

# Modeling *Science*:

Discovering Themes in Large Collections of Documents

David M. Blei

Department of Computer Science  
Princeton University

May 14, 2007

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 outward 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 (*Zellwurst für physiologische chemie*, 3,

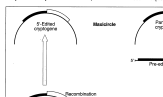
## RNA Editing and the Evolution of Parasites

Larry Simpson and Dmitri A. Maslov

The kinetoplast flagellates, together with their close group of eukaryotes, represent the earliest extant lineage of eukaryotic organisms containing mitochondria (1). Within the kinetoplasts, there are two major groups, the poorly studied bodonids and trypomastixids, which consist of both free-living and parasitic cells, and the better known trypanosomatids, which are obligate parasites (2).

Perhaps because of the sensitivity of the trypanosomatid lineage, these cells possess several unique genetic features (one accompanying Perspective by S. Hall—some of which is RNA editing of mitochondrial transcripts. This RNA editing function (3–7) causes open reading frames to “acquire” the insertion (or occasional deletion) of uridine (U) nucleotides at a few specific sites within the coding region of an mRNA (5′-editing) or at multiple specific sites throughout the mRNA (3′-editing). The

fact that there is disagreement on the nature of the primary parasitic host. The “parasitoid first” model (8, 11) states that the initial parasitism was in the gut of pre-Cambrian invertebrates. Coevolution of parasite and host would have led to a wide distribution of trypanosomatids in insects and leeches. In the theory, diagenetic life cycles (alternating vertebrate and invertebrate hosts) evolved later as a result of the acquisition by some trypanosomatids and diplopods of the ability to feed on the blood



size  
ant  
nuc  
wo  
the  
al  
in  
rho  
mit  
que  
Calk

cre  
mick  
as an  
ctual  
Toga  
the 3  
by 46  
fals p  
size  
tryp  
bran  
apar

**Ecologists have known since the pioneering work of May in the mid 1950s (1) that the population dynamics of animals and plants can be remarkably complex. This complexity arises from two sources. The tangled web of interactions that constitute one 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 are thus persistent oscillatory dynamical systems, the latter characterized by intrinsic sensitivity to initial conditions. If such chaotic dynamics were common to nature, then this would have important ramifications for the management and conservation of natural resources. On page 106 of this issue, Costantini et al. (2) provide the most**

This article is part of the Department of Ecology, Evolution, and Systematics, University of Michigan, 616 Tappan St., Ann Arbor, MI 48106-1103. E-mail: larry@umich.edu

## Chaotic Beetles

Charles Godfray and Michael Hassell

convincing evidence for the existence of complex dynamics and chaos in a biological population is that of the flour beetle, *Tribolium castaneum* (see figure).

If the process were not difficult to demonstrate complex dynamics in populations in the field, by no other means, a chaotic fluctuating population will specifically resemble a stable or cyclic population but hardly the normal-looking perturbations experienced by all species. Given a long enough time series, diagnostic methods, nonlinear mathematics can be used to identify the underlying mechanisms of chaos. In these species, chaotic trajectories correspond to the “strange attractors” common to persistent systems with feedback structure and hence recurrent dimensions. As they



**Contribution and chaos:** The flour beetle, *Tribolium castaneum*, exhibits chaotic population dynamics when the amount of available food is limited to a nontrivially small

resource over the surface of the structure, sets of subsequent trajectories are pulled apart, they stretched and folded, so that it becomes impossible to predict their exact population densities over the future. The strength of the mixing that gives rise to the extreme sensitivity to initial conditions can be measured mathematically by calculating the Lyapunov exponent, which is positive for the chaotic dynamics and nonpositive otherwise. There have been many attempts to measure Lyapunov exponents from time series data, and some chaotic chaotic populations have been identified. In some cases, analysis, and most convincingly, human, childhood diseases, but the statistical difficulties provide any broad generalization (5).

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

SCIENCE • VOL. 271 • 17 JANUARY 1993

303

- Our data are *Science* from 1880-2002, courtesy of JSTOR.
- We have 130K documents, 76M words.
- Goal: Discover a latent thematic structure in this corpus, useful for browsing, search, and similarity assessment.

