

Modeling *Science*

David M. Blei

Department of Computer Science
Princeton University

April 17, 2008

Joint work with John Lafferty (CMU)

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. Novie and I drank of the vanilla extract which was used, and no ill results followed. Still, from this cream we isolated the same poison which I had before found in poisonous cheese (*Zeitschrift für physiologische Chemie*, x,

RNA Editing and the Evolution of Parasites

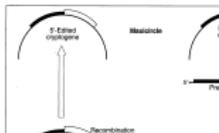
Lary Simpson and Dmitri A. Maslov

The kinoplastid flagellates, together with their sister group of euglenoids, represent the earliest stage of eukaryotic eukaryotes containing mitochondria (1). Within the kinoplastids, there are two major groups, the poorly studied bodonids and trypomastix, which consist of both free-living and parasitic cells, and the better known trypanosomatids, which are obligate parasites (2).

Perhaps because of the antiquity of the trypanosomatid lineage, these cells possess several unique genetic features (see accompanying Perspective by Nisbet)—one of which is RNA editing of mitochondrial transcripts. This RNA editing function (3-7) creates open reading frames ("cryptogenes") by inserting

uracil (U) nucleotides at a few specific sites within the coding region of an mRNA (U-editing) or at multiple sites, thus altering the mRNA (ppv-editing). The

edit, but there is disagreement on the nature of the primary genetic loss. The "vanilla theory" model (2, 11) states that the initial mutation was the loss of the *Centron* insertions. Covalent editing of genetic and base would have led to a wide distribution of trypanosomatids in insect and leishan. In this theory, divergent life cycles following symbiosis and vertebrate hosts evolved later as a result of the separation by new biogeography and dispersal of the ability to feed on the blood



size within genetic weak the ad in genetic nature genetic Cdk1

Chaotic Beetles

Charles Godfrey and Michael Hassel

Ecologists have known about the pioneering work of May in the mid 1970s (1) that the population dynamics of annual and plants can be remarkably 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 as 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. In such chaotic dynamic systems, complex nonlinearities for the management and conservation of natural resources. On page 308 of this issue, Comtet et al. (2) provide the case

The authors are in the Department of Ecology, Université Compiègne, France. Email: charles.godfrey@univ-compiègne.fr

convincing evidence to date of complex chaotic and chaos in a biological population—of the flour beetle, *Tribolium castaneum* (see Perspective by Nisbet). It has proved extremely difficult to demonstrate complex dynamics in populations in the field. In an very recent, a chaotic fluctuating population will specifically resemble a stable or cyclic population half the time the normal underlying behavior, experienced by all species. Given a long enough time series, diagnostic tools from nonlinear mathematics can be used to identify the subtle signatures of chaos in phase space, although experiments come to be an "strange attractor" which geometric objects with fractal structure and hence noninteger dimension. As they



Chaos in the field. *Tribolium castaneum*, an annual plant population, displays chaotic fluctuations in abundance when the amount of conspecifics is allowed to increase in a natural field.

move over the surface of the attractor, sets of adjacent trajectories are pulled apart, then stretched and folded, so that a becomes impossible to predict exact population densities over the year. The strength of the mixing that gives rise to the extreme sensitivity to initial conditions can be measured mathematically estimating the Lyapunov exponent, which is positive for chaotic dynamics and negative for stable systems. These last have many attempts to estimate at initial dimension and Lyapunov exponents from time series data, and some candidate chaotic populations have been identified (some insects, rodents, and more convincingly, human childhood disease). But the statistical difficulties preclude any broad generalizations (3).

An alternative approach to the parameter population results with less from natural populations and then compare their predictions with the observations in the field. This technique has been gaining popularity in recent years. In field, statistical advances in parameter estimation. Global or

SCIENCE • VOL. 271 • 17 JANUARY 1993

303

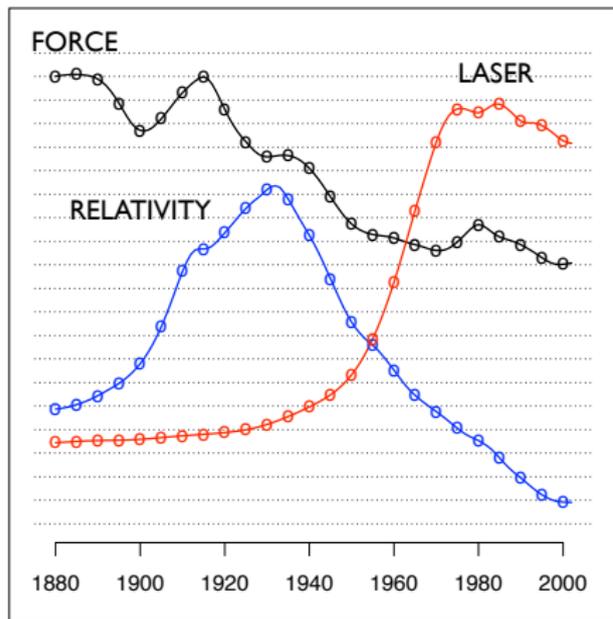
- On-line archives of document collections require better organization. Manual organization is not practical.
- Our goal: to discover the hidden thematic structure with hierarchical probabilistic models called *topic models*.
- Use this structure for browsing, search, and similarity.

Discover topics from a corpus

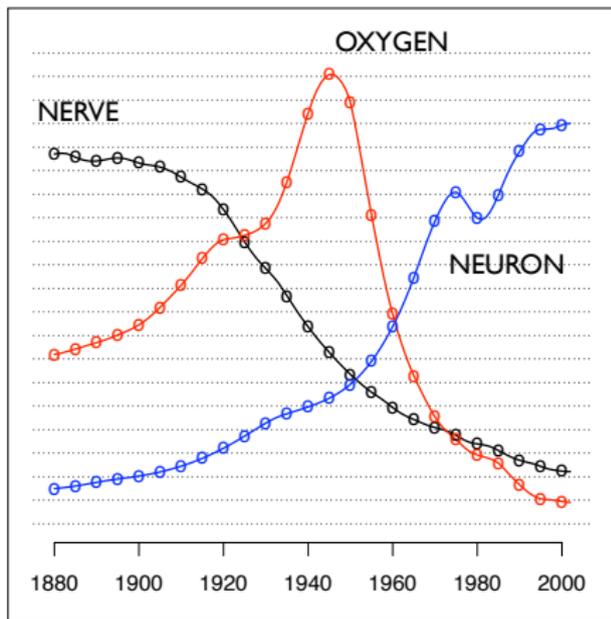
human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living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

Model the evolution of topics over time

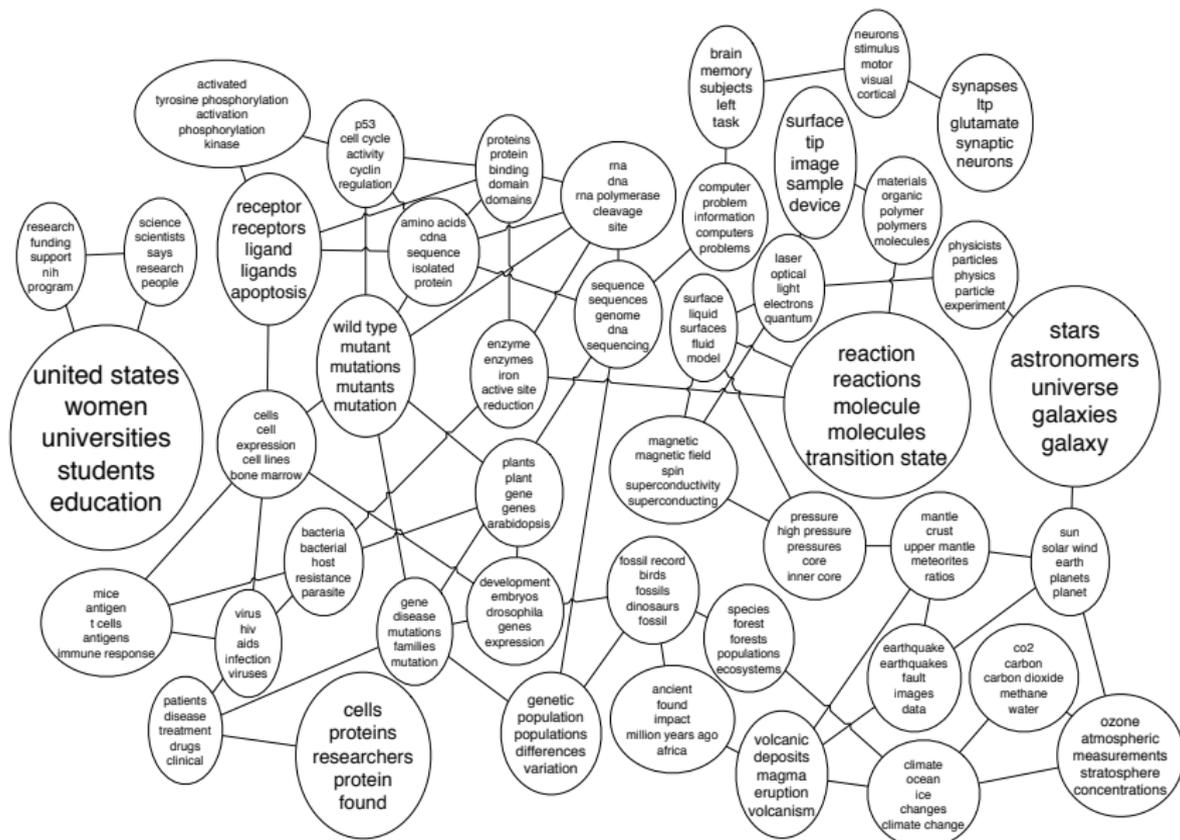
"Theoretical Physics"



