# Modeling *Science*

David M. Blei

Department of Computer Science
Princeton University

October 3, 2007

Joint work with John Lafferty (CMU)

# Modeling *Science*



**Poisoning by ice-cream.**

No chemist certainly would suppose that the same poison exists in all samples of ice-cream which have produced untoward symptoms in man. Mineral poisons, copper, lead, arsenic, and mercury, have all been found in ice cream. In some instances these have been used with criminal intent. In other cases their presence has been accidental. Likewise, that vanilla is sometimes the bearer, at least, of the poison, is well known to all chemists. Dr. Bartley's idea that the poisonous properties of the cream which he examined were due to putrid gelatine is certainly a rational theory. The poisonous principle might in this case arise from the decomposition of the gelatine; or with the gelatine there may be introduced into the milk a ferment, by the growth of which a poison is produced.

But in the cream which I examined, none of the above sources of the poisoning existed. There were no mineral poisons present. No gelatine of any kind had been used in making the cream. The vanilla used was shown to be not poisonous. This showing was made, not by a chemical analysis, which might not have been conclusive, but Mr. Novis and I drank of the vanilla extract which was used, and no ill results followed. Still, from this cream was isolated the same poison which I had before found in poisonous cheese (*Zeitschrift für physiologische chemie, x,*

## RNA Editing and the Evolution of Parasites

Larry Simpson and Dmitri A. Maslov

## Chaotic Beetles

Charles Godfray and Michael Hassell

SCIENCE • VOL. 275 • 17 JANUARY 1997

- Our data are *Science* from 1880-2002, courtesy of JSTOR.
- JSTOR is an on-line archive that scans the original volumes and performs optical character recognition on the scans.
- This process results in 130K documents, 76M words.

# Modeling *Science*



- Discover the hidden thematic structure with hierarchical probabilistic models called *topic models*.
- Use this structure for browsing, search, and similarity assessment.

# Discover topics from a corpus

| human | evolution | disease | computer |
|---|---|---|---|
| genome | evolutionary | host | models |
| dna | species | bacteria | information |
| genetic | organisms | diseases | data |
| genes | life | resistance | computers |
| sequence | origin | bacterial | system |
| gene | biology | new | network |
| molecular | groups | strains | systems |
| sequencing | phylogenetic | control | model |
| map | living | infectious | parallel |
| information | diversity | malaria | methods |
| genetics | group | parasite | networks |
| mapping | new | parasites | software |
| project | two | united | new |
| sequences | common | tuberculosis | simulations |