# Topic models

- Use multinomial distributions over the vocabulary, called *topics*, to describe a collection of documents in a hierarchical model
- Treat documents as arising from a generative probabilistic process that includes hidden themes
- Discover those themes using *posterior inference*
- Useful for many kinds of tasks
  - Organization
  - Classification
  - Collaborative filtering
  - Information retrieval

- Latent Dirichlet allocation
- Dynamic Topic Models
- Correlated Topic Models

## 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.

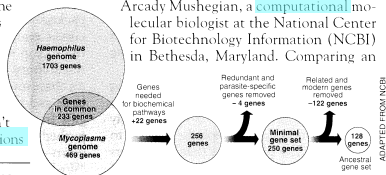
Although the numbers don't match precisely, those predictions

\* 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.

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

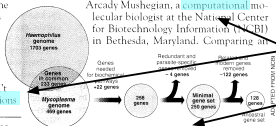
## Seeking Life's Bare (Genetic) Necessities

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.

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Anderson, a biologist at the 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

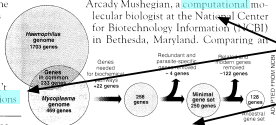
## Seeking Life's Bare (Genetic) Necessities

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.

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Anderson, a biologist at the University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **scientific numbers game**, particularly as more and more **genomes** are completely sequenced 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

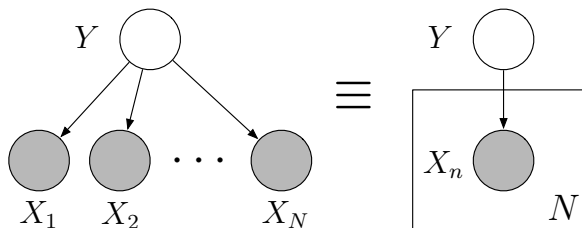


Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

- 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?

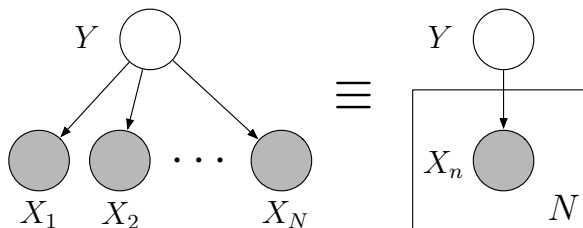
## Graphical models (Aside)



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



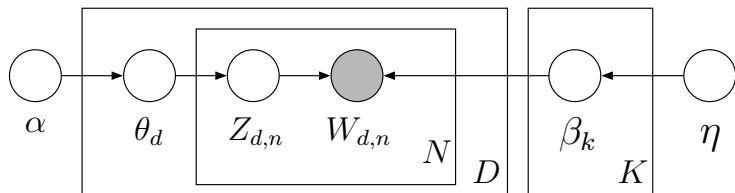
## Graphical models (Aside)



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

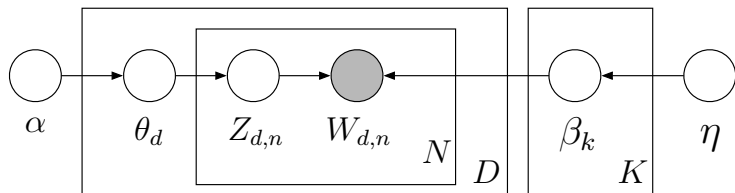
$$p(y, x_1, \dots, x_N) = p(y) \prod_{n=1}^N p(x_n | y)$$

# Latent Dirichlet allocation



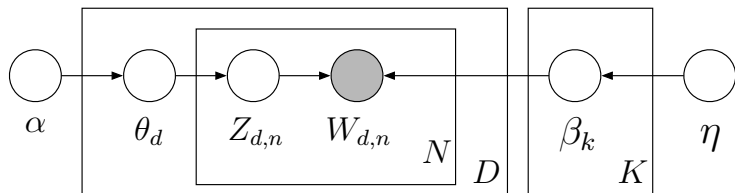
- 1 Draw each topic  $\beta_i \sim \text{Dir}(\eta)$ , for  $i \in \{1, \dots, K\}$ .
- 2 For each document:
  - 1 Draw topic proportions  $\theta_d \sim \text{Dir}(\alpha)$ .
  - 2 For each word:
    - 1 Draw  $Z_{d,n} \sim \text{Mult}(\theta_d)$ .
    - 2 Draw  $W_{d,n} \sim \text{Mult}(\beta_{z_{d,n}})$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

