# Probabilistic Topic Models

David M. Blei

Department of Computer Science
Princeton University
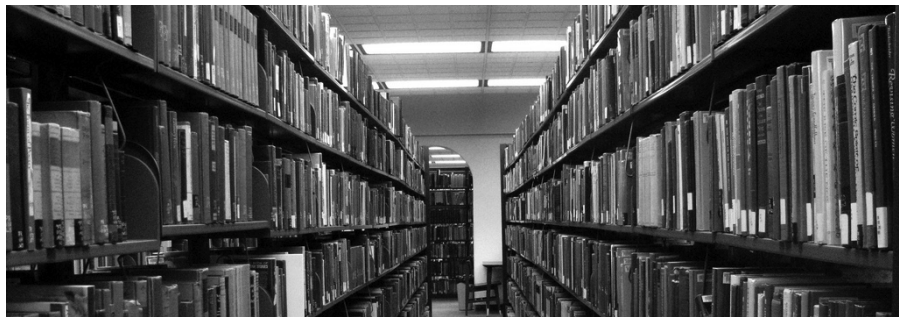
August 22, 2011

## Information overload



As more information becomes available, it becomes more difficult to find and discover what we need.

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

# Topic modeling



Topic modeling provides methods for automatically organizing, understanding, searching, and summarizing large electronic archives.
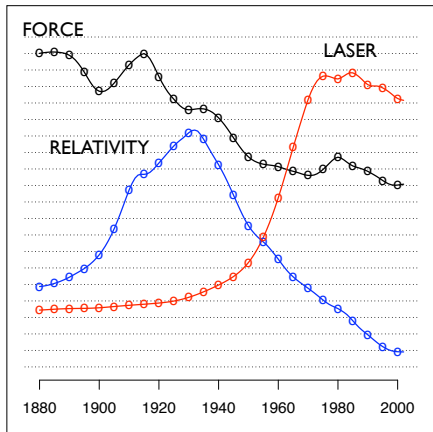
1. Discover the hidden themes that pervade the collection.
2. Annotate the documents according to those themes.
3. Use annotations to organize, summarize, and search the texts.
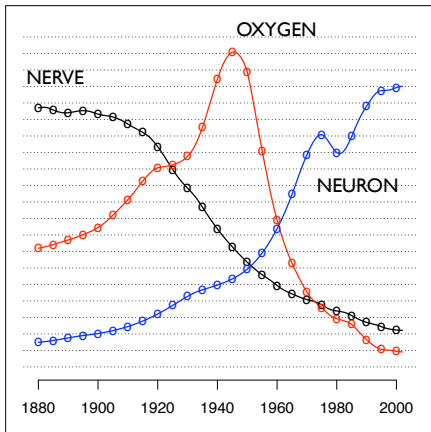
# Discover topics from a corpus

| 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

# Model connections between topics

# Find hierarchies of topics



Quantum lower bounds by polynomials
On the power of bounded concurrency I: finite automata
Dense quantum coding and quantum finite automata
Classical physics and the Church–Turing Thesis

quantum
automata
nc
automaton
languages

online
scheduling
task
competitive
tasks

approximation
s
points
distance
convex

n
functions
polynomial
log
algorithm

routing
adaptive
network
networks
protocols

Nearly optimal algorithms and bounds for multilayer channel routing
How bad is selfish routing?
Authoritative sources in a hyperlinked environment
Balanced sequences and optimal routing

machine
domain
degree
degrees
polynomials

networks
protocol
network
packets
link

learning
learnable
statistical
examples
classes

graph
graphs
edge
minimum
vertices

An optimal algorithm for intersecting line segments in the plane
Recontamination does not help to search a graph
A new approach to the maximum-flow problem
The time complexity of maximum matching by simulated annealing

the,of
a, is
and

constraint
dependencies
local
consistency
tractable

Module algebra
On XML integrity constraints in the presence of DTDs
Closure properties of constraints
Dynamic functional dependencies and database aging

database
constraints
algebra
boolean
relational

logic
logics
query
theories
languages

m
merging
networks
sorting
multiplication

n
algorithm
time
log
bound

consensus
objects
messages
protocol
asynchronous

system
systems
performance
analysis
distributed

learning
knowledge
reasoning
verification
circuit

trees
regular
tree
search
compression

logic
programs
systems
language
sets

networks
queuing
asymptotic
productform
server

Single-class bounds of multi-class queuing networks
The maximum concurrent flow problem
Contention in shared memory algorithms
Linear probing with a nonuniform address distribution

database
transactions
retrieval
concurrency
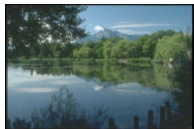restrictions

Magic Functions: In Memoriam: Bernard M. Dwork 1923–1998
A mechanical proof of the Church-Rosser theorem
Timed regular expressions
On the power and limitations of strictness analysis

proof
property
program
resolution
abstract

formulas
firstorder
decision
temporal
queries

# Annotate images



SKY WATER TREE MOUNTAIN PEOPLE



SCOTLAND WATER FLOWER HILLS TREE



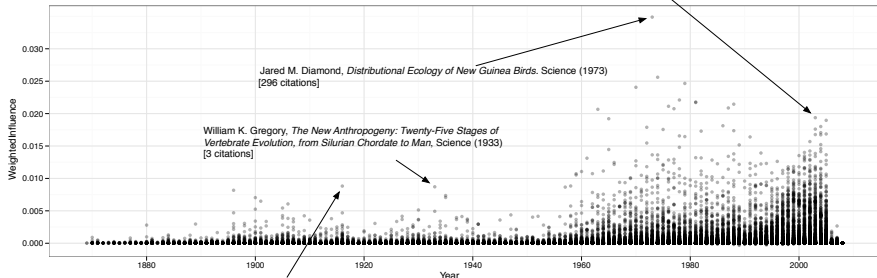SKY WATER BUILDING PEOPLE WATER



FISH WATER OCEAN TREE CORAL



PEOPLE MARKET PATTERN TEXTILE DISPLAY



BIRDS NEST TREE BRANCH LEAVES

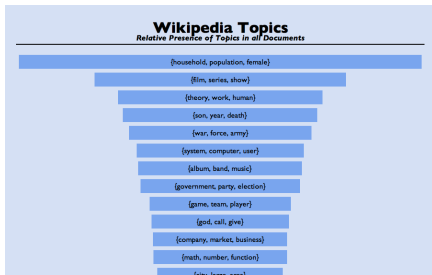# Discover influential articles

# Predict links between articles

| | |
|---|---|
| *Markov chain Monte Carlo convergence diagnostics: A comparative review* | |
| **Minorization conditions and convergence rates for Markov chain Monte Carlo**<br>Rates of convergence of the Hastings and Metropolis algorithms<br>**Possible biases induced by MCMC convergence diagnostics**<br>Bounding convergence time of the Gibbs sampler in Bayesian image restoration<br>Self regenerative Markov chain Monte Carlo<br>Auxiliary variable methods for Markov chain Monte Carlo with applications<br>**Rate of Convergence of the Gibbs Sampler by Gaussian Approximation**<br>Diagnosing convergence of Markov chain Monte Carlo algorithms | RTM ($\psi_e$) |
| Exact Bound for the Convergence of Metropolis Chains<br>Self regenerative Markov chain Monte Carlo<br>**Minorization conditions and convergence rates for Markov chain Monte Carlo**<br>Gibbs-markov models<br>Auxiliary variable methods for Markov chain Monte Carlo with applications<br>Markov Chain Monte Carlo Model Determination for Hierarchical and Graphical Models<br>Mediating instrumental variables<br>A qualitative framework for probabilistic inference<br>Adaptation for Self Regenerative MCMC | LDA + Regression |

# Characterize political decisions



tax credit,budget authority,energy,outlays,tax
county,eligible,ballot,election,jurisdiction
bank,transfer,requires,holding company,industrial
housing,mortgage,loan,family,recipient
energy,fuel,standard,administrator,lamp
student,loan,institution,lender,school
medicare,medicaid,child,chip,coverage
defense,iraq,transfer,expense,chapter
business,administrator,bills,business concern,loan
transportation,rail,railroad,passenger,homeland security
cover,bills,bridge,transaction,following
bills,tax,subparagraph,loss,taxable
loss,crop,producer,agriculture,trade
head,start,child,technology,award
computer,alien,bills,user,collection
science,director,technology,mathematics,bills
coast guard,vessel,space,administrator,requires
child,center,poison,victim,abuse
land,site,bills,interior,river
energy,bills,price,commodity,market
surveillance,director,court,electronic,flood
child,fire,attorney,internet,bills
drug,pediatric,product,device,medical
human,vietnam,united nations,call,people
bills,iran,official,company,sudan
coin,inspector,designee,automobile,lebanon
producer,eligible,crop,farm,subparagraph
people,woman,american,nation,school
veteran,veterans,bills,care,injury
dod,defense,defense and appropriation,military,subtitle

