d,2}$ based on $Z^{-d,2}$, θ^d , β , W

▶ ...

What does this line mean?

- For each $k \in \{1, \dots, K\}$, sample β^k based on Z, W, η
- Sample θ^d based on Z, W, α
- Sample $z_{d,1}$ based on $Z^{-\mathbf{d},\mathbf{1}}$, θ^d , β , W
- Sample $z_{d,2}$ based on $Z^{-d,2}$, θ^d , β , W

▶ ...

What does this line mean?

We need to sample θ^d from $p(\theta|Z, \alpha)$.

• Given Z, we can derive an MLE estimate of θ^d :

$$\theta_k^d = \frac{count(z_{d,i} = k)}{n_d}$$

・ロン ・回 と ・ ヨン ・ ヨン … ヨ

32/1

- For each $k \in \{1, \dots, K\}$, sample β^k based on Z, W, η
- Sample θ^d based on Z, W, α
- Sample $z_{d,1}$ based on $Z^{-\mathbf{d},\mathbf{1}}$, θ^d , β , W
- Sample $z_{d,2}$ based on $Z^{-d,2}$, θ^d , β , W
- ▶ ...

What does this line mean?

We need to sample θ^d from $p(\theta|Z, \alpha)$.

• Given Z, we can derive an MLE estimate of θ^d :

$$\theta_k^d = \frac{count(z_{d,i} = k)}{n_d}$$

But no. We need to sample. What does it mean to sample θ?

- For each $k \in \{1, \dots, K\}$, sample β^k based on Z, W, η
- Sample θ^d based on Z, W, α
- Sample $z_{d,1}$ based on $Z^{-\mathbf{d},\mathbf{1}}$, θ^d , β , W
- Sample $z_{d,2}$ based on $Z^{-d,2}$, θ^d , β , W
- ▶ ...

What does this line mean?

We need to sample θ^d from $p(\theta|Z, \alpha)$.

• Given Z, we can derive an MLE estimate of θ^d :

$$\theta_k^d = \frac{count(z_{d,i} = k)}{n_d}$$

- But no. We need to sample. What does it mean to sample θ?
- ► Under the *Bayesian* philosophy, we do not commit to a single estimate of θ. Instead, we have a distribution p(θ^d|Z, α) of possible θ^d, based on our prior belief α and the data we saw Z.

- For each $k \in \{1, \dots, K\}$, sample β^k based on Z, W, η
- Sample θ^d based on Z, W, α
- Sample $z_{d,1}$ based on $Z^{-\mathbf{d},\mathbf{1}}$, θ^d , β , W
- Sample $z_{d,2}$ based on $Z^{-\mathbf{d},\mathbf{2}}$, θ^d , β , W

▶ ...

What does this line mean? We need to sample θ^d from $p(\theta|Z, \alpha)$.

- For each $k \in \{1, \dots, K\}$, sample β^k based on Z, W, η
- Sample θ^d based on Z, W, α
- Sample $z_{d,1}$ based on $Z^{-\mathbf{d},\mathbf{1}}$, θ^d , β , W
- Sample $z_{d,2}$ based on $Z^{-\mathbf{d},\mathbf{2}}$, θ^d , β , W

▶ ...

What does this line mean? We need to sample θ^d from $p(\theta|Z, \alpha)$.

Because $\theta^d \sim \text{DIRICHLET}(\alpha, K)$, and because dirichlet is *conjugate* to multinomial, we have:

$$\theta^d | Z, \alpha \sim \mathsf{Dirichlet}(\alpha + c^d)$$

where c^d is s *K*-dim vector based on counts from *Z*, with c_k^d is the number of items in document *d* with topic *k*.

- For each $k \in \{1, \dots, K\}$, sample β^k based on Z, W, η
- Sample θ^d based on Z, W, α
- Sample $z_{d,1}$ based on $Z^{-d,1}$, θ^d , β , W
- Sample $z_{d,2}$ based on $Z^{-\mathbf{d},\mathbf{2}}$, θ^d , β , W

▶ ...

We need to sample θ^d from $p(\theta|Z, \alpha)$.

$$\theta^d | Z, \alpha \sim \mathsf{Dirichlet}(\alpha + c^d)$$

- There are algorithms for sampling from Dirichlet, but we don't need to actualy used them.
- Instead, we will use the the collapsed Gibbs sampler.

Recall:

- Inputs: α, η, Κ
- Obeserved variables: words, $W = w_{d,n}$
- Unobserved: $\theta = \theta^1, \dots, \theta^D, \beta = \beta^1, \dots, \beta^K, Z = z_{d,n}$

We need to sample from

 $p(Z, \theta, \beta | W, \alpha, \eta)$

But actually, we are ok with just Z. Can we get rid of θ , β ?

• If θ , β were discrete, we could marginalize over them.

Recall:

- Inputs: α, η, Κ
- Obeserved variables: words, $W = w_{d,n}$
- Unobserved: $\theta = \theta^1, \dots, \theta^D, \beta = \beta^1, \dots, \beta^K, Z = z_{d,n}$

We need to sample from

 $p(Z, \theta, \beta | W, \alpha, \eta)$

But actually, we are ok with just Z. Can we get rid of θ , β ?