"Neuroscience"



Model connections between topics



Outline

- 1 Introduction
- 2 Latent Dirichlet allocation
- 3 Dynamic topic models
- 4 Correlated topic models

Outline

- 1 Introduction
- 2 Latent Dirichlet allocation**
- 3 Dynamic topic models
- 4 Correlated topic models

Probabilistic modeling

- ① 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.
- ② Infer the hidden structure using *posterior inference*
 - What are the topics that describe this collection?
- ③ 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.

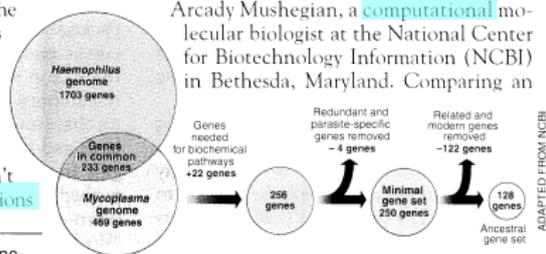
Although the numbers don't match precisely, those predictions

* 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 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



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

ADAPTED FROM NCBI

Simple intuition: Documents exhibit multiple topics.

Generative process

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,² 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.

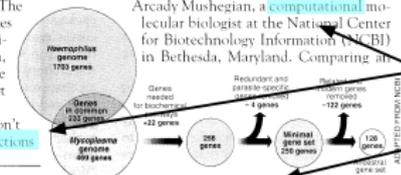
Although the numbers don't match precisely, those **predictions**

² 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 Anderson, a health 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 at



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

- 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

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,¹ 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.

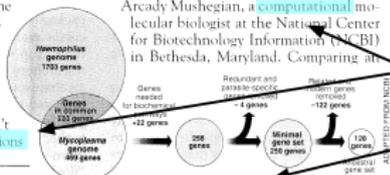
Although the numbers don't match precisely, those **predictions**

¹ 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 Anderson, a postdoctoral fellow at the University of Sussex. "But arriving 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

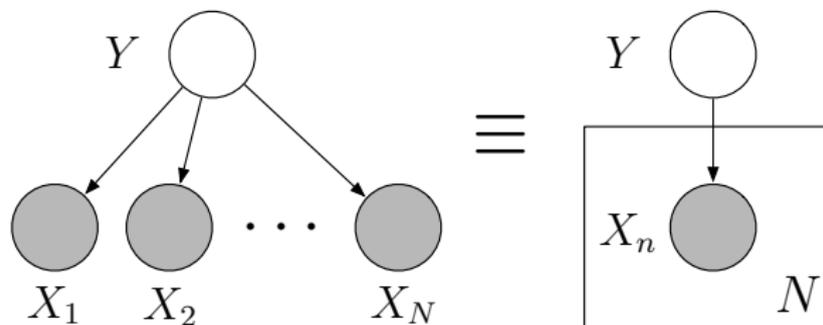
Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing at



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

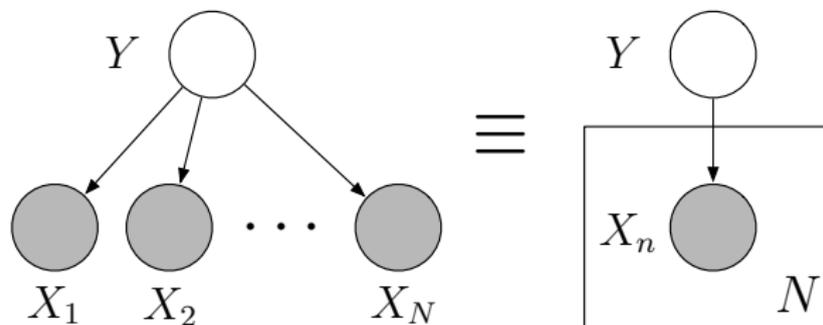
- In reality, we only observe the documents
- Our goal is to infer the underlying topic structure
 - What are the topics?
 - How are the documents divided according to those topics?

Graphical models (Aside)



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

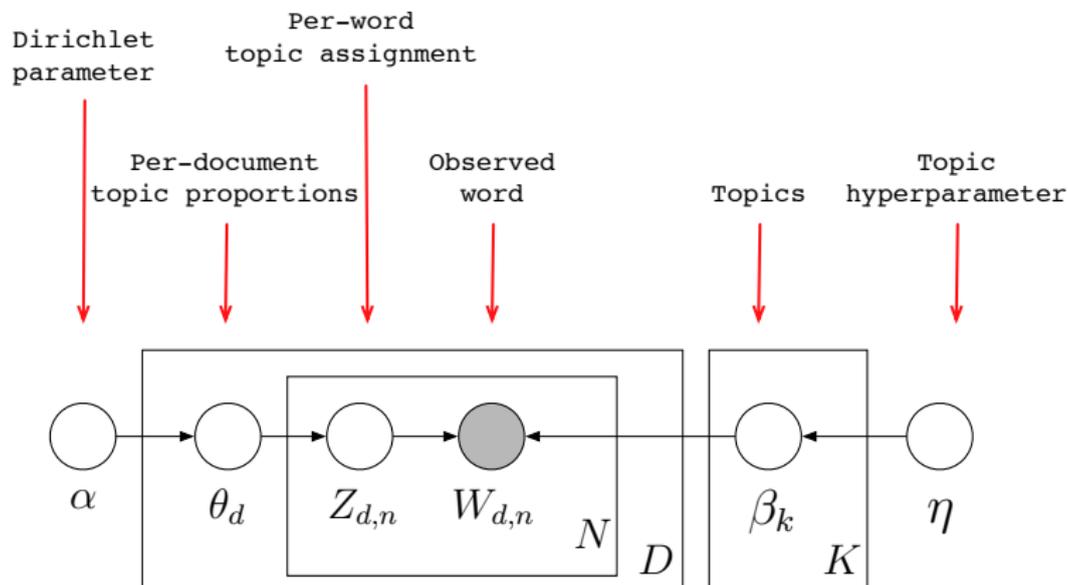
Graphical models (Aside)



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

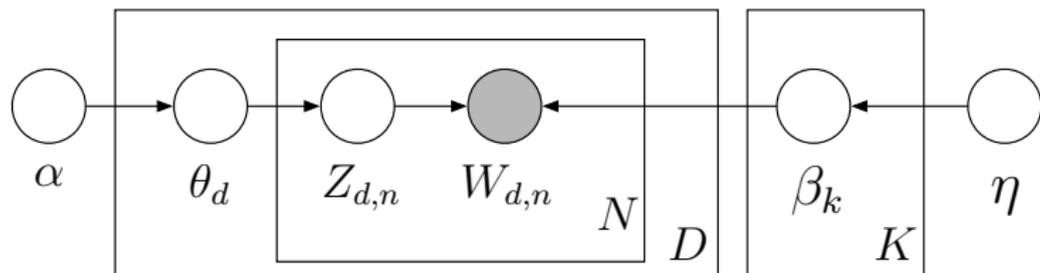
$$p(y, x_1, \dots, x_N) = p(y) \prod_{n=1}^N p(x_n | y)$$

Latent Dirichlet allocation



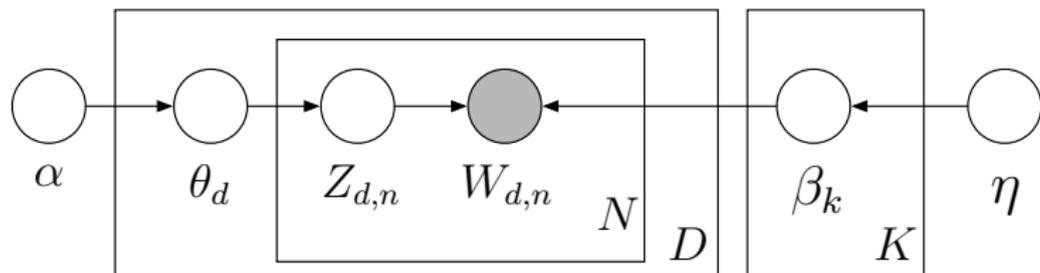
Each piece of the structure is a random variable.