# Latent Dirichlet allocation



- From a collection of documents, infer
  - Per-word topic assignment  $z_{d,n}$
  - Per-document topic proportions  $\theta_d$
  - 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, but we can use:

- 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,<sup>6</sup> 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 Aracely Moshegun, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



<sup>6</sup> Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

**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.

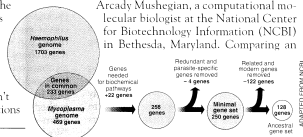
# Example inference

## 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

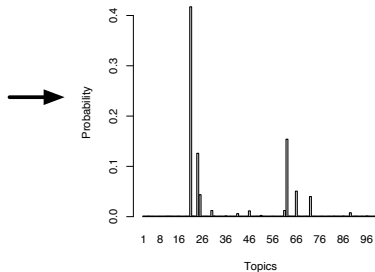
"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes 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



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

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

SCIENCE • VOL. 272 • 24 MAY 1996



# 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

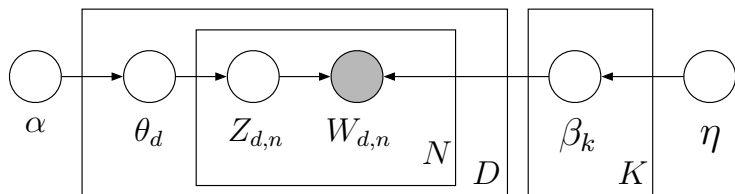
# Latent Dirichlet allocation

- 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).
  - 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>.



# Dynamic Topic Models

# LDA and exchangeability



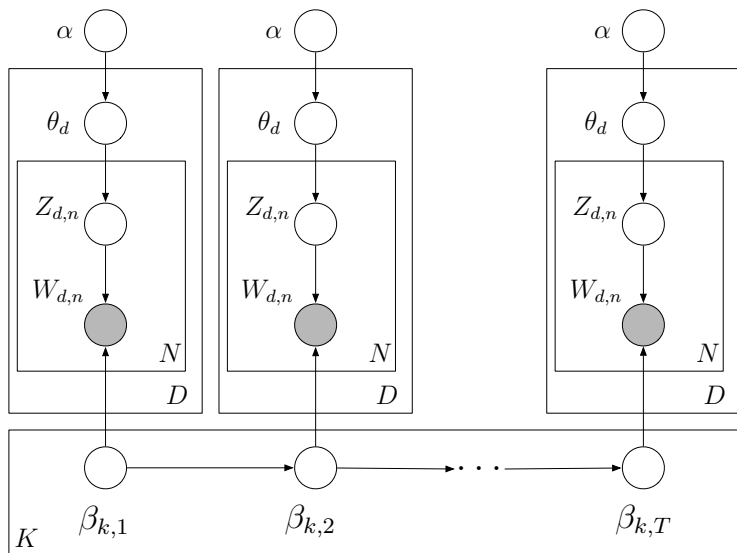
- LDA assumes that documents are exchangeable.
- I.e., their joint probability is invariant to permutation.
- This is too restrictive.



# 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



## Original article

## Topic proportions

TECHVIEW: DNA SEQUENCING

### Sequencing the Genome, Fast

James C. McElhinny and Amanda A. McHurry

Genome sequencing projects reveal the genetic makeup of an organism by reading off the sequence of the DNA bases, which encode all of the information necessary for the life of the organism. The base sequences contain four nucleotides—adenine, thymine, guanine, 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-to-base sequence of DNA easier. An 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, their excitation of a fluorescent dye specific to the base at the end of the molecule yields a base-specific signal that can be optically recorded.