- If θ , β were discrete, we could marginalize over them.
- But they are continuous, so instead we need to integrate

$$p(Z|W,\alpha,\eta) = \int \int p(Z,\theta,\beta|W,\alpha,\eta) \mathrm{d}\theta \mathrm{d}\beta$$

Recall:

- Inputs: α, η, Κ
- Obeserved variables: words, $W = w_{d,n}$
- Unobserved: $\theta = \theta^1, \dots, \theta^D, \beta = \beta^1, \dots, \beta^K, Z = z_{d,n}$

We need to sample from

 $p(Z, \theta, \beta | W, \alpha, \eta)$

But actually, we are ok with just Z. Can we get rid of θ , β ?

- If θ , β were discrete, we could marginalize over them.
- But they are continuous, so instead we need to integrate

$$p(Z|W,\alpha,\eta) = \int \int p(Z,\theta,\beta|W,\alpha,\eta) \mathrm{d}\theta \mathrm{d}\beta$$

$$p(z_{d,i} = k | W, Z^{-\mathbf{d}, \mathbf{i}}, \alpha, \eta) = \int \int p(z_{d,i} = k, \theta, \beta | W, Z^{-\mathbf{d}, \mathbf{i}}, \alpha, \eta) \mathrm{d}\theta \mathrm{d}\beta$$

$$= \int p(z_{d,i} = k|\theta) p(\theta|\alpha) \mathrm{d}\theta \int p(w_{d,i} = \nu|W^{-\mathbf{d},\mathbf{i}}, z_{d,i} = k, Z^{-\mathbf{d},\mathbf{i}}, \beta) p(\beta|\eta) \mathrm{d}\eta$$

You don't really need to know how to integrate! Just remember that for Dirichlet:

$$\int p(x|data,\theta)p(\theta|\alpha)d\theta = \frac{c_x + \alpha}{|data| + K\alpha}$$

Where c_x is the count of event *x* in the data, and $|data| = \sum_{x'} c'_x$ is the number of samples in the data.

Just remember that for Dirichlet:

$$\int p(x|data,\theta)p(\theta|\alpha)d\theta = \frac{c_x + \alpha}{|data| + K\alpha}$$

Where c_x is the count of event *x* in the data, and $|data| = \sum_{x'} c'_x$ is the number of samples in the data. Use this rule twice (once for each \int), and get:

se this fulle twice (once for each f), and get.

$$p(z_{d,i} = k | Z^{-\mathbf{d}, \mathbf{i}}, \alpha, \eta, w_i) = \frac{c_k^d + \alpha}{\sum_{k'} c_{k'}^d + K\alpha} \frac{v_{w_i}^k + \eta}{\sum_{i'} v_{w_{i'}}^k + |V|\eta}$$

 c_k^d number of words in doc d with topic k in $Z^{-\mathbf{d},\mathbf{i}}$

- $v_{w_i}^k$ number of times word w_i is assigned to topic k in $Z^{-\mathbf{d},\mathbf{i}}$
 - K number of topics
- V vocabulary size

Rule of Thumb In MLE land:

$$p(x_n = k | x_1, x_2, \dots, x_{n-1}) = \frac{count(k)}{n-1}$$

In Dirichlet-prior α land:

$$p(x_n = k | x_1, x_2, \dots, x_{n-1}, \alpha) = \frac{count(k) + \alpha}{n - 1 + K\alpha}$$

Derivation in MacKay and Peto (1994)
Collapsed Gibbs Sampler

- Initialize random topics Z
- For many iterations, for each document d, for each word i:
 - ▶ forget z_{d,i} getting Z^{-d,i}
 - sample new assignment for *z*_{d,i} based on equation below.

$$p(z_{d,i} = k | Z^{-\mathbf{d}, \mathbf{i}}, \alpha, \eta, w_i) = \frac{c_k^d + \alpha}{\sum_{k'} c_{k'}^d + K\alpha} \frac{v_{w_i}^k + \eta}{\sum_{i'} v_{w_{i'}}^k + |V|\eta}$$

 c_k^d number of words in doc d with topic k in $Z^{-d,i}$

- $v_{w_i}^k$ number of times word w_i is assigned to topic k in $Z^{-\mathbf{d},\mathbf{i}}$
 - K number of topics
- V vocabulary size

We have topics, are they good?

Internal Evaluation

If we want to compare two different LDA models on the same data:

- Compare the Probability that is assigned to the data by each model.
- Higher probability \rightarrow better model

Internal Evaluation

If we want to compare two different LDA models on the same data:

- Compare the Probability that is assigned to the data by each model.
- Higher probability \rightarrow better model
- But this does not tell us much about how useful the topics are...

Internal Evaluation

If we want to compare two different LDA models on the same data:

- Compare the Probability that is assigned to the data by each model.
- Higher probability \rightarrow better model
- But this does not tell us much about how useful the topics are...