# Organize and browse large corpora

## Wikipedia Topics
*Relative Presence of Topics in all Documents*

- {household, population, female}
- {film, series, show}
- {theory, work, human}
- {son, year, death}
- {war, force, army}
- {system, computer, user}
- {album, band, music}
- {government, party, election}
- {game, team, player}
- {god, call, give}
- {company, market, business}
- {math, number, function}
- {city, large, area}

## {film, series, show}

| words | related documents | related topics |
|---|---|---|
| film | The X-Files | {son, year, death} |
| series | Orson Welles | {work, book, publish} |
| show | Stanley Kubrick | {album, band, music} |
| character | B movie | {woman, child, man} |
| play | Mystery Science Theater 3000 | {law, state, case} |
| make | Monty Python | {black, white, people} |
| episode | Doctor Who | {theory, work, human} |
| movie | Sam Peckinpah | {@card@, make, design} |
| good | Married… with Children | {war, force, army} |
| release | History of film | {god, call, give} |
| feature | The A-Team | {game, team, player} |
| television | Pulp Fiction (film) | {day, year, event} |
| star | Mad (magazine) | {company, market, business} |

## Stanley Kubrick



**Stanley Kubrick** (July 26, 1928 – March 7, 1999) was an American film director, writer, producer, and photographer who lived in England during most of the last four decades of his career. Kubrick was noted for the scrupulous care with which he chose his subjects, his slow method of working, the variety of genres he worked in, his technical perfectionism, and his reclusiveness about his films and personal life. He worked far beyond the confines of the Hollywood system, maintaining almost complete artistic control and making movies according to his own whims and time constraints, but with the rare advantage of big-studio financial support for all his endeavors.

Kubrick's films are characterized by a formal visual style and meticulous attention to detail—his later films often have elements of surrealism and expressionism that eschews structured linear narrative. His films are repeatedly described as slow and methodical, and are often perceived as a reflection of his obsessive and perfectionist nature.[1] A recurring theme in his films is man's inhumanity to man. While often viewed as

### related topics
- {film, series, show}
- {theory, work, human}
- {son, year, death}
- {black, white, people}
- {god, call, give}
- {math, energy, light}

### related documents
- Orson Welles
- B movie
- Mystery Science Theater 3000
- Monty Python
- Doctor Who
- Sam Peckinpah
- The A-Team
- Pulp Fiction (film)
- Buffy the Vampire Slayer (TV series)
- The X-Files
- Sunset Boulevard (film)
- Jack Benny

## {theory, work, human}

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

# This tutorial

- What are topic models?
- What kinds of things can they do?
- How do I compute with a topic model?
- What are some unsanswered questions in this field?
- How can I learn more?

## Related subjects

Topic modeling is a case study in modern machine learning with probabilistic models. It touches on

- Directed graphical models
- Conjugate priors and nonconjugate priors
- Time series modeling
- Modeling with graphs
- Hierarchical Bayesian methods
- Approximate posterior inference (MCMC, variational methods)
- Exploratory and descriptive data analysis
- Model selection and Bayesian nonparametric methods
- Mixed membership models
- Prediction from sparse and noisy inputs

# If you remember one picture...



**Assumptions**

**Data**

**Inference algorithm**

**Discovered structure**

# Organization

- **Introduction to topic modeling**
  - Latent Dirichlet allocation
  - Open source implementations and tools

- **Beyond latent Dirichlet allocation**
  - Modeling richer assumptions
  - Supervised topic modeling
  - Bayesian nonparametric topic modeling

- **Algorithms**
  - Gibbs sampling
  - Variational inference
  - Online variational inference

- **Discussion, open questions, and resources**

# Introduction to Topic Modeling

# Probabilistic modeling

**1** Data are assumed to be observed from a generative probabilistic process that includes hidden variables.

- *In text, the hidden variables are the thematic structure.*

**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 a new document fit into the topic structure?*

# Latent Dirichlet allocation (LDA)



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

COLD SPRING HARBOR, NEW YORK— How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an
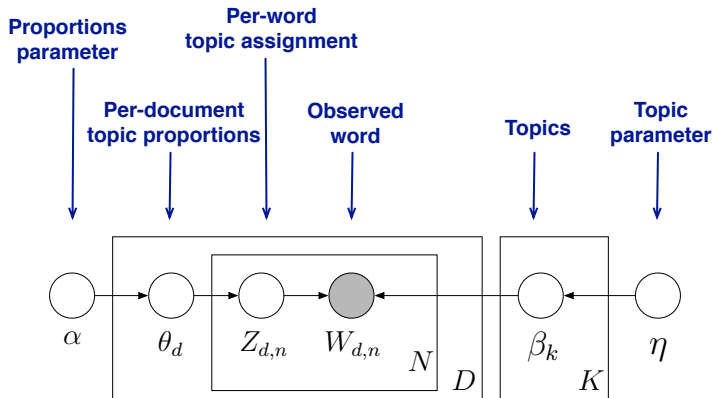
*Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

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

SCIENCE • VOL. 272 • 24 MAY 1996

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

# Generative model for LDA



*Topics*  *Documents*  *Topic proportions and assignments*

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

# The posterior distribution



*Topics*  *Documents*  *Topic proportions and assignments*

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

# The posterior distribution



*Topics*  *Documents*  *Topic proportions and assignments*

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

$$p(\text{topics, proportions, assignments} \mid \text{documents})$$

## LDA as a graphical model



**Proportions parameter** — $\alpha$
**Per-document topic proportions** — $\theta_d$
**Per-word topic assignment** — $Z_{d,n}$
**Observed word** — $W_{d,n}$
**Topics** — $\beta_k$
**Topic parameter** — $\eta$

- Encodes our assumptions about the data
- Connects to algorithms for computing with data
- See *Pattern Recognition and Machine Learning* (Bishop, 2006).

## LDA as a graphical model



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

# LDA as a graphical model



$$\prod_{i=1}^{K} p(\beta_i \mid \eta) \prod_{d=1}^{D} p(\theta_d \mid \alpha) \left( \prod_{n=1}^{N} p(z_{d,n} \mid \theta_d) p(w_{d,n} \mid \beta_{1:K}, z_{d,n}) \right)$$

# LDA



- This joint defines a posterior.

- From a collection of documents, infer
  - Per-word topic assignment $z_{d,n}$
  - Per-document topic proportions $\theta_d$
  - Per-corpus topic distributions $\beta_k$

- Then use posterior expectations to perform the task at hand, e.g., information retrieval, document similarity, exploration, ...

# LDA



Approximate posterior inference algorithms

- 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)
- Online variational inference (Hoffman et al., 2010)

Also see Mukherjee and Blei (2009) and Asuncion et al. (2009).

# Example inference



- **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

* 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.

*Haemophilus genome 1703 genes*
*Genes in common 233 genes*
*Mycoplasma genome 469 genes*

Genes needed for biochemical pathways ~22 genes → *256 genes*

Redundant and parasite-specific genes removed ~4 genes → *Minimal gene set 250 genes*

Related and modern genes removed ~122 genes → *128 genes Ancestral gene set*

ADAPTED FROM NCBI

# Example inference

| | | | |
|---|---|---|---|
| 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 |

# Example inference (II)

## Chaotic Beetles

### Charles Godfray and Michael Hassell

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

The authors are in the Department of Biology, Imperial College at Silwood Park, Ascot, Berks, SL5 7PZ UK. E-mail: m.hassell@ic.ac.uk

convincing evidence to date of complex dynamics and chaos in a biological population—of the flour beetle, *Tribolium castaneum* (see figure).

It has proven extremely difficult to demonstrate complex dynamics in populations in the field. By its very nature, a chaotically fluctuating population will superficially resemble a stable or cyclic population buffeted by the normal random perturbations experienced by all species. Given a long enough time series, diagnostic tools from nonlinear mathematics can be used to identify the telltale signatures of chaos. In phase space, chaotic trajectories come to lie on "strange attractors," curious geometric objects with fractal structure and hence noninteger dimension. As they



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

move over the surface of the attractor, sets of adjacent trajectories are pulled apart, then stretched and folded, so that it becomes impossible to predict exact population densities into the future. The strength of the mixing that gives rise to the extreme sensitivity to initial conditions can be measured mathematically by estimating the Liapunov exponent, which is positive for chaotic dynamics and nonpositive otherwise. There have been many attempts to estimate attractor dimension and Liapunov exponents from time series data, and some candidate chaotic population have been identified (some insects, rodents, and most convincingly, human childhood diseases), but the statistical difficulties preclude any broad generalization (3).