Latent Dirichlet allocation



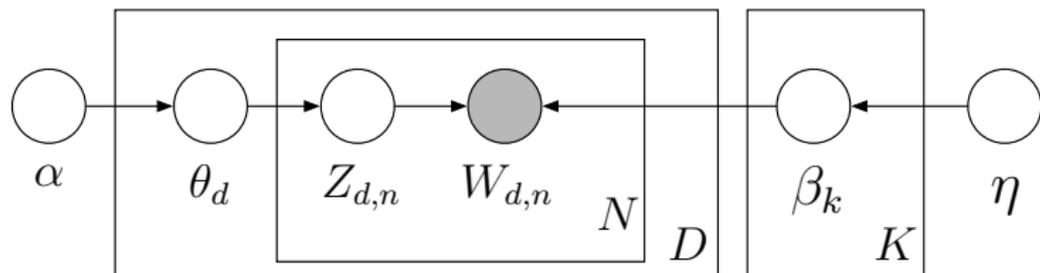
- 1 Draw each topic $\beta_i \sim \text{Dir}(\eta)$, for $i \in \{1, \dots, K\}$.
- 2 For each document:
 - 1 Draw topic proportions $\theta_d \sim \text{Dir}(\alpha)$.
 - 2 For each word:
 - 1 Draw $Z_{d,n} \sim \text{Mult}(\theta_d)$.
 - 2 Draw $W_{d,n} \sim \text{Mult}(\beta_{z_{d,n}})$.

Latent Dirichlet allocation



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

Latent Dirichlet allocation

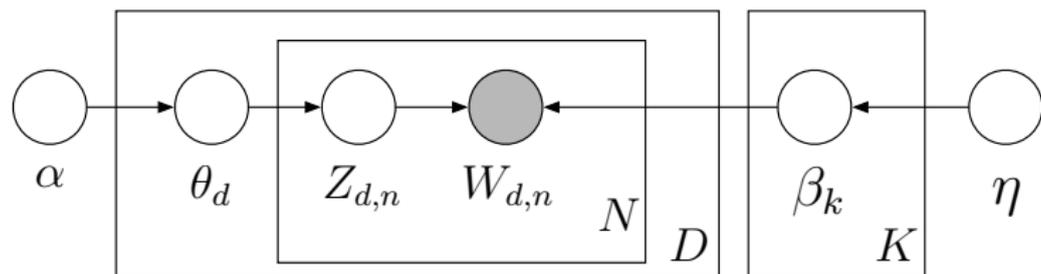


- Computing the posterior is intractable:

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

- 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

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.

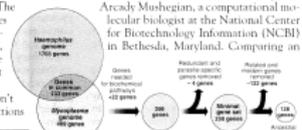
Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 25,000 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.

- **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

Seeking Life's Bare (Genetic) Necessities

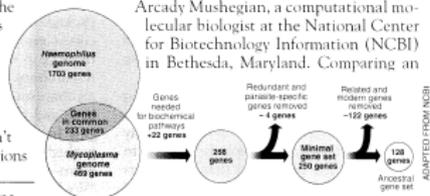
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known 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

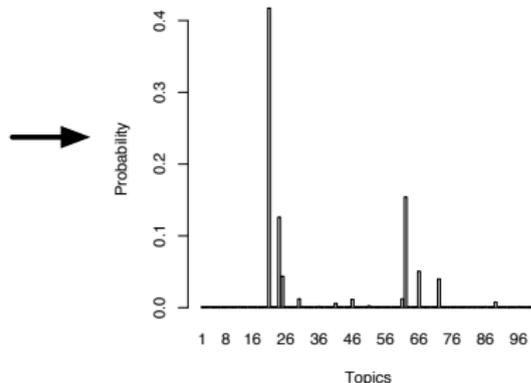
* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

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 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



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



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 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

LDA summary

- *Modular*: It can be embedded in more complicated models.
 - E.g., syntax and semantics; authorship; word sense
- *General*: The data generating distribution can be changed.
 - E.g., images; social networks; population genetics data
- Variational inference is fast; lets us to analyze large data sets.

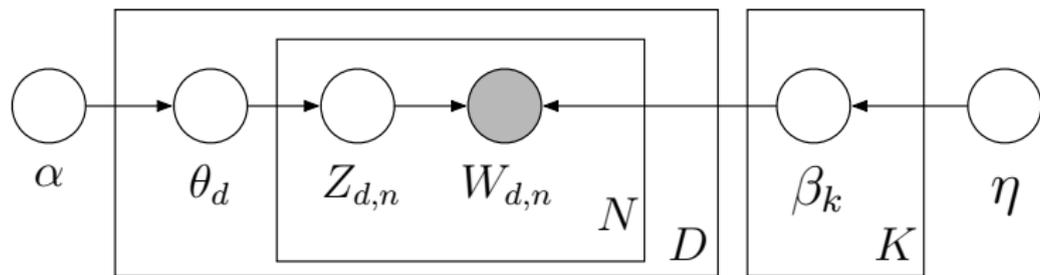
- See Blei et al., 2003 for details and a quantitative comparison.
- Code to play with LDA is freely available on my web-site, <http://www.cs.princeton.edu/~blei>.

- But, LDA makes certain assumptions about the data.
- When are they appropriate?

Outline

- 1 Introduction
- 2 Latent Dirichlet allocation
- 3 Dynamic topic models**
- 4 Correlated topic models

LDA and exchangeability

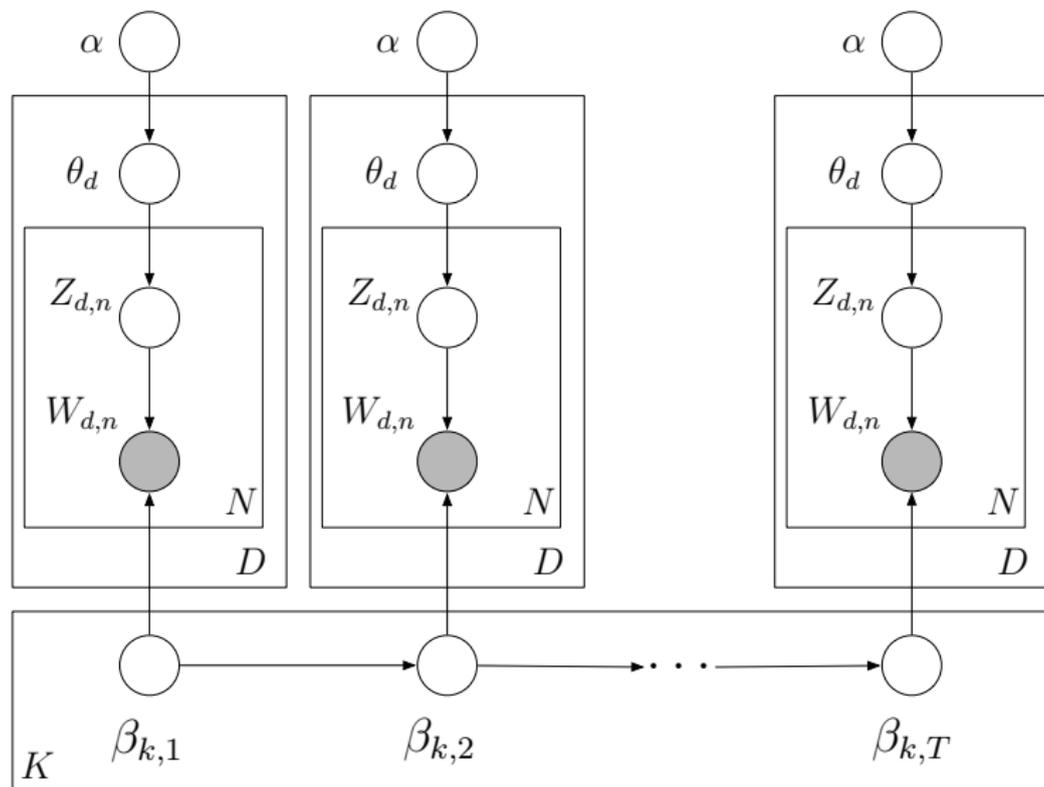


- LDA assumes that documents are exchangeable.
- I.e., their joint probability is invariant to permutation.
- This is too restrictive.