External (task-based) Evaluation

- Use the LDA topics as features in another task
- Measure the accuracy of the other task

Internal Evaluation

If we want to compare two different LDA models on the same data:

- Compare the Probability that is assigned to the data by each model.
- Higher probability \rightarrow better model
- But this does not tell us much about how useful the topics are...

External (task-based) Evaluation

- Use the LDA topics as features in another task
- Measure the accuracy of the other task
- Good! But we need to have a task that we can automatically measure.

Human Evaluation

If we just want to know if our topics are "good" we can ask people.

But what is a good topic?

Human Evaluation

If we just want to know if our topics are "good" we can ask people.

- But what is a good topic?
- "Intruder Detection"
 - Take top words from a topic.
 - Insert a random word which is high in another topic.
 - Can a human identify the random word?
 - Yes \rightarrow good topic

Other Applications of LDA

Change the definition of Document

Selectional Preferences Take parsed corpus:

Documents each Verb is a document Words each subject of a verb is a "word" in the document Topics each topic is one "kind" of arguments

Topic t	Arg1	Relations which assign	Arg2
	-	highest probability to t	-
18	The residue - The mixture - The reaction	was treated with, is	EtOAc - CH2Cl2 - H2O - CH.sub.2Cl.sub.2
	mixture - The solution - the mixture - the re-	treated with, was	- H.sub.2O - water - MeOH - NaHCO3 -
	action mixture - the residue - The reaction -	poured into, was	Et2O - NHCl - CHCl.sub.3 - NHCl - drop-
	the solution - The filtrate - the reaction - The	extracted with, was	wise - CH2Cl.sub.2 - Celite - Et.sub.2O -
	product - The crude product - The pellet -	purified by, was di-	Cl.sub.2 - NaOH - AcOEt - CH2C12 - the
	The organic layer - Thereto - This solution	luted with, was filtered	mixture - saturated NaHCO3 - SiO2 - H2O
	- The resulting solution - Next - The organic	through, is disolved in,	- N hydrochloric acid - NHCl - preparative
	phase - The resulting mixture - C.)	is washed with	HPLC - to0 C
151	the Court - The Court - the Supreme Court	will hear, ruled in, de-	the case - the appeal - arguments - a case -
	- The Supreme Court - this Court - Court	cides, upholds, struck	evidence - this case - the decision - the law
	- The US Supreme Court - the court - This	down, overturned,	- testimony - the State - an interview - an
	Court - the US Supreme Court - The court	sided with, affirms	appeal - cases - the Court - that decision -
	- Supreme Court - Judge - the Court of Ap-		Congress - a decision - the complaint - oral
	peals - A federal judge		arguments - a law - the statute
211	President Bush - Bush - The President -	hailed, vetoed, pro-	the bill - a bill - the decision - the war - the
	Clinton - the President - President Clinton	moted, will deliver,	idea - the plan - the move - the legislation -
	- President George W. Bush - Mr. Bush -	favors, denounced,	legislation - the measure - the proposal - the
	The Governor - the Governor - Romney -	defended	deal - this bill - a measure - the program -
	McCain - The White House - President -		the law - the resolution - efforts - the agree-
	Schwarzenegger - Obama		ment - gay marriage - the report - abortion
224	Google - Software - the CPU - Clicking -	will display, to store, to	data - files - the data - the file - the URL -
	Excel - the user - Firefox - System - The	load, processes, cannot	information - the files - images - a URL - the
	CPU - Internet Explorer - the ability - Pro-	find, invokes, to search	information - the IP address - the user - text
	gram - users - Option - SQL Server - Code	for, to delete	- the code - a file - the page - IP addresses -
1	- the OS - the BIOS		PDF files - messages - pages - an IP address

Table 1: Example argument lists from the inferred topics. For each topic number t we list the most

[6/11]



Model is slightly different - topic generates two groups of things.

(how would you change the Gibbs sampler?)

Change the definition of Document

Beyond NLP

Dataset of users who watched movies

Documents each user is a document Words each movie is a word Topics each topic is a "taste" or "genre"

- High topic-word prob: movie belong to genre
- High topic-doc prob: user likes genre

Can recommend new movies to users



- LDA can be **embedded in more complicated models**, embodying further intuitions about the structure of the texts.
- E.g., it can be used in models that account for syntax, authorship, word sense, dynamics, correlation, hierarchies, and other structure.

Summary

Unsupervised Learning

- Define generative story
- Include hidden ("latent") variables
- Find probable assignments to latent variables
- Can use Gibbs sampling

Unsupervised Learning

- Define generative story
- Include hidden ("latent") variables
- Find probable assignments to latent variables
- Can use Gibbs sampling

Topic Modeling / LDA

- A very powerful and useful model. Use it
- Generative story for LDA
- Dirichlet distributions \rightarrow can encourage sparsity
- Examples of LDA usage
- Gibbs sampler for LDA (briefly)
 - relevant for every model with dirichlet
- Evaluation: quantify human judgement ("intruder detection")
- Creative definition of documents

Next Time



Next Time



WHAT DO WE WANT? Natural language processing! WHEN DO WE WANT IT? Sorry, when do we want what?