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

# Example inference (II)

| | | | |
|---|---|---|---|
| problem | model | selection | species |
| problems | rate | male | forest |
| mathematical | constant | males | ecology |
| number | distribution | females | fish |
| new | time | sex | ecological |
| mathematics | number | species | conservation |
| university | size | female | diversity |
| two | values | evolution | population |
| first | value | populations | natural |
| numbers | average | population | ecosystems |
| work | rates | sexual | populations |
| time | data | behavior | endangered |
| mathematicians | density | evolutionary | tropical |
| chaos | measured | genetic | forests |
| chaotic | models | reproductive | ecosystem |

# Held out perplexity



Nematode abstracts

Associated Press

Legend:
- Smoothed Unigram
- Smoothed Mixt. Unigrams
- LDA
- Fold in pLSI

$$\text{perplexity} = \exp\left\{\frac{-\sum_d \log p(\mathbf{w}_d)}{\sum_d N_d}\right\}$$

# Used to explore and browse document collections

## Aside: The Dirichlet distribution

- The Dirichlet distribution is an exponential family distribution over the simplex, i.e., positive vectors that sum to one

$$p(\theta \mid \vec{\alpha}) = \frac{\Gamma\left(\sum_i \alpha_i\right)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta_i^{\alpha_i - 1}.$$

- It is **conjugate** to the multinomial. Given a multinomial observation, the posterior distribution of $\theta$ is a Dirichlet.

- The parameter $\alpha$ controls the mean shape and sparsity of $\theta$.

- The topic proportions are a $K$ dimensional Dirichlet.
  The topics are a $V$ dimensional Dirichlet.

$\alpha = 10$

$\alpha = 100$

$\alpha = 0.001$

## Why does LDA "work"?

Why does the LDA posterior put "topical" words together?

- Word probabilities are maximized by dividing the words among the topics. (More terms means more mass to be spread around.)

- In a mixture, this is enough to find clusters of co-occurring words.

- In LDA, the Dirichlet on the topic proportions can encourage sparsity, i.e., a document is penalized for using many topics.

- Loosely, this can be thought of as softening the strict definition of "co-occurrence" in a mixture model.

- This flexibility leads to sets of terms that more tightly co-occur.

# Summary of LDA



- LDA can
    - visualize the hidden thematic structure in large corpora
    - generalize new data to fit into that structure

- Builds on Deerwester et al. (1990) and Hofmann (1999)
  It is a *mixed membership model* (Erosheva, 2004).
  Relates to *multinomial PCA* (Jakulin and Buntine, 2002)

- Was independently invented for genetics (Pritchard et al., 2000)

## Implementations of LDA

There are many available implementations of topic modeling—

| | |
|---|---|
| **LDA-C**$^*$ | A C implementation of LDA |
| **HDP**$^*$ | A C implementation of the HDP ("infinite LDA") |
| **Online LDA**$^*$ | A python package for LDA on massive data |
| **LDA in R**$^*$ | Package in R for many topic models |
| **LingPipe** | Java toolkit for NLP and computational linguistics |
| **Mallet** | Java toolkit for statistical NLP |
| **TMVE**$^*$ | A python package to build browsers from topic models |

$^*$ available at www.cs.princeton.edu/~blei/

## Example: LDA in R (Jonathan Chang)

perspective identifying tumor suppressor genes in human...
letters global warming report leslie roberts article global....
research news a small revolution gets under way the 1990s....
a continuing series the reign of trial and error draws to a close...
making deep earthquakes in the laboratory lab experimenters...
quick fix for freeways thanks to a team of fast working...
feathers fly in grouse population dispute researchers...

....

245 1897:1 1467:1 1351:1 731:2 800:5 682:1 315:6 3668:1 14:1
260 4261:2 518:1 271:6 2734:1 2662:1 2432:1 683:2 1631:7
279 2724:1 107:3 518:1 141:3 3208:1 32:1 2444:1 182:1 250:1
266 2552:1 1993:1 116:1 539:1 1630:1 855:1 1422:1 182:3 2432:1
233 1372:1 1351:1 261:1 501:1 1938:1 32:1 14:1 4067:1 98:2
148 4384:1 1339:1 32:1 4107:1 2300:1 229:1 529:1 521:1 2231:1
193 569:1 3617:1 3781:2 14:1 98:1 3596:1 3037:1 1482:12 665:2

....

```
docs <- read.documents("mult.dat")
K <- 20
alpha <- 1/20
eta <- 0.001
model <- lda.collapsed.gibbs.sampler(documents, K, vocab, 1000, alpha, eta)
```

| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| dna | protein | water | says | mantle |
| gene | cell | climate | researchers | high |
| sequence | cells | atmospheric | new | earth |
| genes | proteins | temperature | university | pressure |
| sequences | receptor | global | just | seismic |
| human | fig | surface | science | crust |
| genome | binding | ocean | like | temperature |
| genetic | activity | carbon | work | earths |
| analysis | activation | atmosphere | first | lower |
| two | kinase | changes | years | earthquakes |

| 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|
| end | time | materials | dna | disease |
| article | data | surface | rna | cancer |
| start | two | high | transcription | patients |
| science | model | structure | protein | human |
| readers | fig | temperature | site | gene |
| service | system | molecules | binding | medical |
| news | number | chemical | sequence | studies |
| card | different | molecular | proteins | drug |
| circle | results | fig | specific | normal |
| letters | ind | university | sequences | drugs |

| 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|
| years | species | protein | cells | space |
| million | evolution | structure | cell | solar |
| ago | population | proteins | virus | observations |
| age | evolutionary | two | hiv | earth |
| university | university | amino | infection | stars |
| north | populations | binding | immune | university |
| early | natural | acid | human | mass |
| fig | studies | residues | antigen | sun |
| evidence | genetic | molecular | infected | astronomers |
| record | biology | structural | viral | telescope |

| 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|
| fax | cells | energy | research | neurons |
| manager | cell | electron | science | brain |
| science | gene | state | national | cells |
| aaas | genes | light | scientific | activity |
| advertising | expression | quantum | scientists | fig |
| sales | development | physics | new | channels |
| member | mutant | electrons | states | university |
| recruitment | mice | high | university | cortex |
| associate | fig | laser | united | neuronal |
| washington | biology | magnetic | health | visual |

# Open source document browser (with Allison Chaney)

## Wikipedia Topics
*Relative Presence of Topics in all Documents*

{household, population, female}
{film, series, show}
{theory, work, human}
{son, year, death}
{war, force, army}
{system, computer, user}
{album, band, music}
{government, party, election}
{game, team, player}
{god, call, give}
{company, market, business}
{math, number, function}
{city, large, area}

## {film, series, show}

| words | related documents | related topics |
|---|---|---|
| film | The X-Files | {son, year, death} |
| series | Orson Welles | {work, book, publish} |
| show | Stanley Kubrick | {album, band, music} |
| character | B movie | {woman, child, man} |
| play | Mystery Science Theater 3000 | {law, state, case} |
| make | Monty Python | {black, white, people} |
| episode | Doctor Who | {theory, work, human} |
| movie | Sam Peckinpah | {@card@, make, design} |
| good | Married... with Children | {war, force, army} |
| release | History of film | {god, call, give} |
| feature | The A-Team | {game, team, player} |
| television | Pulp Fiction (film) | {day, year, event} |
| star | Mad (magazine) | {company, market, business} |

## Stanley Kubrick

**Stanley Kubrick** (July 26, 1928 – March 7, 1999) was an American film director, writer, producer, and photographer who lived in England during most of the last four decades of his career. Kubrick was noted for the scrupulous care with which he chose his subjects, his slow method of working, the variety of genres he worked in, his technical perfectionism, and his reclusiveness about his films and personal life. He worked far beyond the confines of the Hollywood system, maintaining almost complete artistic control and making movies according to his own whims and time constraints, but with the rare advantage of big-studio financial support for all his endeavors.

Kubrick's films are characterized by a formal visual style and meticulous attention to detail—his later films often have elements of surrealism and expressionism that eschews structured linear narrative. His films are repeatedly described as slow and methodical, and are often perceived as a reflection of his obsessive and perfectionist nature.[1] A recurring theme in his films is man's inhumanity to man. While often viewed as

**related documents**

Orson Welles
B movie
Mystery Science Theater 3000
Monty Python
Doctor Who
Sam Peckinpah
The A-Team
Pulp Fiction (film)
Buffy the Vampire Slayer (TV series)
The X-Files
Sunset Boulevard (film)
Jack Benny

**related topics**

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

## {theory, work, human}

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

# Why develop these kinds of models?



- Organizing and finding patterns in data has become important in the sciences, humanties, industry, and culture.

- LDA can be embedded in more complicated models that capture richer assumptions about the data.

- Algorithmic improvements let us fit models to massive data.

# Bigger Picture: Probabilistic modeling

**Assumptions**



**Data**



**Inference algorithm**

**Discovered structure**

- Research in modeling separates these basic activities
- Though linked, we can work on each piece separately

# Beyond Latent Dirichlet Allocation

**So far...**



- LDA is a simple topic model

- Can be used to find topics that describe a corpus

- Each document exhibits multiple topics

- How can we build on this simple model of text?

# LDA is extendible



- LDA can be **embedded in more complicated models**, embodying further intuitions about the structure of the texts.