The latest sequencer to be launched is Perkin-Elmer's next-generation ABI Prism 3700 DNA Analyzer, which, like the Molecular Dynamics Megalace CE 3500 launched last year, incorporates a capillary tube to hold the sequencer gel rather than a traditional slab-shaped gel. Specific to this is the ABI 3700 has been patented by inventor Craig Venter of Celera Genomics Corporation, and unlike the CE 3500, which produces one sequence for the entire 376-nucleotide (nt) length of the human genome in 2 hours, the significant feature of the ABI 3700 machine is that, with less than 1 hour of human labor per day, one sequencer can produce 148 samples per day. Assuming that each sample gives an average of 400 base pairs (bp) of usable sequence data (its read length) and a window from an average human genome is covered by an average of 18 overlapping independent reads (1), the 75 million samples that Celera Genomics will require (2) will require about 430 days, which affords some margin of error for unexpected developments.

At the Sequencer Centre, we have finished 146 Mb of genomic sequencing from a variety of genomes, including human, in a variety



Fig. 3. Comparison of read length histograms for the ABI 3700 and Megalace CE 3500. The Megalace CE 3500 produces one sequence for the entire 376-nt length of the human genome in 2 hours, while the ABI 3700 produces 148 samples per day. The ABI 3700 produces 148 samples per day, which affords some margin of error for unexpected developments.

of the Sequencer Centre in December 1999—also in our Research and Development department for evaluation. Thus, the ABI 3700 will ultimately be added to our present capacity to track our goal.

The ABI 3700 DNA sequencer is built into a four-lane gel cabinet, which contains in its base all the reagents required for its operation. The reagent containers are readily accessible for replenishment, which is required every day under high-throughput operation. At bench level, the user enters a 4-parameter batch, in which reagent quantities of DNA samples are located. The operator places the prepared plates in the sequencer, closes the front of the machine and programs it by using a personal computer. A robotic arm transfers DNA sam-

ples from the plates into wells that open up to the capillaries. This and the rest of the sequencing operation is fully automatic. The reagent can inventory process for 1000 plates of DNA samples unattended, making approximately 10 hours of bench operation unnecessary to repeat. This rate falls short of the average capabilities of four ABI 3700 plates in 32 hours.

The main innovation of the ABI 3700 is the use of a sheath flow fluorescence detection system (4). Instead of the DNA being sequenced, 2000s pass the end of the capillary within a sheath flow. A sheath fluid flows over the ends of the capillaries, drawing the DNA fragments in their center from the capillaries through a fluid layer that automatically eliminates all but the target DNA. The capillary fluorescence is detected with a special CCD (charge-coupled device) camera. This arrangement means that there are no moving parts in the detection system, other than a shutter to burst

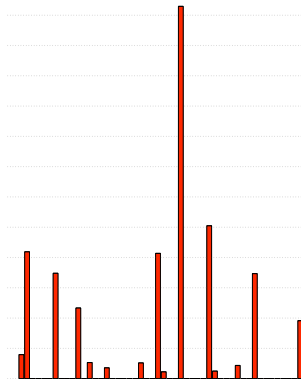
the CCD detector.

We have evaluated these machines for their performance, speed, and reliability in comparison to the more commonly used slab gel sequencing machines. In automated sequencing, there are two methods for counting the gel matrix. One is to perform a gel matrix between two finely separated glass slides (4, 5) or, less—the slab gel method. The other is to inject a polymer matrix into a capillary (internal diameter 100  $\mu$ m). Most sequencing facilities use the slab gel method, because multilane sequencing have only recently become commercially available.

With either type of system, the data is read as many bases possible for a given sample of DNA—that is, long read lengths are desirable. In fact, a system that could read longer is more than half the speed of another system is preferable, if both systems cost the same. This is because sequencing relatively long fragments requires more sequencing fragments to achieve the same coverage.

We have directly compared the ABI 3700 sequencer to the ABI 3770L slab gel sequencer by evaluating the sample-to-sample time with both machines with human DNA samples. These samples were also obtained from both machines with human DNA samples. These samples were also prepared and sequenced with our standard protocols for Perkin-Elmer Big Dye Terminator chemistry.