# Annotate unlabeled images
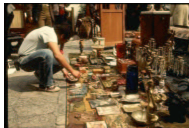


SKY WATER TREE
MOUNTAIN PEOPLE

SCOTLAND WATER
FLOWER HILLS TREE

SKY WATER BUILDING
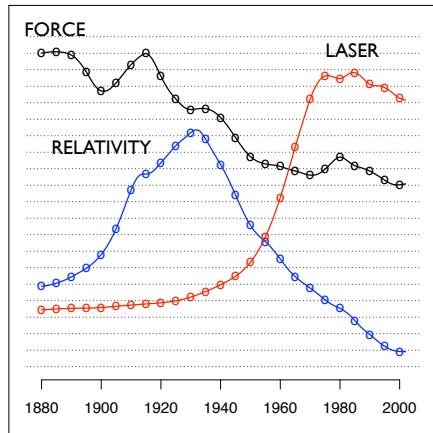PEOPLE WATER

FISH WATER OCEAN
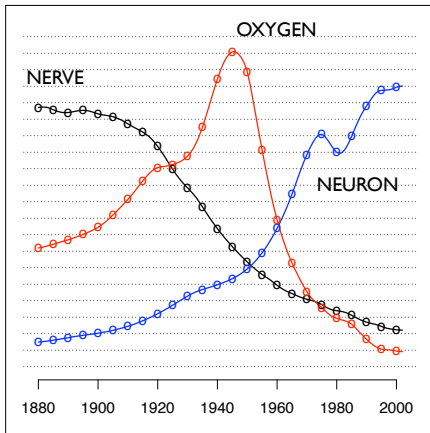TREE CORAL

PEOPLE MARKET PATTERN
TEXTILE DISPLAY

BIRDS NEST TREE
BRANCH LEAVES

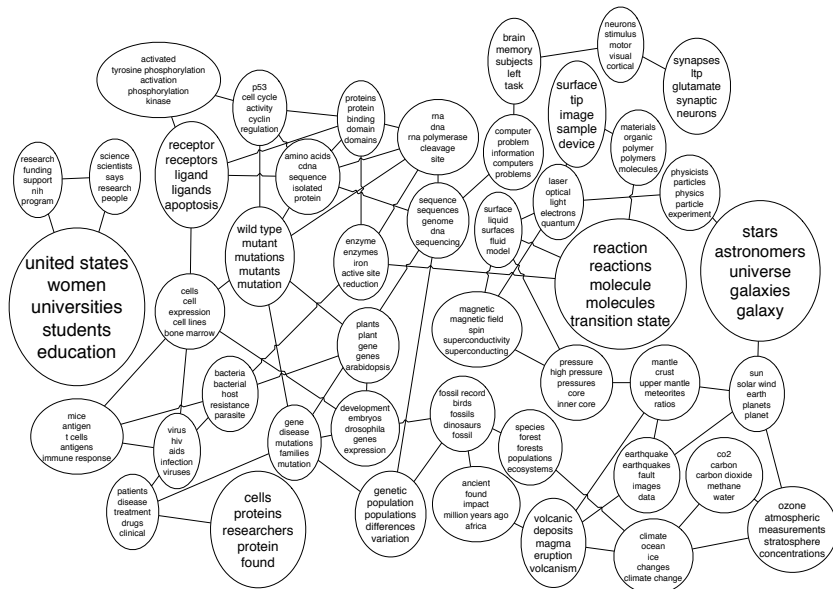# Model the evolution of topics over time

# Model connections between topics

# Outline

# Outline

# Probabilistic modeling

- Treat data as observations that arise from a generative probabilistic process that includes hidden variables
  - For documents, the hidden variables reflect the thematic structure of the collection.
- Infer the hidden structure using *posterior inference*
  - What are the topics that describe this collection?
- Situate new data into the estimated model.
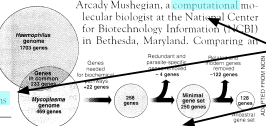  - How does this query or new document fit into the estimated topic structure?

**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

"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 more 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

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

SCIENCE • VOL. 272 • 24 MAY 1996

**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

**Simple intuition**: Documents exhibit multiple topics.

# Generative process



**Seeking Life's Bare (Genetic) Necessities**

**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

- Cast these intuitions into a generative probabilistic process
- Each document is a random mixture of corpus-wide topics
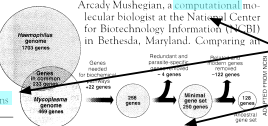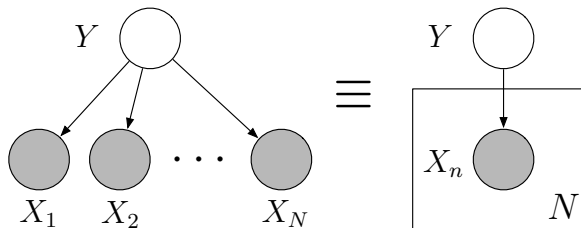- Each word is drawn from one of those topics

# Generative process



**Seeking Life's Bare (Genetic) Necessities**

- In reality, we only observe the documents
- Our goal is to infer the underlying topic structure
  - What are the topics?
  - How are the documents divided according to those topics?

# 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

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

# Latent Dirichlet allocation

# Latent Dirichlet allocation



1. Draw each topic $\beta_i \sim \mathrm{Dir}(\eta)$, for $i \in \{1, \ldots, K\}$.
2. For each document:
   1. Draw topic proportions $\theta_d \sim \mathrm{Dir}(\alpha)$.
   2. For each word:
      1. Draw $Z_{d,n} \sim \mathrm{Mult}(\theta_d)$.
      2. Draw $W_{d,n} \sim \mathrm{Mult}(\beta_{z_{d,n}})$.