- LDA models can include syntax, authorship, word sense, dynamics, correlation, hierarchies, ...

# LDA is extendible



- The **data generating distribution** can be changed.

- LDA models can be built for images, social networks, music, purchase histories, computer code, genetic data, click-through-data, neural spike trains, ...

# LDA is extendible



- The **LDA posterior** can be used in creative ways

- It can be used for information retrieval, collaborative filtering, document similarity, visualization, ...

## Beyond latent Dirichlet allocation

- Modeling richer assumptions
  - Correlated topic models
  - Dynamic topic models
  - Measuring scholarly impact
- Supervised topic models
  - Supervised LDA
  - Relational topic models
  - Ideal point topic models
- Bayesian nonparametric topic models

# Modeling richer assumptions

- Correlated topic models
- Dynamic topic models
- Measuring scholarly impact

# The hidden assumptions of the Dirichlet



- 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 (Aitchison, 1980).
- The natural parameters of the multinomial are drawn from a multivariate Gaussian distribution.

$$
\begin{aligned}
X &\sim \mathcal{N}_{K-1}(\mu, \Sigma) \\
\theta_i &\propto \exp\{x_i\}
\end{aligned}
$$

## The correlated topic model (CTM) (Blei and Lafferty, 2007)



Noconjugate prior on topic proportions

- 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

# Held out log probability in a CTM



- Analyzed held-out log probability on *Science*, 1960.
- CTM supports more topics and provides a better fit than LDA.

- activated tyrosine phosphorylation activation phosphorylation kinase
- p53 cell cycle activity cyclin regulation
- proteins protein binding domain domains
- brain memory subjects left task
- neurons stimulus motor visual cortical
- surface tip image sample device
- synapses ltp glutamate synaptic neurons
- rna dna rna polymerase cleavage site
- research funding support nih program
- science scientists says research people
- receptor receptors ligand ligands apoptosis
- amino acids cdna sequence isolated protein
- computer problem information computers problems
- materials organic polymer polymers molecules
- physicists particles physics particle experiment
- united states women universities students education
- wild type mutant mutations mutants mutation
- enzyme enzymes iron active site reduction
- sequence sequences genome dna sequencing
- surface liquid surfaces fluid model
- laser optical light electrons quantum
- reaction reactions molecule molecules transition state
- stars astronomers universe galaxies galaxy
- cells cell expression cell lines bone marrow
- plants plant gene genes arabidopsis
- magnetic magnetic field spin superconductivity superconducting
- bacteria bacterial host resistance parasite
- mice antigen t cells antigens immune response
- virus hiv aids infection viruses
- gene disease mutations families mutation
- development embryos drosophila genes expression
- fossil record birds fossils dinosaurs fossil
- pressure high pressure pressures core inner core
- mantle crust upper mantle meteorites ratios
- sun solar wind earth planets planet
- patients disease treatment drugs clinical
- cells proteins researchers protein found
- genetic population populations differences variation
- species forest forests populations ecosystems
- ancient found impact million years ago africa
- earthquake earthquakes fault images data
- co2 carbon carbon dioxide methane water
- ozone atmospheric measurements stratosphere concentrations
- volcanic deposits magma eruption volcanism
- climate ocean ice changes climate change

# Dynamic topic models (Blei and Lafferty, 2006)

**1789**



My fellow citizens: I stand here today humbled by the task before us, grateful for the trust you have bestowed, mindful of the sacrifices borne by our ancestors...

*Inaugural addresses*

**2009**



AMONG the vicissitudes incident to life no event could have filled me with greater anxieties than that of which the notification was transmitted by your order...

- LDA assumes that the order of documents does not matter.
- Not appropriate for corpora that span hundreds of years
- We may want to track how language changes over time.

# Dynamic topic models



Topics drifting in time

# Modeling evolving topics



- Use a logistic normal distribution to model topics evolving over time (Aitchison, 1980)

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

$$\beta_{t,k} \,|\, \beta_{t-1,k} \;\sim\; \mathcal{N}(\beta_{t-1,k}, I\sigma^2)$$
$$p(w \,|\, \beta_{t,k}) \;\propto\; \exp\left\{\beta_{t,k}\right\}$$

# Analyzing a document

## Original article



## Topic proportions

# Analyzing a document

**Original article**　　　　**Most likely words from top topics**



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

# Analyzing a topic



| 1880 | 1890 | 1900 | 1910 | 1920 | 1930 | 1940 |
|------|------|------|------|------|------|------|
| electric | electric | apparatus | air | apparatus | tube | air |
| machine | power | steam | water | tube | apparatus | tube |
| power | company | power | engineering | air | glass | apparatus |
| engine | steam | engine | apparatus | pressure | air | glass |
| steam | electrical | engineering | room | water | mercury | laboratory |
| two | machine | water | laboratory | glass | laboratory | rubber |
| machines | two | construction | engineer | gas | pressure | pressure |
| iron | system | engineer | made | made | made | small |
| battery | motor | room | gas | laboratory | gas | mercury |
| wire | engine | feet | tube | mercury | small | gas |

| 1950 | 1960 | 1970 | 1980 | 1990 | 2000 |
|------|------|------|------|------|------|
| tube | tube | air | high | materials | devices |
| apparatus | system | heat | power | high | device |
| glass | temperature | power | design | power | materials |
| air | air | system | heat | current | current |
| chamber | heat | temperature | system | applications | gate |
| instrument | chamber | chamber | systems | technology | high |
| small | power | high | devices | devices | light |
| laboratory | high | flow | instruments | design | silicon |
| pressure | instrument | tube | control | device | material |
| rubber | control | design | large | heat | technology |

# Visualizing trends within a topic

# Evaluating the DTM on all of *Science*



See the browser at http://topics.cs.princeton.edu/Science/

## Time-corrected document similarity

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

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

- 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)

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

## Measuring scholarly impact (Gerrish and Blei, 2009)



- Influential articles reflect future changes in language use.
- The "influence" of an article is a latent variable.
- Influential articles affect the drift of the topics that they discuss.
- The posterior gives a retrospective estimate of influential articles.

# Measuring scholarly impact

# Measuring scholarly impact



- Each document has an influence score $I_d$.

- Each topic drifts in a way that is biased towards the documents with high influence.

- The posterior of $I_{1:D}$ can be examined to retrospectively find articles that best explain future changes in language.

# Measuring scholarly impact



- This measure of impact only uses the words of the documents. It correlates strongly with citation counts.

- High impact, high citation: "The Mathematics of Statistical Machine Translation: Parameter Estimation" (Brown et al., 1993)

- "Low" impact, high citation: "Building a large annotated corpus of English: the Penn Treebank" (Marcus et al., 1993)

# Measuring scholarly impact at large scale

**(with S. Gerrish, A. Chaney and D. Mimno)**



- PNAS, *Science*, and *Nature* from 1880–2005
- 350,000 Articles
- 163M observations
- Year-corrected correlation is 0.166

| 1880 | 1890 | 1900 | 1910 | 1920 | 1930 | 1940 |
|------|------|------|------|------|------|------|
| species | species | species | species | species | species | species |
| evolution | genus | genus | genus | evolution | genus | genus |
| genera | evolution | genera | evolution | genus | evolution | evolution |
| origin | distribution | evolution | genera | genera | distribution | distribution |
| genus | genera | distribution | distribution | distribution | genera | genera |
| distribution | fossil | fossil | origin | fossil | fossil | mammals |
| darwin | origin | origin | fossil | origin | primitive | primitive |
| fossil | primitive | environment | primitive | primitive | mammals | origin |
| primitive | mammals | fauna | fauna | mammals | origin | fossil |
| families | fauna | primitive | environment | environment | environment | ecological |
| europe | environment | generic | generic | evolutionary | fauna | evolutionary |
| fauna | ancestors | mammals | mammals | generic | africa | relationships |
| ancestors | darwin | new_species | represented | fauna | modern | generic |
| mammals | extinct | nomenclature | ancient | families | fossils | taxonomic |
| africa | families | extinct | flora | pleistocene | ancient | modern |

| 1950 | 1960 | 1970 | 1980 | 1990 | 2000 |
|------|------|------|------|------|------|
| species | species | species | species | species | species |
| genus | evolution | evolution | evolution | evolution | supplementary |
| evolution | genus | evolutionary | evolutionary | evolutionary | evolution |
| distribution | evolutionary | populations | modern | diversity | materials |
| genera | distribution | ecological | populations | tree | evolutionary |
| modern | genera | relationships | patterns | patterns | e-mail |
| relationships | fossil | fossil | organisms | biological | diversity |
| fossil | relationships | modern | diversity | sites | file |
| mammals | modern | genus | biological | e-mail | biological |
| africa | populations | organisms | sites | relationships | patterns |
| african | environment | biological | ecological | origin | organisms |
| origin | mammals | patterns | fossil | populations | tree |
| evolutionary | primitive | distribution | relationships | organisms | relationships |
| ecological | areas | evolved | origin | phylogenetic | sites |
| biological | ecological | population | ecology | modern | populations |

## Summary: Modeling richer assumptions

- The Dirichlet assumptions on topics and topic proportions makes strong conditional independence assumptions about the data.