The authors are at the Sequencer Centre, Wellcome Trust Genome Centre, Hinxton, Cambs, CB3 0ET, UK. E-mail: jcm@jgphd.well.ac.uk



## Original article

## Most likely words from top topics

TECHNOLOGY

### Sequencing the Genome, Fast

James C. McElhin and Amanda A. McPherson

Genomic sequencing projects reveal the genetic makeup of an organism by reading off the sequence of the DNA bases, which encode all of the information necessary for the life of the organism. The base sequence contains four nucleotides—adenine, thymine, guanine, 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-to-base sequence of DNA easier. An 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 attached to the base at the end of the molecule yields a base-specific signal that can be electronically recorded.

The latest sequencer to be launched is Perkin-Elmer's next-generation ABI Prism 3700 DNA Analyzer, which like the Molecular Dynamics MegaBACE 3000 launched last year, incorporates a capillary tube to hold the sequencer gel rather than a traditional slab-shaped gel. Thanks to advances in the ABI 3700, its base-pairing accuracy is comparable to that of the original 3700. The MegaBACE 3000 produces 100 million reads per day, while the ABI 3700 machine says that, with less than 1 hour of hands-on time, it can sequence 148 samples per day. Assuming that each sample gives an average of 400 base pairs (bp) of usable sequence data (its read length) and that within that section human genome is covered by an average of 10 overlapping independent reads (1), the 75 million reads per day that Celera Genomics will require (100,000 ABI 3700 machines, days with ~250 machines, that works out to less than 2 years) can be accomplished.

At the Sanger Centre, we have finished 146 Mb of genomic sequencing from a water-



**Fig. 3.** Comparison of read length histograms for sequencing centers with the ABI 3700 capillary sequencer and the Sanger Centre. The Sanger Centre histogram shows the distribution of read lengths for the 100 million reads per day that the Sanger Centre produces. The ABI 3700 histogram shows the distribution of read lengths for the 100 million reads per day that the ABI 3700 produces. The Sanger Centre histogram shows a peak around 1000 bp, while the ABI 3700 histogram shows a peak around 300 bp.

to the Sanger Centre in December 1999—was in our Research and Development department for evaluation. Thus, the ABI 3700 will ultimately be added to our present capacity to track our goal.

The ABI 3700 DNA sequencer is built into a four-module cabinet, which incorporates in its base all the reagents required for its operation. The reagent containers are readily accessible to the technician, which is required every day for high-throughput operation. At Perkin-Elmer, we have completed a four-month test, in which we processed 400 samples on the machine and program it by using a personal computer. A robotic arm transfers DNA sam-

ples from the plates into wells that open to the capillaries. This, and the rest of the sequencing operation, is fully automatic. The reagent use economy is great. For the 3700 plates of DNA sample, an automated, taking approximately 10 hours to fill up, the reagents are required. This re-usable tube of the design simplification of the 3700 plates to 32 holes.

The main innovation of the ABI 3700 is the use of a sheath flow fluorescence detection system (4). Instead of the DNA fragments occupying 2000 pmol per well of the flow, they enter a sheath flow. A sheath flow flows over the ends of the capillaries, drawing the DNA fragments in the center from the capillaries through a sheath flow. This automatically eliminates well-to-well cross-talk. The capillary fluorescence is detected with a special CCD (charge-coupled device) camera. This arrangement means that there are no moving parts in the detection system, other than a shutter to burst

the CCD detector.

We have evaluated these machines for their performance, operation ease of use, and reliability in comparison to the more commonly used slab gel sequencing machines. In automated sequencing, there are two methods for creating the gel matrix. One is to perform a gel matrix between two finely separated glass slides (4, 5) and the other is to use a polymer matrix into a capillary (internal diameter 100  $\mu$ m). Most sequencing facilities use the slab gel method, because slab gel sequencing machines have only recently become commercially available.

With either type of system, the data is so small as many bases possible for a given sample of DNA—its long read lengths are

not available. In fact, a system that could read longer is more likely to be used for speed of analysis system is preferable, if the system cost the same. This is because sequencing relatively long fragments is an important parameter when evaluating new sequencing technologies.

We have directly compared the ABI 3700 sequencer to the ABI 3770X slab gel sequencer by evaluating the sample data obtained from both machines with human DNA samples. These samples were also obtained from phage and *Yersinia enterocolitica* samples and sequenced with our standard protocols for Perkin-Elmer Big Dye Terminator chemistry.