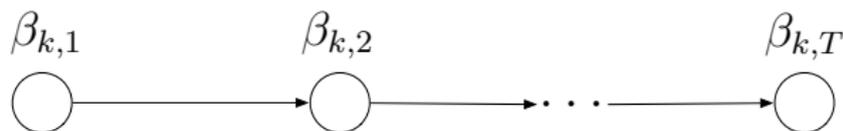
Dynamic topic model

- 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 evolving topics (Aitchison, 1980)
- A state-space model on the natural parameter of the topic multinomial (West and Harrison, 1997)

$$\beta_{t,k} | \beta_{t-1,k} \sim \mathcal{N}(\beta_{t-1,k}, I\sigma^2)$$

$$p(w | \beta_{t,k}) = \exp \left\{ \beta_{t,k} - \log \left(1 + \sum_{v=1}^{V-1} \exp\{\beta_{t,k,v}\} \right) \right\}$$

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

TECHVIEW

Sequencing the Genome, Fast

James C. Mullikin and Amanda A. McMurray

Combining sequencing accuracy and speed is a major challenge in the field of genomics. The first step in sequencing a genome is to break it up into small fragments, called clones, which are then sequenced individually. The first step in sequencing a genome is to break it up into small fragments, called clones, which are then sequenced individually. The first step in sequencing a genome is to break it up into small fragments, called clones, which are then sequenced individually.



Fig. 1 A comparison of sequencing techniques as of 2000. The number of clones sequenced per day is shown for the Sanger method (red line), the Maxam-Gilbert method (blue line), and the shotgun method (green line). The shotgun method shows a significant increase in the number of clones sequenced per day, especially after 1995.

Fig. 2 A comparison of sequencing techniques as of 2000. The number of clones sequenced per day is shown for the Sanger method (red line), the Maxam-Gilbert method (blue line), and the shotgun method (green line). The shotgun method shows a significant increase in the number of clones sequenced per day, especially after 1995.

TECHVIEW: DNA SEQUENCING

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 information necessary for the life of the organism. The base sequence contains four nucleotides—adenine, thymidine, guanosine, and cytosine—which are linked together into long double-helical chains. Over the last two decades, automated DNA sequencers have made the process of obtaining the base-by-base sequence of DNA...

- Analyze JSTOR's entire collection from *Science* (1880-2002)
- Restrict to 30K terms that occur more than ten times
- The data are 76M words in 130K documents

Original article

Most likely words from top topics



TECHVIEW: DNA SEQUENCING

Sequencing the Genome, Fast

James C. Mullikin and Amanda A. McPherson

Genome sequencing projects reveal the genetic makeup of an organism by reading the sequence of the DNA bases, which encode all of the information necessary for the life of the organism. The basic sequence consists of four nucleotides—adenine, thymine, guanine, and cytosine—which are joined together in long double-helical chains. Over the last two decades, automated DNA sequencers have made the process of obtaining the base-by-base sequence of DNA easier. By application of an electric field across a gel matrix, these sequencers separate fluorescently labeled DNA molecules that differ in size by one base. As the molecules move past a given point in the gel, laser excitation of a fluorescent dye specific to the base at the end of the molecule yields a base-specific signal that can be automatically recorded.

The latest sequencer to be launched is Perkin-Elmer's next-generation ABI Prism 7700 DNA sequencer, which, like the Molecular Dynamics MegabACE 2000 launched last year, incorporates a capillary tube to hold the sequencing gel rather than a traditional slab-gel design. It offers a major improvement to the ABI 3700 but has gone further than the MegabACE 2000 because of its use of a capillary tube to hold the sequencing gel rather than a traditional slab-gel design. It offers a major improvement to the ABI 3700 but has gone further than the MegabACE 2000 because of its use of a capillary tube to hold the sequencing gel rather than a traditional slab-gel design. It offers a major improvement to the ABI 3700 but has gone further than the MegabACE 2000 because of its use of a capillary tube to hold the sequencing gel rather than a traditional slab-gel design.



Fig. 5. Comparison of read-length histograms for the sequencer column with the ABI 3700 sequencer and the ABI 7700. The top histogram shows the sequencer column with the ABI 3700 sequencer and the ABI 7700. The bottom histogram shows the sequencer column with the ABI 7700 sequencer. Both sets of data are from runs with the dye chemistry described. Read length is computed as the number of bases per read where the predicted error rate is less than or equal to 1.0% (2). The "stitch" Q value was calculated for each type of run.

to the Sanger Centre in December 1998... are in our Research and Development department for evaluation. Thus, the ABI 3700 will ultimately be added to our product capacity to meet our end... The ABI 3700 DNA sequencer is built into a flow-stationed column, which consists in its base all the reagents required for its operation. The reagent solutions are manually accessed for the replenishment, which is required every day under high-throughput operation. As bench length within the column is a fixed parameter, the number of protein plates of DNA samples are loaded. The operator places the prepared plates on to protein, above the fluid of the reaction and programs it by using a personal computer. A robotic arm transfers DNA sam-

TECHSIGHT

ples from the plates into wells that open to the capillaries. This and the rest of the sequencing operation is fully automated. The machine can currently process four to seven plates of DNA samples automatically, taking approximately 15 hours before operator intervention is required. This new laboratory of the design specification of four to six plates a 12-hour run.

The main innovation of the ABI 3700 is the use of a short flow fluorescence detection system (4). Detection of the DNA fluorescence occurs in the path of the capillary within a fixed silica capillary. A laser beam flows over the ends of the capillaries, drawing the DNA fragments as they emerge from the capillaries through a fixed laser beam that automatically synchronizes with all of the samples. The emitted fluorescence is detected with a special CCD (charge-coupled device) detector. This arrangement means that there are no moving parts in the detection system, other than a shutter in front of the CCD device.

We have evaluated these machines for their performance, significant ease of use, and reliability in comparison to the more commonly used slab gel sequencing methods. In automated sequencing, there are two methods for connecting the gel matrix between two fixed support glass plates (4) or one to one—the slab gel method. The other is to inject a polymer matrix into a capillary (5). Most sequencing facilities use the slab gel method, because microfluidic sequencing have only recently become commercially available.

With either type of process, the aim is to avoid as many hours as possible for a given sample of DNA—that is, long read lengths are desirable. In fact, a system that could read faster on many bases but at the expense of another system is preferable, if both systems cost the same. This is because sequencing relatively fewer long-range fragments is easier than sequencing many short ones. So, read length an important parameter when evaluating sequencing technologies.

We have directly compared the ABI 3700 sequencer to the ABI 3770X slab sequencer by evaluating the sequence data obtained with both machines with human DNA samples. These samples were subcloned into plasmids of 1 to 11 kb and prepared and sequenced with our standard procedure for Perkin-Elmer Big Dye Terminator chemistry.