- The **correlated topic model** uses a logistic normal on the topic proportions to find patterns in how topics tend to co-occur.
  - See also Li and McCallum (2007) for another approach.
  - See http://www.cs.princeton.edu/∼blei/ for code.

- The **dynamic topic model** uses a logistic normal in a linear dynamic model to capture how topics change over time.
  - Documents can exhibit sequential structure.
  - Opens the door to a citation-free model of scholarly impact.
  - See also Wang and Blei (2010) for a continuous time variant

- What's the catch? The Dirichlet is easier to compute with than the logistic normal. (Stay tuned.)

# Supervised topic models

- Supervised LDA
- Relational topic models
- Ideal point topic models

## Supervised LDA (Blei and McAuliffe, 2007)

- LDA is an unsupervised model. How can we build a topic model that is good at the task we care about?

- Many data are paired with **response variables**.
  - User reviews paired with a number of stars
  - Web pages paired with a number of "likes"
  - Documents paired with links to other documents
  - Images paired with a category

- **Supervised topic models** are topic models of documents and responses, fit to find topics predictive of the response.

## Supervised LDA



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

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

## Supervised LDA



- The response variable *y* is drawn *after* the document because it depends on $z_{1:N}$, an assumption of **partial exchangeability**.

- Consequently, *y* is necessarily conditioned on the words.

- In a sense, this blends generative and discriminative modeling.

## Prediction

- Fit sLDA parameters to documents and responses. This gives:
    - topics $\beta_{1:K}$
    - coefficients $\eta_{1:K}$

- We have a new document $w_{1:N}$ with unknown response value.

- We predict $y$ using the SLDA expected value:

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

## Example: Movie reviews



| least | bad | more | awful | his | both |
|-------|-----|------|-------|-----|------|
| problem | guys | has | featuring | their | motion |
| unfortunately | watchable | than | routine | character | simple |
| supposed | its | films | dry | many | perfect |
| worse | not | director | offered | while | fascinating |
| flat | one | will | charlie | performance | power |
| dull | movie | characters | paris | between | complex |

-30    -20    -10    10    20

| have | not | one | however |
|------|-----|-----|---------|
| like | about | from | cinematography |
| you | movie | there | screenplay |
| was | all | which | performances |
| just | would | who | pictures |
| some | they | much | effective |
| out | its | what | picture |

- 10-topic sLDA model on movie reviews (Pang and Lee, 2005).

- Response: number of stars associated with each review

- Each component of coefficient vector $\eta$ is associated with a topic.

# Held out correlation

## Diverse response types with GLMs

- Want to work with response variables that don't live in the reals.
    - binary / multiclass classification
    - count data
    - waiting time

- Model the response response with a generalized linear model

$$p(y \mid \zeta, \delta) = h(y, \delta) \exp \left\{ \frac{\zeta y - A(\zeta)}{\delta} \right\} \ ,$$

where $\zeta = \eta^\top \bar{z}$.

- Complicates inference, but allows for flexible modeling.

# Image classification and annotation (Wang et al., 2009)



*highway*

car, sign, road

*inside city*

buildings, car, sidewalk

*street*

tree, car, sidewalk

*tall building*

trees, buildings
occluded, window

- Uses GLM sLDA for multiclass classification.

- Uses ideas from Blei and Jordan (2004) for annotation.

# Supervised topic models



- SLDA enables model-based regression where the predictor "variable" is a text document.

- It can easily be used wherever LDA is used in an unsupervised fashion (e.g., images, genes, music).

- SLDA is a supervised dimension-reduction technique, whereas LDA performs unsupervised dimension reduction.

# Relational topic models (Chang and Blei, 2010)



- Many data sets contain **connected observations**.

- For example:
  - Citation networks of documents
  - Hyperlinked networks of web-pages.
  - Friend-connected social network profiles

## Relational topic models (Chang and Blei, 2010)



- Research has focused on finding communities and patterns in the link-structure of these networks.

- We adapt sLDA to pairwise response variables.
  This adaptation leads to a model of **content and connection**.

- RTMs find related hidden structure in both types of data.

# Relational topic models



- Adapt fitting algorithm for sLDA with binary GLM response

- RTMs allow predictions about new and unlinked data. These predictions are out of reach for traditional network models.

# Predicting links from documents

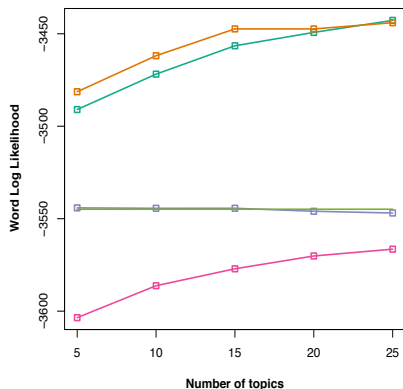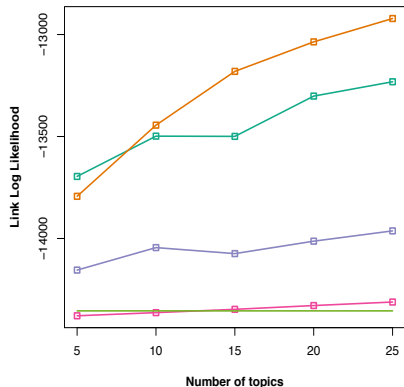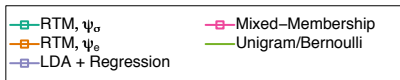| | |
|---|---|
| *Markov chain Monte Carlo convergence diagnostics: A comparative review* | |
| **Minorization conditions and convergence rates for Markov chain Monte Carlo** <br> Rates of convergence of the Hastings and Metropolis algorithms <br> **Possible biases induced by MCMC convergence diagnostics** <br> Bounding convergence time of the Gibbs sampler in Bayesian image restoration <br> Self regenerative Markov chain Monte Carlo <br> Auxiliary variable methods for Markov chain Monte Carlo with applications <br> **Rate of Convergence of the Gibbs Sampler by Gaussian Approximation** <br> Diagnosing convergence of Markov chain Monte Carlo algorithms | RTM ($\psi_e$) |
| Exact Bound for the Convergence of Metropolis Chains <br> Self regenerative Markov chain Monte Carlo <br> **Minorization conditions and convergence rates for Markov chain Monte Carlo** <br> Gibbs-markov models <br> Auxiliary variable methods for Markov chain Monte Carlo with applications <br> Markov Chain Monte Carlo Model Determination for Hierarchical and Graphical Models <br> Mediating instrumental variables <br> A qualitative framework for probabilistic inference <br> Adaptation for Self Regenerative MCMC | LDA + Regression |

Given a new document, which documents is it likely to link to?

# Predicting links from documents

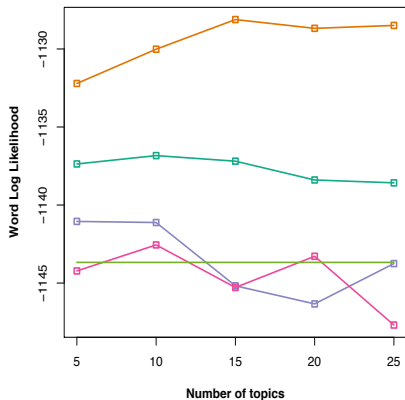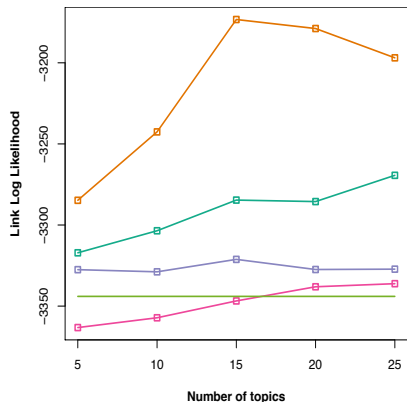| | |
|---|---|
| *Competitive environments evolve better solutions for complex tasks* | |
| **Coevolving High Level Representations**<br>A Survey of Evolutionary Strategies<br>**Genetic Algorithms in Search, Optimization and Machine Learning**<br>**Strongly typed genetic programming in evolving cooperation strategies**<br>Solving combinatorial problems using evolutionary algorithms<br>A promising genetic algorithm approach to job-shop scheduling...<br>Evolutionary Module Acquisition<br>An Empirical Investigation of Multi-Parent Recombination Operators... | RTM ($\psi_e$) |
| A New Algorithm for DNA Sequence Assembly<br>Identification of protein coding regions in genomic DNA<br>Solving combinatorial problems using evolutionary algorithms<br>A promising genetic algorithm approach to job-shop scheduling...<br>A genetic algorithm for passive management<br>The Performance of a Genetic Algorithm on a Chaotic Objective Function<br>Adaptive global optimization with local search<br>Mutation rates as adaptations | LDA + Regression |