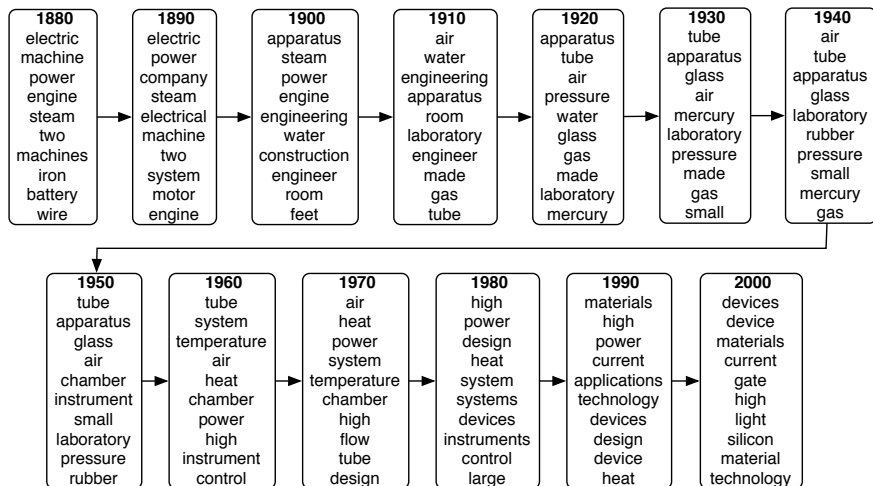
sequence  
genome  
genes  
sequences  
human  
gene  
dna  
sequencing  
chromosome  
regions  
analysis  
data  
genomic  
number

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 Centre, Hinxton, Cambs CB3 0ET, UK. E-mail: jcm@phd.sanger.ac.uk

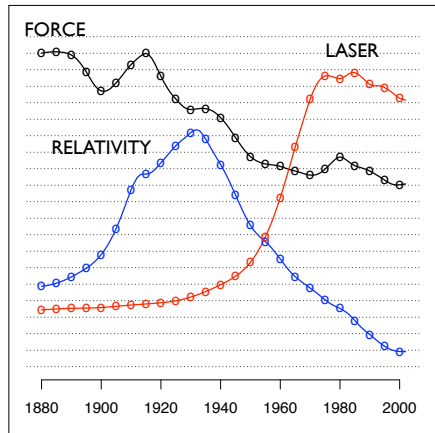
# Analyzing a topic



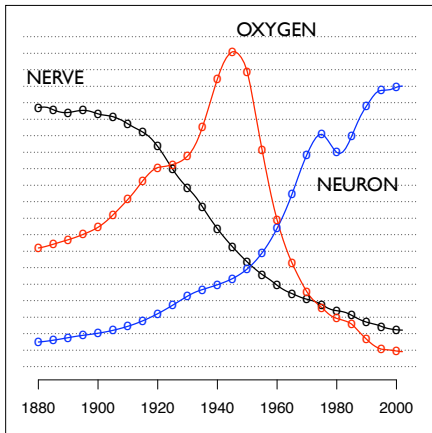


# Visualizing trends within a topic

## "Theoretical Physics"



## "Neuroscience"





# Time-corrected document similarity

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

plexus, the visuographic organization of the medial occipital-parietal cortex was compared with electrophysiological recording techniques in five owl monkeys (O). The monkeys were anesthetized with chloral hydrate and prepared for recording. Temporal and photostimulus microstimulation were used to record from small clusters of neurons or occasionally from single neurons in temporal projections parallel to the medial surface of occipital-parietal cortex. Receptive fields were plotted by using stimulus spots or scintigraphs and by using the surfaces of a translucent plastic hemisphere situated in front of the contralateral eye. The position of the optic disk was projected onto the plastic hemisphere with the method of Fernald and Chase (5). The (indicated) eye usually was

concentrated with an eyeglass shield. Electrode tracks and recording sites were reconstructed from histological sections and photographs of the intact brain.