# Latent Dirichlet allocation



- From a collection of documents, infer
  - Per-word topic assignment $z_{d,n}$
  - Per-document topic proportions $\theta_d$
  - Per-corpus topic distributions $\beta_k$
- Use posterior expectations to perform the task at hand, e.g., information retrieval, document similarity, etc.

# Latent Dirichlet allocation



- Computing the posterior is intractable:

$$\frac{p(\theta \mid \alpha) \prod_{n=1}^{N} p(z_n \mid \theta) p(w_n \mid z_n, \beta_{1:K})}{\int_{\theta} p(\theta \mid \alpha) \prod_{n=1}^{N} \sum_{z=1}^{K} p(z_n \mid \theta) p(w_n \mid z_n, \beta_{1:K})}$$

- Several approximation techniques have been developed.

# Latent Dirichlet allocation



- Mean field variational methods (Blei et al., 2001, 2003)
- Expectation propagation (Minka and Lafferty, 2002)
- Collapsed Gibbs sampling (Griffiths and Steyvers, 2002)
- Collapsed variational inference (Teh et al., 2006)

# Example inference



**Seeking Life's Bare (Genetic) Necessities**

**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

- **Data**: The OCR'ed collection of *Science* from 1990–2000
  - 17K documents
  - 11M words
  - 20K unique terms (stop words and rare words removed)
- **Model**: 100-topic LDA model using variational inference.

# Example inference

# Example topics

| human | evolution | disease | computer |
|---|---|---|---|
| genome | evolutionary | host | models |
| dna | species | bacteria | information |
| genetic | organisms | diseases | data |
| genes | life | resistance | computers |
| sequence | origin | bacterial | system |
| gene | biology | new | network |
| molecular | groups | strains | systems |
| sequencing | phylogenetic | control | model |
| map | living | infectious | parallel |
| information | diversity | malaria | methods |
| genetics | group | parasite | networks |
| mapping | new | parasites | software |
| project | two | united | new |
| sequences | common | tuberculosis | simulations |

## LDA discussion

- LDA is a powerful model for
    - Visualizing the hidden thematic structure in large corpora
    - Generalizing new data to fit into that structure
- LDA is a mixed membership model (Erosheva, 2004) that builds on the work of Deerwester et al. (1990) and Hofmann (1999).
    - For document collections and other grouped data, this might be more appropriate than a simple finite mixture
    - See Blei et al., 2003 for a quantitative comparison.
- *Modular*: It can be embedded in more complicated models.
- *General*: The data generating distribution can be changed.
- Variational inference is fast; allows us to analyze large data sets.
- Code to play with LDA is freely available on my web-site, http://www.cs.princeton.edu/∼blei.

# Outline

# LDA and exchangeability



- LDA assumes that documents are exchangeable.
- I.e., their joint probability is invariant to permutation.
- This is too restrictive.

# Documents are not exchangeable

"Instantaneous Photography" (1890)

"Infrared Reflectance in Leaf-Sitting Neotropical Frogs" (1977)



- Documents about the same topic are not exchangeable.
- Topics evolve over time.

# Dynamic topic model

- Divide corpus into sequential slices (e.g., by year).
- Assume each slice's documents exchangeable.
  - Drawn from an LDA model.
- Allow topic distributions evolve from slice to slice.

# Dynamic topic models

# Modeling evolving topics



- Use a logistic normal distribution to model topics evolving over time (Aitchison, 1980)
- A state-space model on the natural parameter of the topic multinomial (West and Harrison, 1997)

$$
\begin{aligned}
\beta_{t,k} \,|\, \beta_{t-1,k} &\sim \mathcal{N}(\beta_{t-1,k}, I\sigma^2) \\
p(w \,|\, \beta_{t,k}) &= \exp\left\{\beta_{t,k} - (1 + \sum_{v=1}^{V-1} \exp\{\beta_{t,k,v}\})\right\}
\end{aligned}
$$

## Posterior inference

- Our goal is to compute the posterior distribution,

$$p(\beta_{1:T,1:K}, \theta_{1:T,1:D}, \mathbf{z}_{1:T,1:D} \mid \mathbf{w}_{1:T,1:D}).$$

- Exact inference is impossible
  - Per-document mixed-membership model
  - Non-conjugacy between $p(w \mid \beta_{t,k})$ and $p(\beta_{t,k})$
- MCMC is not practical for the amount of data.
- Solution: Variational inference

# Variational inference

- Define a family of distributions $q$ on the latent variables indexed by free *variational parameters*.