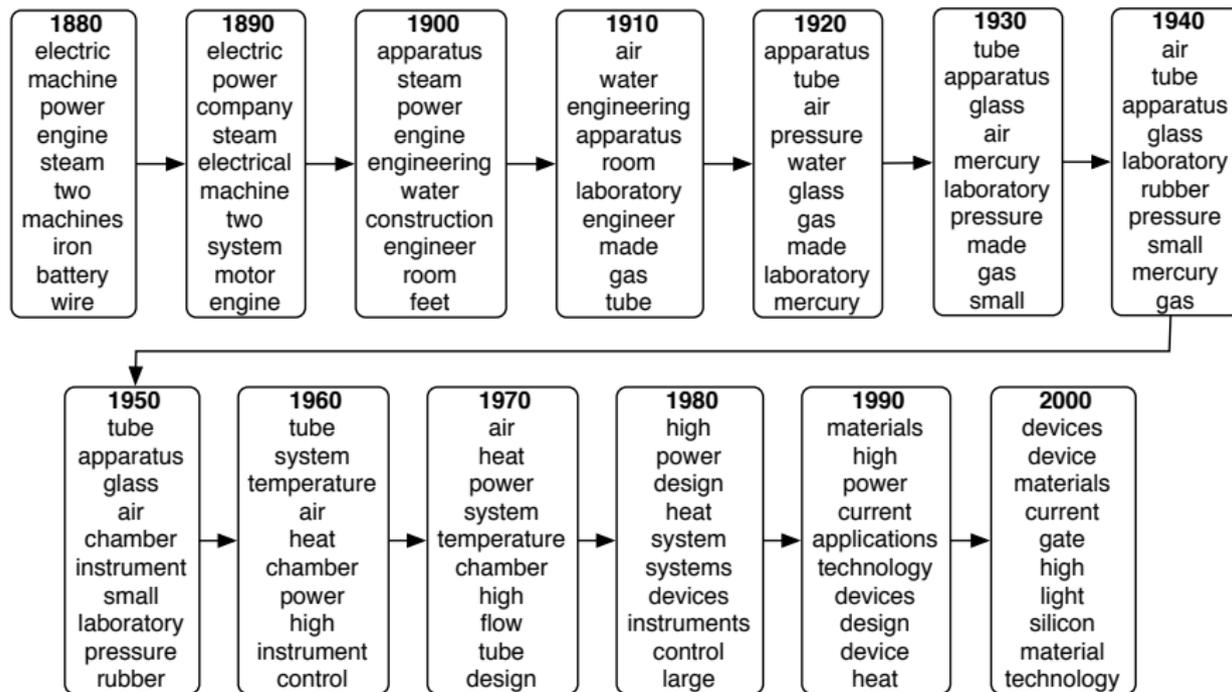
sequence
genome
genes
sequences
human
gene
dna
sequencing
chromosome
regions
analysis
data
genomic
number

devices
device
materials
current
high
gate
light
silicon
material
technology
electrical
fiber
power
based

data
information
network
web
computer
language
networks
time
software
system
words
algorithm
number
internet

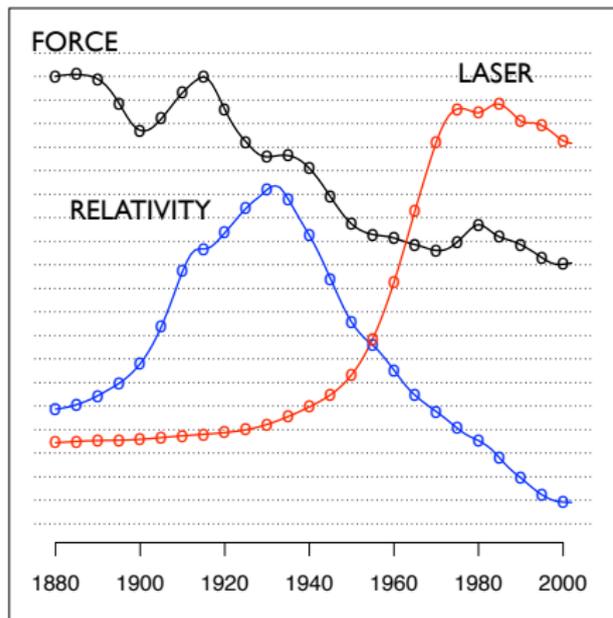
The authors are at the Sanger Centre, Wellcome Trust Genome Campus, Hinxton, Cambridgeshire CB10 1TA, UK (e-mail: jcm@sanger.ac.uk).

Analyzing a topic

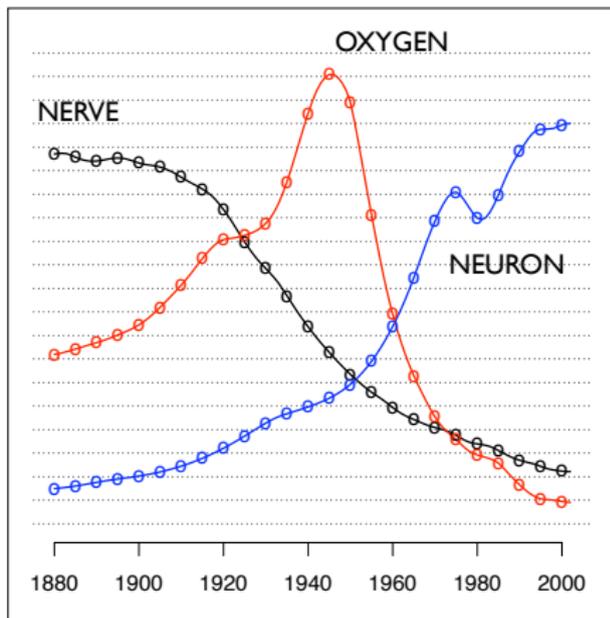


Visualizing trends within a topic

"Theoretical Physics"



"Neuroscience"



Time-corrected document similarity

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

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

- 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)

366

SCIENCE.

Published in other forms, which were submitted to the authors in the fall of 1880, and by comparison of the numbers, we suppose the original must be printed from a second number. After publication Professor Agassiz now writes that the figures under his name are not identical with his. We therefore request our readers to consider these engravings.

Professor George F. Barker, Professor G. C. Mack and Professor J. S. Huxley are preparing more valuable reports of their respective parts, and promise them at an early day.

THE BRAIN OF THE ORANG.*

BY GEORGE F. BARKER.

The brain of the Orang has been figured by Verhagen, Sandberg, Schroeder van der Kolk, and Van der Graaf, in Holland, etc. On account, however, of the low illustration and also of the imperfection of the sections, I send a series of the representations of presenting several views of my Orang's brain (fig. 1, to 11, which were prepared from the skull and a few bones also).

The cerebrum with its high level of convolution, and a ridge on the surface of the left hemisphere had been designated by Verhagen, Schroeder van der Kolk, and myself as the *convex surface*. It weighed exactly 300 grams. The brain of the Orang, Schroeder van der Kolk, and myself, were more than those of either of the Chimpanzees which I examined. In these the brain was more elongated. The general character of the table is likewise in



FIG. 1.

The brain of the Orang, Chimpanzee, and man are the same. There are only two main differences between the brain of the Orang and that of man. The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded.

* From the Proceedings of the Academy of Natural Sciences, Philadelphia, 1880, p. 100.

in the Orang, the post-central sulcus does not reach the occipital lobe, as it does in the Chimpanzee and the human brain. The sulcus is more deeply marked in the human brain. The sulcus is more deeply marked in the human brain. The sulcus is more deeply marked in the human brain. The sulcus is more deeply marked in the human brain.

According to Verhagen, the distribution of the brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded.

The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded.

The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded.

The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded.

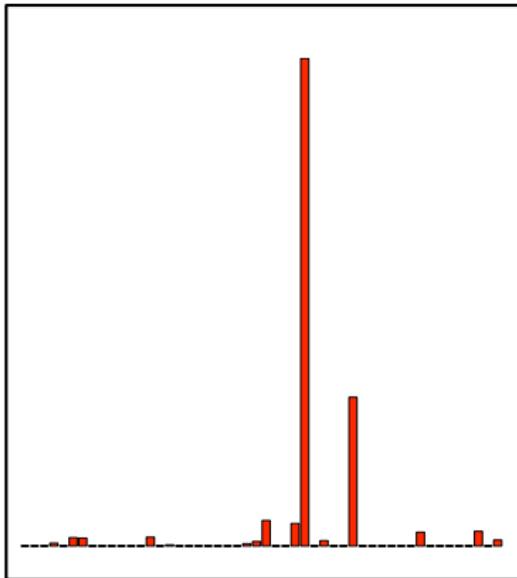
The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded. The brain of the Orang is more elongated and the posterior part of the brain is more rounded.



FIG. 2.