Figure 1 illustrates the data from an almost complete mapping of the medial cortex obtained in the other four experiments. Through a comparison of the recorded surface of occipital-parietal cortex at a distance of approximately 1 mm from the medial surface (as previously published in part 1) we found that the receptive fields recorded adjacent to the medial area in the visual area (P 11) was located in the lower quadrant near the horizontal midline about 20° to 40° from the center (5). This, as is shown in Fig. 1, and

also in Fig. 2, which illustrates the organization of the other cortical visual areas, has not been mapped in the owl monkey. The border between the medial area and the central visual area corresponds to a complete portion of the horizontal midline. In other experiments in the dorso-lateral areas we found that receptive fields located near the anterior border with the medial area began near the vertical midline as in the lower quadrant and proceeded to a broad base in the periphery toward the horizontal meridian (5). Thus, as is shown in Figs. 1 and 2, the common border between the dorso-lateral 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 by pro-

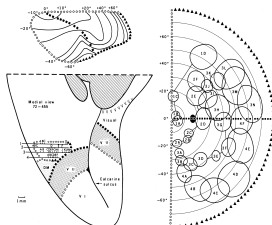
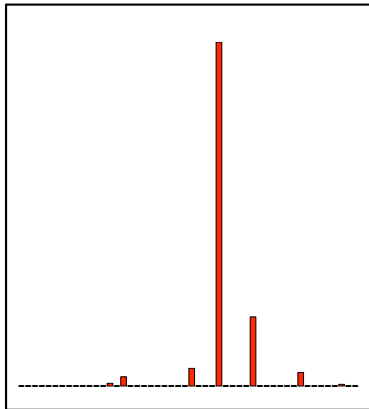


Fig. 1. Medial wall recording penetrations and receptive fields for the medial visual area in owl monkey 71,495. The diagram on the left is a view of the posterior half of the medial wall of cerebral cortex of the left hemisphere with the brain on and cerebellum removed. Arrows are used to indicate the site of the electrode. Microphotograph penetrations are represented, and recording sites are indicated by their case numbers (O). The corresponding receptive fields are shown in the same position on the right. In the upper half is an expanded view of the visuographic organization of the medial area (P 11) in the lower quadrant near the horizontal midline. The receptive fields are indicated by their case numbers (O). P 11 is the horizontal area, P 12 is the dorsal area, P 13 is the dorsal area, P 14 is the dorsal area, P 15 is the dorsal area, P 16 is the dorsal area, P 17 is the dorsal area, P 18 is the dorsal area, P 19 is the dorsal area, P 20 is the dorsal area, P 21 is the dorsal area, P 22 is the dorsal area, P 23 is the dorsal area, P 24 is the dorsal area, P 25 is the dorsal area, P 26 is the dorsal area, P 27 is the dorsal area, P 28 is the dorsal area, P 29 is the dorsal area, P 30 is the dorsal area, P 31 is the dorsal area, P 32 is the dorsal area, P 33 is the dorsal area, P 34 is the dorsal area, P 35 is the dorsal area, P 36 is the dorsal area, P 37 is the dorsal area, P 38 is the dorsal area, P 39 is the dorsal area, P 40 is the dorsal area, P 41 is the dorsal area, P 42 is the dorsal area, P 43 is the dorsal area, P 44 is the dorsal area, P 45 is the dorsal area, P 46 is the dorsal area, P 47 is the dorsal area, P 48 is the dorsal area, P 49 is the dorsal area, P 50 is the dorsal area, P 51 is the dorsal area, P 52 is the dorsal area, P 53 is the dorsal area, P 54 is the dorsal area, P 55 is the dorsal area, P 56 is the dorsal area, P 57 is the dorsal area, P 58 is the dorsal area, P 59 is the dorsal area, P 60 is the dorsal area, P 61 is the dorsal area, P 62 is the dorsal area, P 63 is the dorsal area, P 64 is the dorsal area, P 65 is the dorsal area, P 66 is the dorsal area, P 67 is the dorsal area, P 68 is the dorsal area, P 69 is the dorsal area, P 70 is the dorsal area, P 71 is the dorsal area, P 72 is the dorsal area, P 73 is the dorsal area, P 74 is the dorsal area, P 75 is the dorsal area, P 76 is the dorsal area, P 77 is the dorsal area, P 78 is the dorsal area, P 79 is the dorsal area, P 80 is the dorsal area, P 81 is the dorsal area, P 82 is the dorsal area, P 83 is the dorsal area, P 84 is the dorsal area, P 85 is the dorsal area, P 86 is the dorsal area, P 87 is the dorsal area, P 88 is the dorsal area, P 89 is the dorsal area, P 90 is the dorsal area, P 91 is the dorsal area, P 92 is the dorsal area, P 93 is the dorsal area, P 94 is the dorsal area, P 95 is the dorsal area, P 96 is the dorsal area, P 97 is the dorsal area, P 98 is the dorsal area, P 99 is the dorsal area, P 100 is the dorsal area.