Given a new document, which documents is it likely to link to?
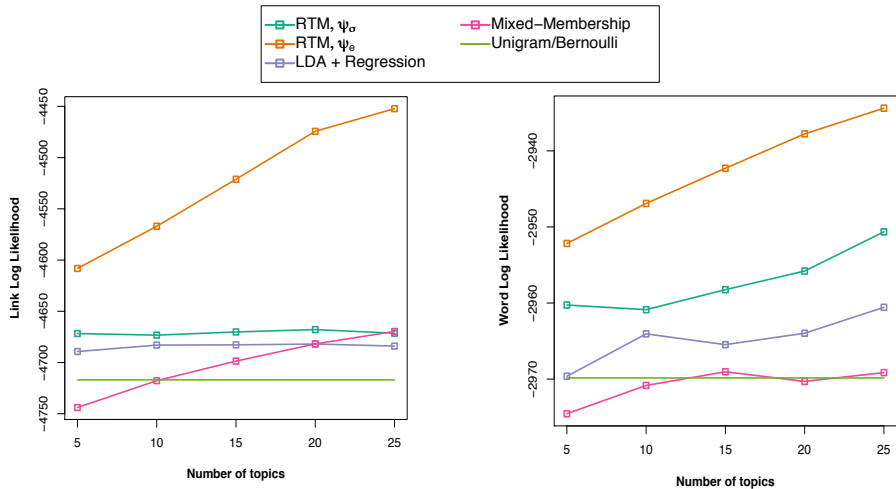
# Predictive performance of each type



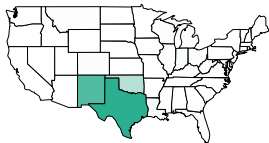**Cora corpus** (McCallum et al., 2000)

# Predictive performance of each type



**WebKB corpus** (Craven et al., 1998)

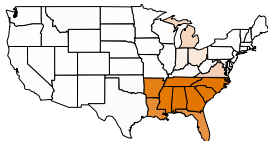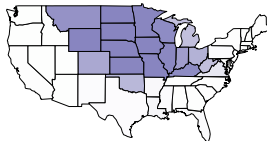# Predictive performance of each type



**PNAS corpus** (courtesy of JSTOR)

# Spatially consistent topics



Topic 1

Topic 2

Topic 3

Topic 4

Topic 5

- For exploratory tasks, RTMs can be used to "guide" the topics

- Documents are geographically-tagged news articles from Yahoo! Links are the adjacency matrix of states

- RTM finds **spatially consistent** topics.

# Relational topic models



- RTMs let us analyze connected documents, modeling both content and connections.

- Most network models cannot predict with new and unlinked data.

- RTMs allow for such predictions
    - links given the new words of a document
    - words given the links of a new document

# The ideal point model



$$p(v_{ij}) = f(d(x_i, a_j))$$

- A model devised to uncover voting patterns (Clinton et al., 2004).
- We observe roll call data $v_{ij}$.
- Bills attached to discrimination parameters $a_j$.
  Senators attached to ideal points $x_i$.

# The ideal point model



- Posterior inference reveals the political spectrum of senators
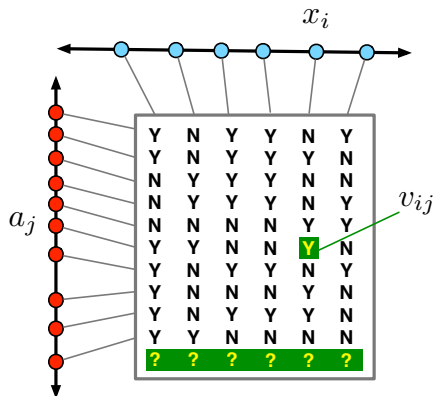- Widely used in quantitative political science.

# The ideal point model is limited for prediction



$$p(v_{ij}) = f(d(x_i, a_j))$$

- We can predict a missing vote.
- But we cannot predict all the missing votes from a bill.
- Cf. the limitations of collaborative filtering

# Ideal point topic models (Gerrish and Blei, 2010)



Use supervised topic modeling assumptions as a predictive
mechanism from bill texts to bill discrimination.

# Ideal point topic models



Bill content (topic model)    Bill sentiment variables    Observed votes    Legislator ideal points

# Ideal point topics



In addition to senators and bills, IPTM places **topics** on the spectrum.

# Prediction on completely held-out votes



Versus the LASSO, the IPTM correctly predicted 126,000 more votes.
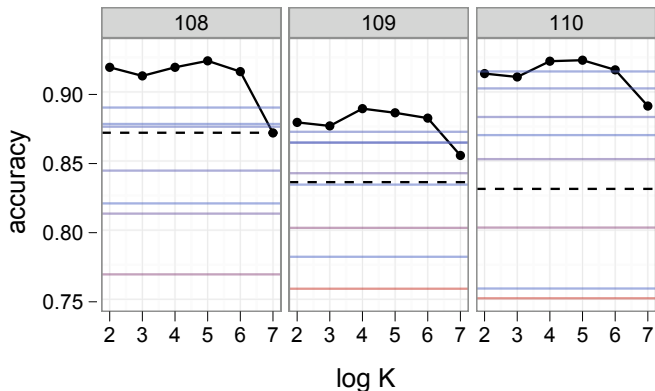
- Many of documents are associated with response variables.

- **Supervised LDA** embeds LDA in a generalized linear model that is conditioned on the latent topic assignments.

- **Relational topic models** use sLDA assumptions with pair-wise responses to model networks of documents.

- **Ideal point topic models** demonstrates how the response variables can themselves be latent variables. In this case, they are used downstream in a model of legislative behavior.

- Note that sLDA and the RTM (and others) are implemented in Jonathan Chang's excellent R package "lda."

# Still other ways to build on LDA

**New applications**—

- Syntactic topic models (Boyd-Graber and Blei 2009)
- Topic models on images (Fei-fei and Perona 2005 and others)
- Topic models on social network data (Airoldi et al. 2008)
- Topic models on music data (Hoffman et al. 2008)
- **Topic models for user recommendation (Wang and Blei, 2011)**

**Testing and relaxing assumptions**—

- Spike and slab priors (Wang and Blei 2009 and Williamson et al. 2010)
- Models of word contagion (Elkan 2006)
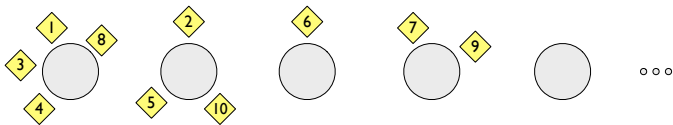- N-gram topic models (Wallach 2006)

# Bayesian nonparametric topic models

- Why Bayesian nonparametric models?
- The Chinese restaurant process
- Chinese restaurant process mixture models
- The Chinese restaurant franchise
- Bayesian nonparametric topic models

# Why Bayesian nonparametric models?

- Topic models assume that the number of topics is fixed.

- It can be determined by cross validation and other model selection techniques.

- Bayesian nonparametric methods skirt model selection—
  - The data determine the number of topics during inference
  - Future data can exhibit new topics

- This is really a field unto itself, but it has found wide application in topic modeling.
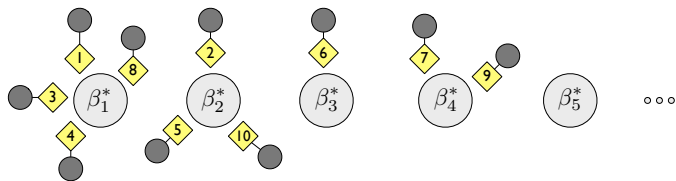
# The Chinese restaurant process (CRP)



- *N* customers arrive to an infinite-table restaurant. Each sits down according to how many people are sitting at each table,

$$p(z_i = k \mid z_{1:(i-1)}, \alpha) \propto \begin{cases} n_k & \text{for} \quad k \leq K \\ \alpha & \text{for} \quad k = K + 1. \end{cases}$$
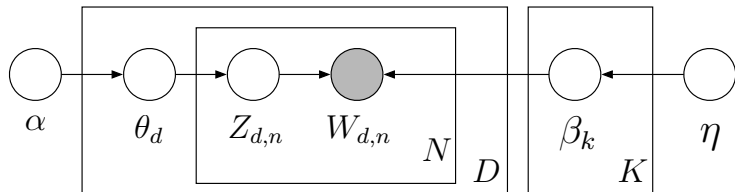
- The resulting seating plan provides a partition

- This distribution is **exchangeable**: Seating plan probabilities are the same regardless of the order of customers (Pitman, 2002).

# CRP mixture models



- Associate each table with a topic ($\beta^*$).
  Associate each customer with a data point (grey node).

- The number of clusters is infinite a priori; the data determines the number of clusters in the posterior.

- Further: the next data point might sit at new table.