- Find the member closest in $\mathrm{KL}(q||p)$ to the true posterior.

- Equivalently, maximize the Jensen's bound on the marginal likelihood of the data, within the variational family.

- See Jordan et al. (1999) and Wainwright and Jordan (2003).

- (More details at the end of the talk, if you are interested.)

## Science data



TECHVIEW: DNA S E Q U E N C I N G

Sequencing the Genome, Fast

James C. Mullikin and Amanda A. McMurray

Genome sequencing projects reveal
    the genetic makeup of an organism
    by reading off the sequence of the
DNA bases, which encodes all of the infor-
mation necessary for the life of the organ-
ism. The base sequence contains four nu-
cleotides-adenine, thymidine, guanosine,
and cytosine-which are linked together
into long double-helical chains. Over the
last two decades, automated DNA se-
quencers have made the process of obtain-
ing the base-by-base sequence of DNA...

- Analyze JSTOR's entire collection from *Science* (1880-2002)
- No reliable punctuation, meta-data, or references
- Restrict to 30K terms that occur more than ten times
- The data are 76M words in 130K documents

# Analyzing a document

**Original article**



**Topic proportions**

# Analyzing a document

**Original article**



TECHVIEW: DNA SEQUENCING

**Sequencing the Genome, Fast**

James C. Mullikin and Amanda A. McMurray

**Most likely words from top topics**

| | | |
|---|---|---|
| sequence | devices | data |
| genome | device | information |
| genes | materials | network |
| sequences | current | web |
| human | high | computer |
| gene | gate | language |
| dna | light | networks |
| sequencing | silicon | time |
| chromosome | material | software |
| regions | technology | system |
| analysis | electrical | words |
| data | fiber | algorithm |
| genomic | power | number |
| number | based | internet |

# Analyzing a topic

# Visualizing trends within a topic

## Time-corrected document similarity

- Consider the expected Hellinger distance between the topic proportions of two documents,

$$d_{ij} = \mathrm{E}\left[\sum_{k=1}^{K} (\sqrt{\theta_{i,k}} - \sqrt{\theta_{j,k}})^2 \,|\, \mathbf{w}_i, \mathbf{w}_j\right]$$

- Uses the latent structure to define similarity
- Time has been factored out because the topics associated to the components are different from year to year.
- Similarity based only on topic proportions

# Time-corrected document similarity

The Brain of the Orang (1880)

## Representation of the Visual Field on the Medial Wall of Occipital-Parietal Cortex in the Owl Monkey (1976)

# Quantitative comparison

- Compute the probability of each year's documents conditional on all the previous year's documents,

$$p(\mathbf{w}_t \mid \mathbf{w}_1, \ldots, \mathbf{w}_{t-1})$$

- Compare exchangeable and dynamic topic models

# Quantitative comparison

## Dynamic topic models discussion

- The DTM is a hierarchical model of sequential document collections;
- Exchangeability assumptions should be taken seriously.
- Variational methods allow large scale posterior inference.
- Examining the latent structure yields useful browsing tools
- Some open issues
  - Model selection: choosing the number of topics
  - Variational inference: what are the hidden assumptions?

# Outline

# The hidden assumptions of the Dirichlet distribution



- The Dirichlet is an exponential family distribution on the *simplex*, positive vectors that sum to one.
- However, the near independence of components makes it a poor choice for modeling topic proportions.
- An article about *fossil fuels* is more likely to also be about *geology* than about *genetics*.

# The logistic normal distribution



- The logistic normal is a distribution on the simplex that can model dependence between components.
- The natural parameters of the multinomial are drawn from a multivariate Gaussian distribution.

$$
\begin{aligned}
X &\sim \mathcal{N}_{K-1}(\mu, \Sigma) \\
\theta_i &= \exp\{x_i - \log(1 + \sum_{j=1}^{K-1} \exp\{x_j\})\}
\end{aligned}
$$

# Correlated topic model (CTM)



- Draw topic proportions from a logistic normal, where topic occurrences can exhibit correlation.
- Use for:
  - Providing a "map" of topics and how they are related
  - Better prediction via correlated topics

# Summary

- Topic models provide useful descriptive statistics for analyzing and understanding the latent structure of large text collections.