19 FEBRUARY 1976

111



# Browser of Science

## 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 mining 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 model, all information as well as information presented by sets, or, if typically a word, associated with that information. In principle chosen from a controlled thesaurus, but being constructed for unrestricted topics to derive the terms under consideration to a text content.

Because the terms for content representation introduce a term-weighting scheme, and lower weights to A powerful term-weighting scheme is the well-known term frequency-inverse document frequency (TF-IDF), which is a function of the term frequency ( $f_{ij}$ ) in document  $j$  with a low frequency ( $f_{.j}$ ). Such terms are distinguished by their occurrence.

When all texts are represented by weighted vectors  $D_j = (d_{1j}, d_{2j}, \dots, d_{nj})$ , the weight assigned to each term in the similarity measure between pairs of vectors is  $\cos(\theta_{ij})$ . Thus, given

SCIENCE • VOL. 271

## "Automatic Analysis, Theme Generation, and Summarization of Machine-Readable Texts" (1994)

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.

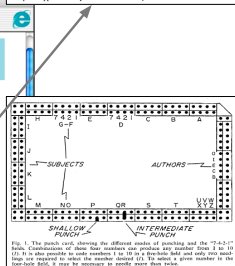


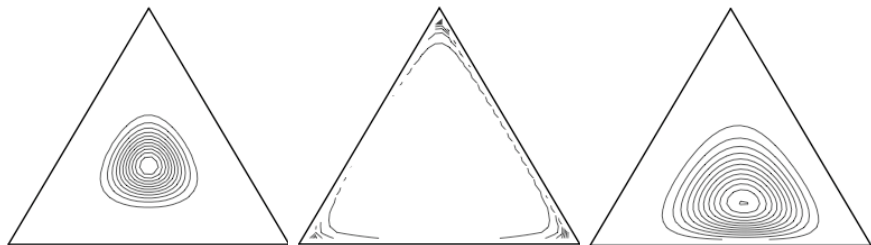
Fig. 1. The punch card, showing the different codes of punching and the "2-3-1" grids. Combinations of these three numbers can produce any number from 1 to 10 (1). It is also possible to code numbers 1 to 10 in a 5-hole field and only five punch-cards are required to select the number desired (2). To select a given number in the 10-hole field, it may be necessary to punch more than twice.

### THE STORING OF PAMPHLETS.

On reading Professor Minot's explanation of his method of storing pamphlets as given in the issue of December 9th 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.

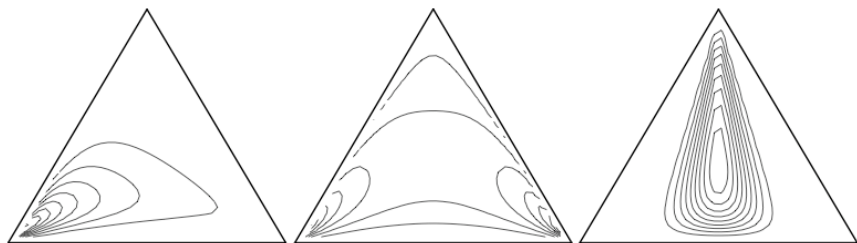
# Correlated Topic Models

# 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.

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

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





# Summary

- Topic models provide useful descriptive statistics for understanding the latent thematic structure of text data.
- But, models come with hidden assumptions, e.g.,
  - Exchangeability
  - Component-wise independence
- Current research
  - Choosing the number of topics
  - Continuous time dynamic topic models
  - Topic models for prediction
  - Inferring the impact of a document
- Download code and papers at <http://www.cs.princeton.edu/~blei>.

“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)