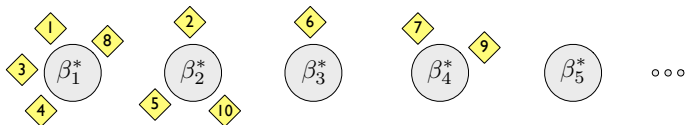- Exchangeability makes inference easy (see Neal, 2000).

# The CRP is not a mixed-membership model



- Mixture models draw each data point from one component.

- The advantage of LDA is that it's a **mixed membership model**.

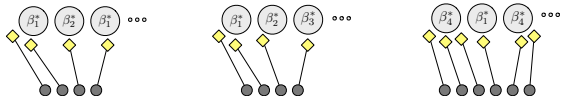- This is addressed by the **Chinese restaurant franchise**.

# The Chinese restaurant franchise (Teh et al., 2006)

**Corpus level restaurant**

*At the corpus level, topics are drawn from a prior.*
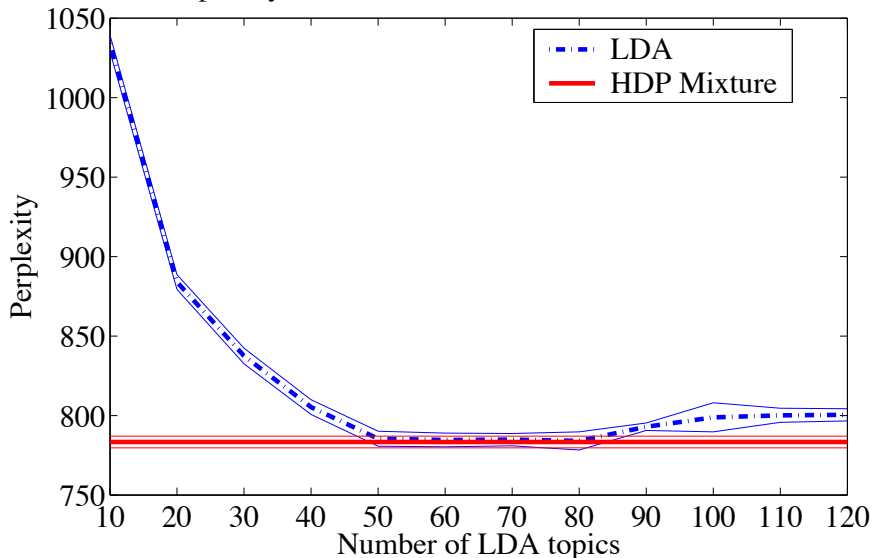


**Document level restaurants**

*Each document-level table is associated with a customer at the corpus level restaurant.*
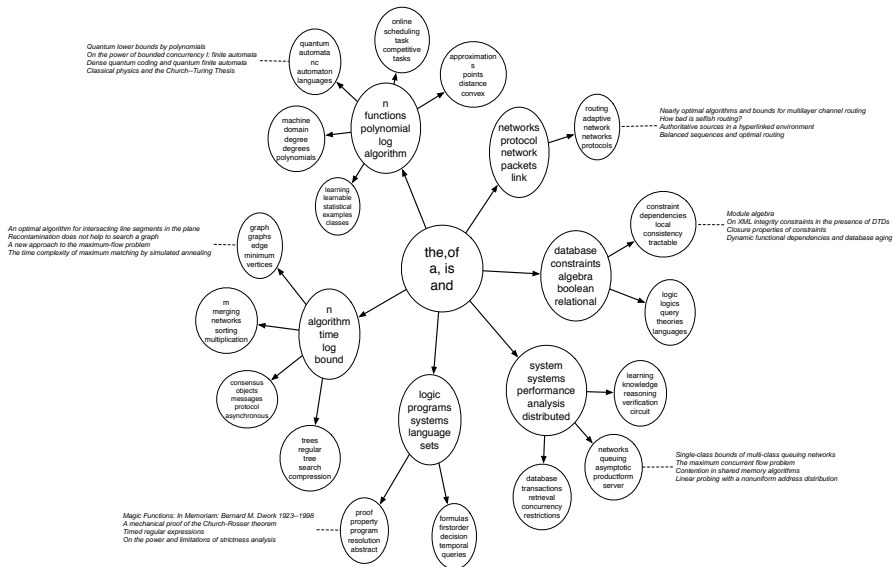


*Each word is associated with a customer at the document's restuarant. It is drawn from the topic that it's table is associated with*

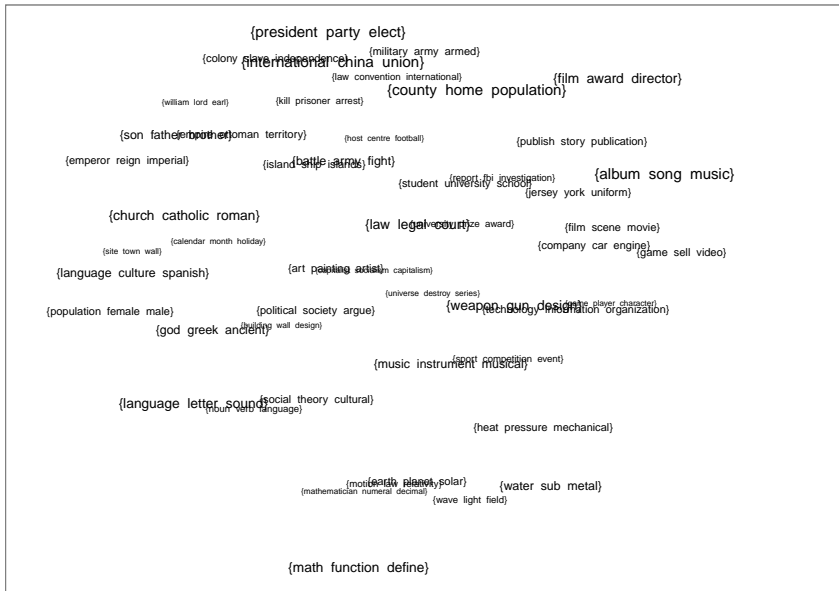# The CRF selects the "right" number of topics



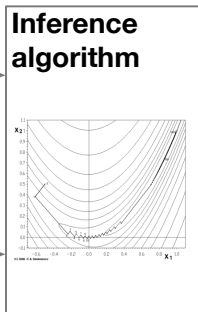Perplexity on test abstacts of LDA and HDP mixture

## Summary: Bayesian nonparametrics

- **Bayesian nonparametrics** is a growing field (Hjort et al., 2011).

- BNP methods can define priors over combinatorial structures.

- In the posterior, the documents determine the particular form of the structure that is best for the corpus at hand.

- These models are also interpretable as **random distribution models**, such as the Dirichlet process (Fergusen 1973, Antoniak 1974).

- Recent innovations:
    - Improved inference methods (Blei and Jordan, 2005)
    - Dependent models, such as time series models (MacEachern 1999, Dunson 2010)
    - Models for predictions (Hannah et al. 2011)
    - Models for matrix factorization and other non-mixtures (Griffiths and Ghahramani, 2011)

# Algorithms

## So far...



- We can express many kinds of assumptions about a corpus.

- Next: How can we analyze it under those assumptions?

# Posterior inference



*Topics*      *Documents*      *Topic proportions and assignments*

- Posterior inference is the main computational problem.
- Inference links observed data to statistical assumptions.
- Inference on large data is crucial for topic modeling applications.

# Posterior inference



*Topics*  ·  *Documents*  ·  *Topic proportions and assignments*

- Our goal is to compute the distribution of the hidden variables conditioned on the documents

$$p(\text{topics, proportions, assignments} \mid \text{documents})$$

# Posterior inference for LDA



- The joint distribution of the latent variables and documents is

$$\prod_{i=1}^{K} p(\beta_i \mid \eta) \prod_{d=1}^{D} p(\theta_d \mid \alpha) \left( \prod_{n=1}^{N} p(z_{d,n} \mid \theta_d) p(w_{d,n} \mid \beta_{1:K}, z_{d,n}) \right).$$

- The posterior of the latent variables given the documents is

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D,1:N} \mid w_{1:D,1:N}).$$

## Posterior inference for LDA



- This is equal to

$$\frac{p(\beta_{1:K}, \theta_{1:D}, \mathbf{z}_{1:D}, \mathbf{w}_{1:D})}{\int_{\beta_{1:K}} \int_{\theta_{1:D}} \sum_{\mathbf{z}_{1:D}} p(\beta_{1:K}, \theta_{1:D}, \mathbf{z}_{1:D}, \mathbf{w}_{1:D})}.$$

- We can't compute the denominator, the marginal $p(\mathbf{w}_{1:D})$.
- This is the crux of the inference problem.

## Posterior inference for LDA



- There is a large literature on approximating the posterior.

- We will focus on
  - Gibbs sampling
  - Mean-field variational methods (batch and online)

# Markov chain Monte Carlo

- Construct a **Markov chain** on the hidden variables, whose limiting distribution is the posterior.