- More generally, probabilistic graphical models are a useful way to express assumptions about the hidden structure of complicated data.

- Variational methods allow us to perform posterior inference to automatically infer that structure from large data sets.

- Current research
  - Choosing the number of topics
  - Continuous time dynamic topic models
  - Topic models for prediction
  - Inferring the impact of a document

"We should seek out unfamiliar summaries of observational material, and establish their useful properties... And still more novelty can come from finding, and evading, still deeper lying constraints." (Tukey, 1962)

## Diversion: Variational inference

- Let $x_{1:N}$ be observations and $z_{1:M}$ be latent variables
- Our goal is to compute the posterior distribution

$$p(z_{1:M} \mid x_{1:N}) = \frac{p(z_{1:M}, x_{1:N})}{\int p(z_{1:M}, x_{1:N}) dz_{1:M}}$$

- For many interesting distributions, the marginal likelihood of the observations is difficult to efficiently compute

## Variational inference

- Use Jensen's inequality to bound the log prob of the observations:

$$\log p(x_{1:N}) \geq \mathrm{E}_{q_\nu}[\log p(z_{1:M}, x_{1:N})] - \mathrm{E}_{q_\nu}[\log q_\nu(z_{1:M})].$$

- We have introduced a distribution of the latent variables with free *variational parameters* $\nu$.

- We optimize those parameters to tighten this bound.

- This is the same as finding the member of the family $q_\nu$ that is closest in KL divergence to $p(z_{1:M} \,|\, x_{1:N})$.

# Mean-field variational inference

- Complexity of optimization is determined by the factorization of $q_\nu$
- In *mean field variational inference* we choose $q_\nu$ to be fully factored

$$q_\nu(z_{1:M}) = \prod_{m=1}^{M} q_{\nu_m}(z_m).$$

- The latent variables are independent.
    - Each is governed by its own variational parameter $\nu_m$.
- In the true posterior they can exhibit dependence
  (often, this is what makes exact inference difficult).

## MFVI and conditional exponential families

- Suppose the distribution of each latent variable conditional on the observations and other latent variables is in the exponential family:

$$p(z_m \,|\, \mathbf{z}_{-m}, \mathbf{x}) = h_m(z_m) \exp\{g_m(\mathbf{z}_{-m}, \mathbf{x})^T z_m - a_m(g_i(\mathbf{z}_{-m}, \mathbf{x}))\}$$

- Assume $q_\nu$ is fully factorized, and each factor is in the same exponential family:

$$q_{\nu_m}(z_m) = h_m(z_m) \exp\{\nu_m^T z_m - a_m(\nu_m)\}$$

## MFVI and conditional exponential families

- Variational inference is the following coordinate ascent algorithm

$$\nu_m = \mathrm{E}_{q_\nu}[g_m(\mathbf{Z}_{-m}, \mathbf{x})]$$

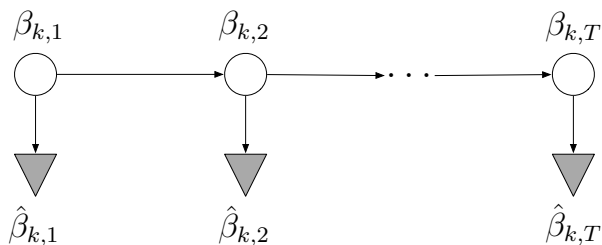- Notice the relationship to Gibbs sampling

# Variational family for the DTM



- Distribution of $\theta$ and $z$ is fully-factorized (Blei et al., 2003)
- Distribution of $\{\beta_{1,k}, \ldots, \beta_{T,k}\}$ is a *variational Kalman filter*
- Gaussian state-space model with free *observations* $\hat{\beta}_{k,t}$.
- Fit observations such that the corresponding posterior over the chain is close to the true posterior.

# Variational family for the DTM



- Given a document collection, use coordinate ascent on all the variational parameters until the KL converges.
- Yields a distribution close to the true posterior of interest
- Take expectations w/r/t the simpler variational distribution