The brain of the Chimpanzee is more elongated and the posterior part of the brain is more rounded. The brain of the Chimpanzee is more elongated and the posterior part of the brain is more rounded. The brain of the Chimpanzee is more elongated and the posterior part of the brain is more rounded.

* From the Proceedings of the Academy of Natural Sciences, Philadelphia, 1880, p. 100.



Time-corrected document similarity

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

project, the stereoscopic representation of the medial occipital-parietal cortex was re-plotted with stereotaxiological coordinates in five owl monkeys (15). The monkeys were anesthetized with sodium pentobarbital and prepared for recording. Longitudinal and planar-section microelectrodes were used to record from small clusters of neurons or occasionally from single neurons at temporal positions parallel to the medial surface of occipital-parietal cortex. Receptive fields were plotted by moving a target spot or microstimulus into and out on the surface of a translucent plastic hemisphere centered in front of the animal's eyes. The position of the optokinetic was projected onto the flat surface of the hemisphere with the method of Frazar and Chase (6). The optokinetic eye usually was

covered with an opaque shield. Electrode and recording sites were recorded from histological sections and photographs of the entire brain.

Figure 1 illustrates the data from one monkey, complete mapping of the individual data obtained in the other four experiments revealed the same pattern of receptive organization. Temporal positions 1 through 4 are parallel to the medial surface of occipital-parietal cortex at a distance of approximately 1 cm from the medial surface. In previously published reports, we found that the receptive fields recorded adjacent to the medial area in the medial visual area (V II) were located in the lower quadrant and the horizontal meridian about 5P or 6P from the center of P. Thus, as is shown in Fig. 1, and

also in Fig. 2, which illustrates the organization of the other medial visual areas that have been mapped in the owl monkey, the border between the medial area and the caudal visual area corresponds to a peripheral portion of the horizontal meridian. In other experiments in the owl monkey we found that receptive fields recorded near the caudal border with the medial area began near the vertical meridian in the lower quadrant and proceeded in a broad loop in the periphery toward the horizontal meridian (7). Thus, as is shown in Figs. 1 and 2, the caudal border between the dorsal and the medial areas corresponds to part of the lower field vertical meridian and the peripheral portions of the lower visual quadrant. Therefore, the medial area is adjacent to por-

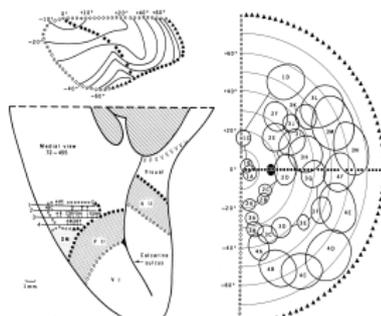
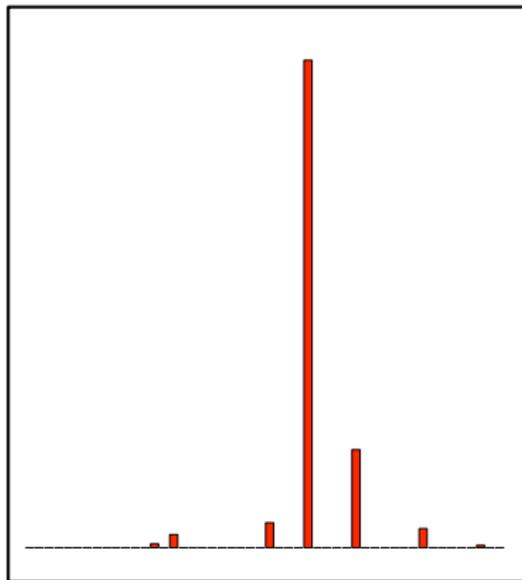


Fig. 1. Microelectrode recording positions and receptive field data for the medial visual area in owl monkey 12-475. The diagram on the top left is a view of the posterior wall of the medial wall of occipital-parietal cortex of the owl monkey with the electrode positions indicated. Stereotaxiological coordinates are given in degrees below the horizontal meridian and in degrees above the vertical meridian. The diagram on the middle left is a view of the medial wall of the occipital-parietal cortex of the owl monkey with the electrode positions indicated. Stereotaxiological coordinates are given in degrees below the horizontal meridian and in degrees above the vertical meridian. The diagram on the right is a view of the visual field with the electrode positions indicated. The diagram on the right is a view of the visual field with the electrode positions indicated. The diagram on the right is a view of the visual field with the electrode positions indicated. The diagram on the right is a view of the visual field with the electrode positions indicated.

11 FEBRUARY 1976

11



Automatic Analysis, Theme Generation, and Summarization of Machine-Readable Texts

Gerard Salton, James Allan, Chris Buckley.

Vast amounts of text material are now available in machine-readable form. Here, approaches are outlined for manipulating and accessing subject areas in accordance with user needs. In particular, methods for determining text themes, traversing texts selectively, and extracting summary text content.

Many kinds of texts are currently available in machine-readable form and are amenable to automatic processing. Because the available databases are large and cover many different subject areas, automatic aids must be provided to users interested in accessing the data. It has been suggested that links be placed between related pieces of text, connecting, for example, particular text paragraphs to other paragraphs covering related subject matter. Such a linked text structure, often called hypertext, makes it possible for the reader to start with particular text passages and use the linked structure to find related text elements (1). Unfortunately, until now, viable methods for automatically building large hypertext structures and for using such structures in a sophisticated way have not been available. Here we give methods for constructing text relation maps and for using text relations to access and use text databases. In particular, we outline procedures for determining text themes, traversing texts selectively, and extracting summary statements that reflect text content.

Text Analysis and Retrieval: The Smart System

The Smart system is a sophisticated text retrieval tool, developed over the past 30 years, that is based on the vector space

The authors are in the Department of Computer Science, Cornell University, Ithaca, NY 14853-7501, USA.

model of retrieval, all information as well as information presented by sets, or vectors, is typically a word, associated with a particular document. In principle, chosen from a thesaurus, but being constructed such that they are unrestrictive to derive the terms under consideration. Terms assigned to a text content.

Because the term for content representation introduces a term-weighting scheme, high weights to terms and lower weights to terms. A powerful term-weighting scheme is the well-known term frequency-inverse document frequency ($f_i \cdot \frac{1}{\sqrt{d_i}}$), which is a low frequency (f_i). Such terms that occur frequently in documents.

When all texts are represented by weighted vectors, the weight assigned to each term is the cosine of the angle between pairs of vectors. Thus, $\cos(\theta) = \frac{D_1 \cdot D_2}{\|D_1\| \|D_2\|}$.

SCIENCE • VOL. 273

TOPIC	PROB
data computer system information network	0.30
information library text index libraries	0.19
two three four different single	0.16

DOCUMENT	SCORE
"Global Text Matching for Information Retrieval" (1991)	0.2570
"Automatic Text Analysis" (1970)	0.3110
"Gauging Similarity with n-Grams: Language-Independent Categorization of Text" (1995)	0.3210
"Developments in Automatic Text Retrieval" (1991)	0.3480
"Simple and Rapid Method for the Coding of Punched Cards" (1962)	0.3610
"Data Processing by Optical Coincidence" (1961)	0.4290
"Pattern-Analyzing Memory" (1976)	0.4320
"The Storing of Pamphlets" (1899)	0.4440
"A Punched-Card Technique for Computing Means, Standard Deviations, and the Product-Moment Correlation Coefficient and for Listing Scattergrams" (1946)	0.4550

Global Text Matching for Information Retrieval

GERARD SALTON* and CHRIS BUCKLEY

An approach is outlined for the retrieval of natural language texts in response to available search requests and for the recognition of content similarity between text excerpts. The proposed retrieval process is based on flexible text matching procedures carried out in a number of different text environments and is applicable to large text collections covering unrestricted subject matter. For unrestricted text environments this system appears to outperform other currently available methods.

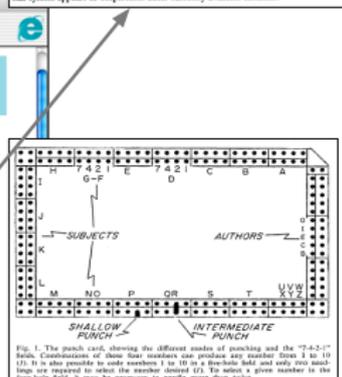


Fig. 1. The punch card, showing the different modes of punching and the "7-4-2-1" code. Combinations of these four numbers can produce any number from 1 to 10. It is also possible to code numbers 11 to 10 in a five-hole field and code zero markings are required to select the number desired (1). To select a given number in the code hole field, a mark is necessary to denote more than zero.

THE STORING OF PAMPHLETS.

On reading Professor Minot's explanation of his method of storing pamphlets as given in the issue of December 30th I feel inclined to add a word in commendation of the method. I began using these boxes six or seven years ago and now have 152 upon my shelves. About one-half are devoted to Experiment Station bulletins, the boxes being labeled by States and arranged alphabetically. The other half is used for miscellaneous pamphlets on subjects pertaining to my line of work. The boxes have proved perfectly satisfactory in every way, and as a simple time-saving device they are worth many times the cost. My system of pamphlet arrangement differs in some ways from that adopted by Professor Minot and has been adopted only after trial of several other methods.

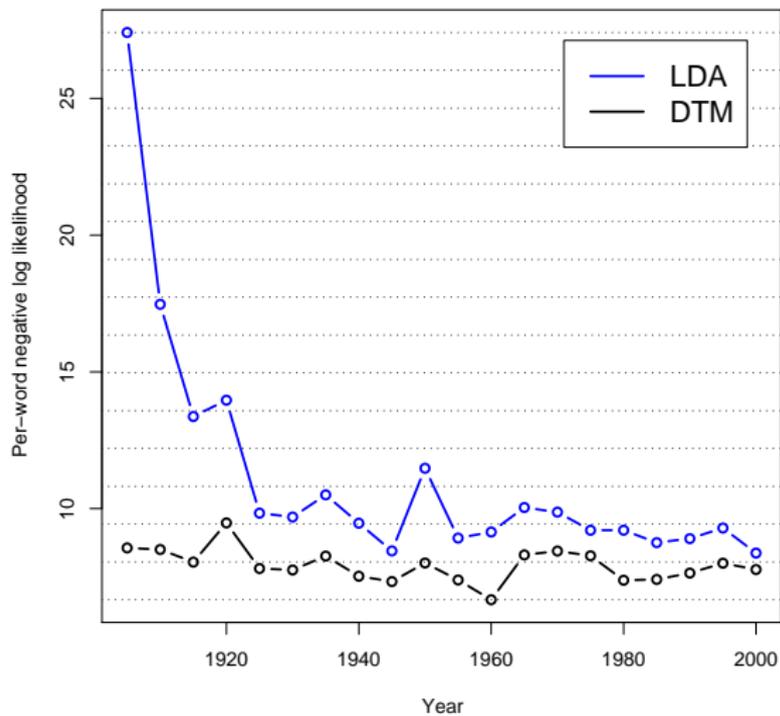
Quantitative comparison

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

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

- Compare exchangeable and dynamic topic models

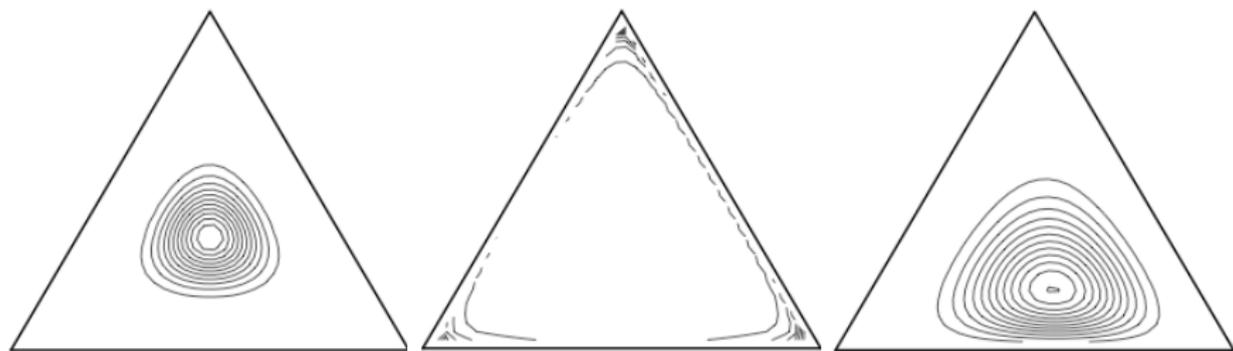
Quantitative comparison



Outline

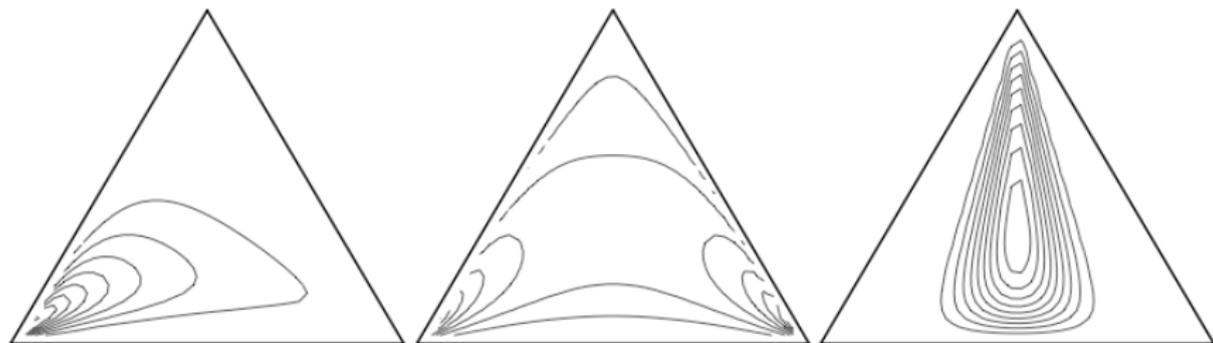
- 1 Introduction
- 2 Latent Dirichlet allocation
- 3 Dynamic topic models
- 4 Correlated topic models**

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

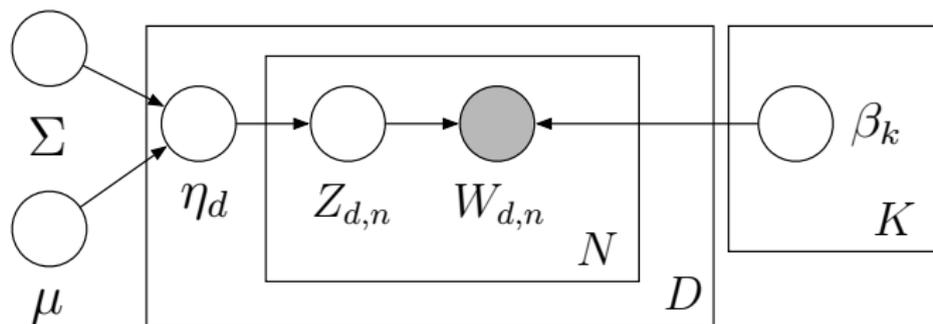


- The logistic normal is a distribution on the simplex that can model dependence between components.
- The natural parameters of the multinomial are drawn from a multivariate Gaussian distribution.

$$X \sim \mathcal{N}_{K-1}(\mu, \Sigma)$$

$$\theta_i = \exp\{x_i - \log(1 + \sum_{j=1}^{K-1} \exp\{x_j\})\}$$

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

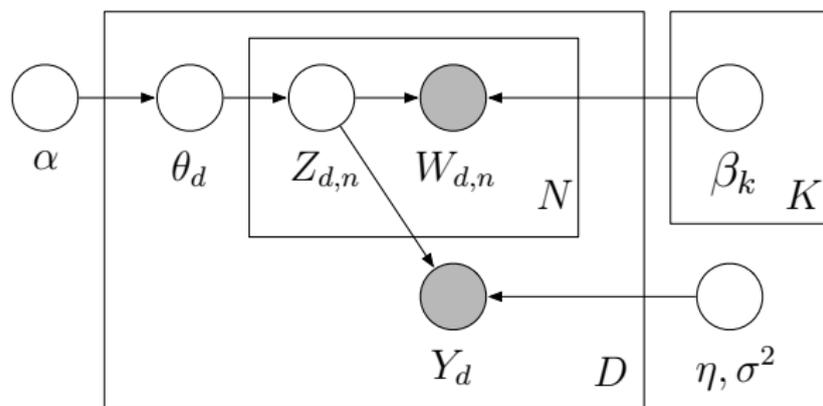
- Topic models provide useful descriptive statistics for analyzing and understanding the latent structure of large text collections.
- Probabilistic graphical models are a useful way to express assumptions about the hidden structure of complicated data.
- Variational methods allow us to perform posterior 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.”
(John Tukey, *The Future of Data Analysis*, 1962)

Supervised topic models (with Jon McAuliffe)

- Most topic models are *unsupervised*. They are fit by maximizing the likelihood of a collection of documents.
- Consider documents paired with response variables.
For example:
 - Movie reviews paired with a number of stars
 - Web pages paired with a number of “diggs”
- We develop *supervised topic models*, models of documents and responses that are fit to find topics predictive of the response.

Supervised LDA



- 1 Draw topic proportions $\theta \mid \alpha \sim \text{Dir}(\alpha)$.
- 2 For each word
 - 1 Draw topic assignment $z_n \mid \theta \sim \text{Mult}(\theta)$.
 - 2 Draw word $w_n \mid z_n, \beta_{1:K} \sim \text{Mult}(\beta_{z_n})$.
- 3 Draw response variable $y \mid z_{1:N}, \eta, \sigma^2 \sim \text{N}(\eta^\top \bar{z}, \sigma^2)$, where

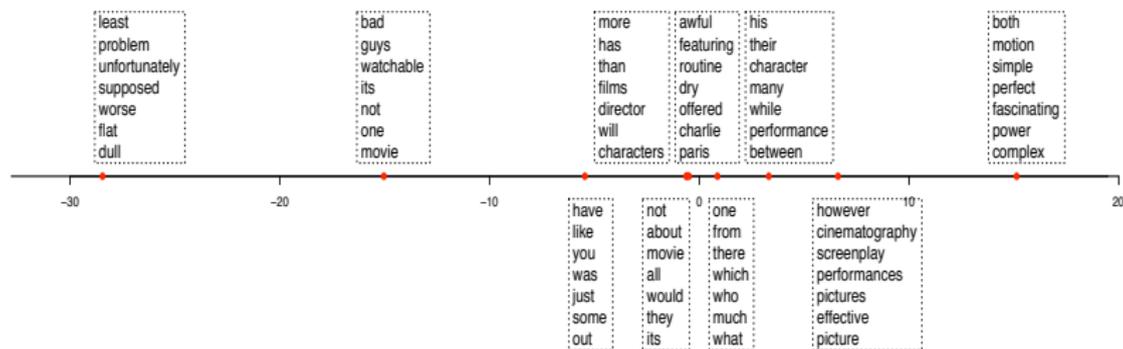
$$\bar{z} = (1/N) \sum_{n=1}^N z_n.$$

- SLDA is used as follows.
 - Fit coefficients and topics from a collection of document-response pairs.
 - Use the fitted model to predict the responses of previously unseen documents,

$$E[Y | w_{1:N}, \alpha, \beta_{1:K}, \eta, \sigma^2] = \eta^\top E[\bar{Z} | w_{1:N}, \alpha, \beta_{1:K}].$$

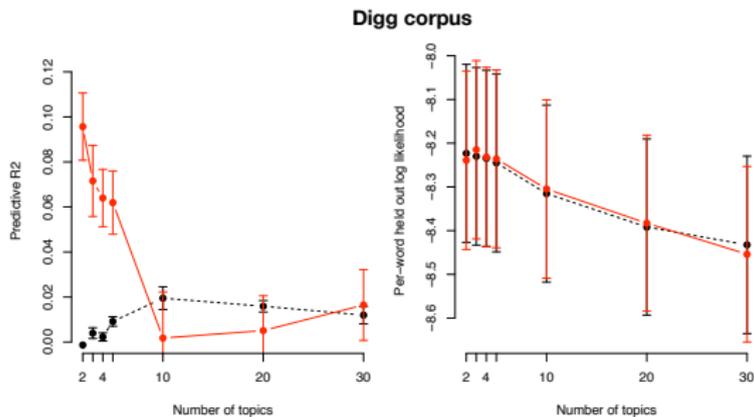
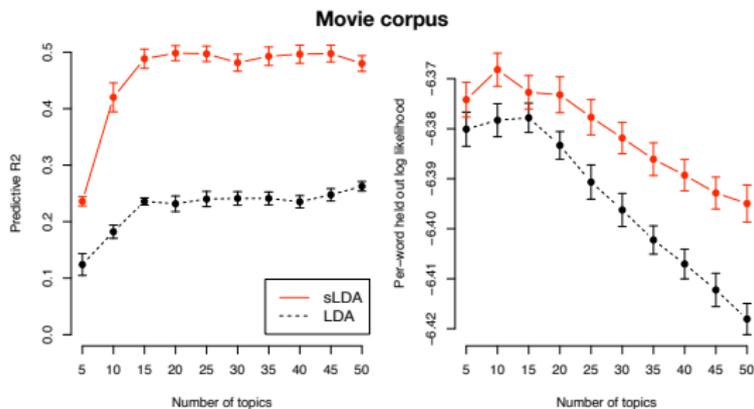
- The process enforces that the document is generated first, followed by the response. The response is generated from the particular topics that were realized in generating the document.

Example: Movie reviews



- We fit a 10-topic sLDA model to movie review data (Pang and Lee, 2005).
 - The documents are the words of the reviews.
 - The responses are the number of stars associated with each review (modeled as continuous).
- Each component of coefficient vector η is associated with a topic.

Simulations



Diversion: Variational inference

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

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

- 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:

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

- We have introduced a distribution of the latent variables with free *variational parameters* ν .
- We optimize those parameters to tighten this bound.
- This is the same as finding the member of the family q_ν that is closest in KL divergence to $p(z_{1:M} | x_{1:N})$.

Mean-field variational inference

- Complexity of optimization is determined by factorization of q_v
- In *mean field variational inference* q_v is fully factored

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

- The latent variables are independent.
 - Each is governed by its own variational parameter v_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:

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

- Assume q_v is fully factorized and each factor is in the same exponential family:

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

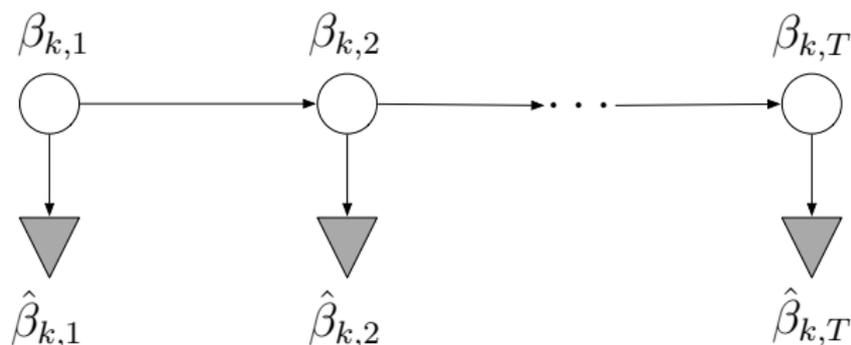
MFVI and conditional exponential families

- Variational inference is the following coordinate ascent algorithm

$$\nu_m = E_{q_v} [g_m(\mathbf{Z}_{-m}, \mathbf{x})]$$

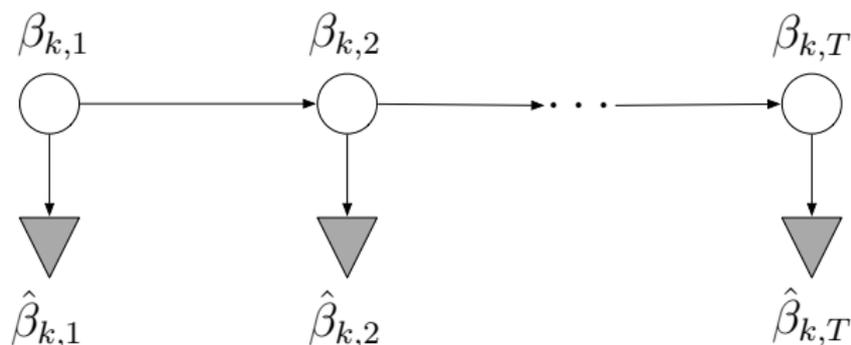
- Notice the relationship to Gibbs sampling

Variational family for the DTM



- Distribution of θ and z is fully-factorized (Blei et al., 2003)
- Distribution of $\{\beta_{1,k}, \dots, \beta_{T,k}\}$ 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