- Collect **independent samples** from that distribution; approximate the posterior with them

- In **Gibbs sampling** the chain is defined by the conditional distribution of each hidden variable given observations and the current setting of the other hidden variables.

# Local and global variables



- Local variables are local to each document
  - Topic proportions $\theta_d$
  - Topic assignemnts $z_{d,n}$

- Global variables are shared by the corpus
  - Topics $\beta_k$

## Local and global variables



- Assume the topics are fixed.

- Even "local inference" is intractable,

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

# Local Gibbs sampling for LDA



- We observe words $\mathbf{w} = w_{1:N}$. The Markov chain is defined on $\{\theta, z_{1:N}\}$, the topic proportions and topic assignments.

- Some notation—

$$
\begin{aligned}
n(z_{1:N}) &= \sum_{n=1}^{N} z_n \\
m_k(\mathbf{z}_{1:D}, \mathbf{W}) &= \sum_{d=1}^{D} \sum_{n=1}^{N} z_{d,n}^k w_{d,n}.
\end{aligned}
$$

- $n(z_{1:N})$ are topic counts;
  $m_k(z_{1:N}, \mathbf{W})$ are within-topic word counts.

## Local Gibbs sampling for LDA



A simple Gibbs sampler is

$$
\begin{aligned}
\theta \,|\, \mathbf{w}, z_{1:N} &\sim \mathrm{Dir}(\gamma) \\
z_n \,|\, \theta, \mathbf{w} &\sim \mathrm{Mult}(\phi_n)
\end{aligned}
$$

where

$$
\begin{aligned}
\gamma &= \alpha + n(z_{1:N}) \\
\phi_n &\propto \theta \cdot p(w_n \,|\, \beta_{1:K}).
\end{aligned}
$$

## Collapsed local Gibbs sampling



- The topic proportions $\theta$ can be integrated out,

$$p(z_n \mid z_{-n}, \mathbf{w}) = p(w_n \mid \beta_{1:K}) \cdot \int_\theta p(z_n \mid \theta) p(\theta \mid z_{-n}) d\theta$$

- A collapsed Gibbs sampler constructs a chain on $z_{1:N}$,

$$z_n \mid z_{-n}, \mathbf{w} \sim \mathrm{Mult}(\phi_n),$$

where $\phi_n \propto p(w_n \mid \beta_{1:K})(n(z_{-n}) + \alpha)$.

# 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 wit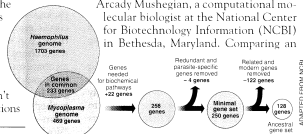h just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

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 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.

# Sampling the topics



- We observe the corpus $\mathbf{W} = \mathbf{w}_{1:D}$.
- We define the chain on $\{\mathbf{z}_{1:D}, \theta_{1:D}, \beta_{1:K}\}$.
- First, sample latent variables $(\mathbf{z}_d, \theta_d)$ for each document.
- Then, sample each topic from

$$\beta_k \,|\, \mathbf{z}_{1:D}, \mathbf{W} \sim \mathrm{Dir}(\lambda_k),$$

where

$$\lambda_k := \eta + m_k(\mathbf{z}_{1:D}, \mathbf{W}).$$

Recall $m_k(\mathbf{z}_{1:D}, \mathbf{W})$ are words counts for topic $k$.

# Collapsed Gibbs sampling with topics



- We can integrate out the topics $\beta_{1:K}$ too.

- The sampler is defined on the topic assigments $\mathbf{z}_{1:D}$,

$$
p(z_{n,d} = k \,|\, \mathbf{z}_{-(n,d)}, \mathbf{W}) \propto \left( \frac{m_k(\mathbf{z}_{-(n,d)}, \mathbf{W}) + \eta}{\sum_v m_k^v(\mathbf{z}_{-(n,d)}) + V\eta} \right) (n_k(z_{-i}) + \alpha)
$$

- This is an excellent Gibbs sampler for LDA. It was developed by Giffiths and Steyvers (2002) and is widely used.

## Example topic inference

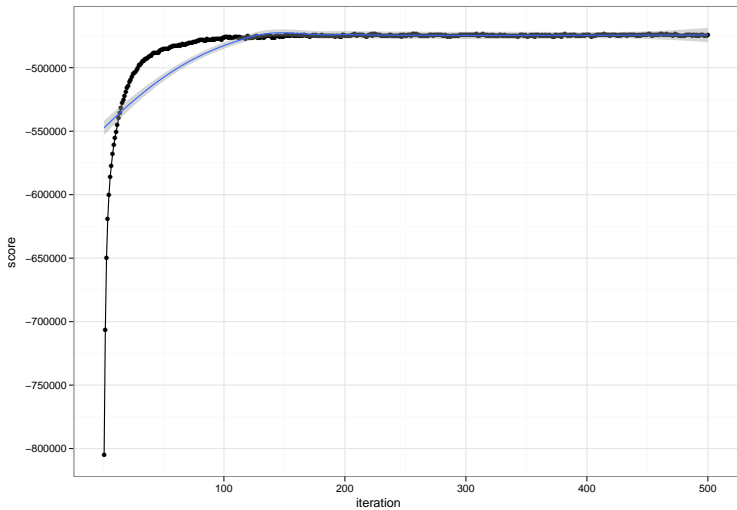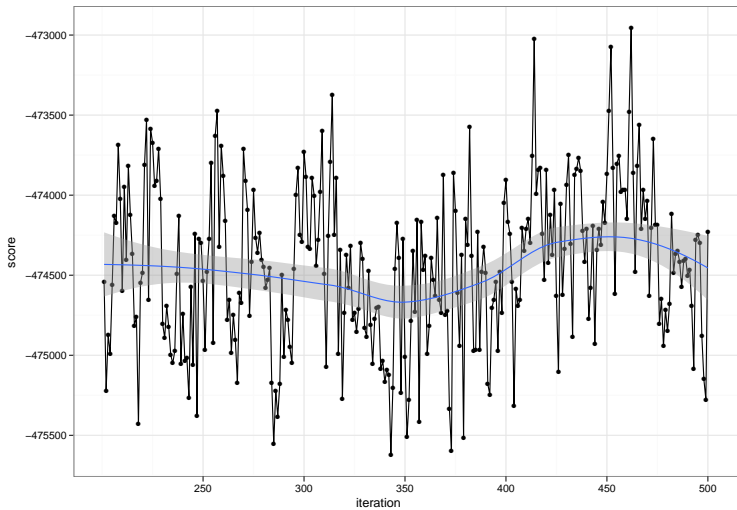| | | | |
|---|---|---|---|
| human | evolution | disease | computer |
| genome | evolutionary | host | models |
| dna | species | bacteria | information |
| genetic | organisms | diseases | data |
| genes | life | resistance | computers |
| sequence | origin | bacterial | system |
| gene | biology | 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 |

## Gibbs sampling for LDA in practice

- In practice:
    1. Obtain a corpus of documents **W**
    2. Run the Gibbs sampler for some number of iterations.
    3. Store states at some lag, or store the MAP state.

- Look at counts like $m_k(\mathbf{z}_{1:D}, \mathbf{W})$ to investigate the topics; look at $n(\mathbf{z}_d)$ to investigate how each document exhibits them.

- **A good habit: Assess the convergence of the chain.**
    - Monitor the log probability of the state & observations. (Its exponential is proportional to the posterior.)
    - Do something fancier, e.g., Raftery and Lewis (1992).

# Assessing convergence example
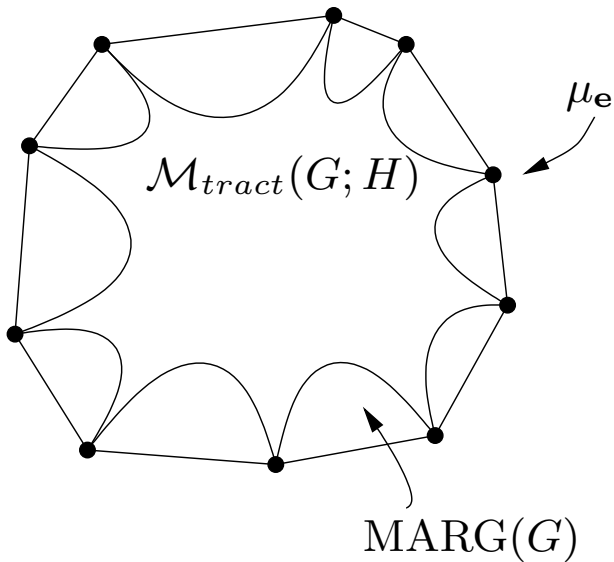
# Assessing convergence example

## Gibbs sampling for LDA

- Simple algorithm for sampling from a complex distribution.

- Works well in practice. Is the best first algorithm to try.

- However
  - Can be slow for very large data sets
  - It is difficult to handle nonconjugacy; it is hard to generalize to the dynamic topic model and correlated topic model.

## Variational inference

- Variational inference replaces sampling with **optimization**.

- The main idea—

  - Place a distribution over the hidden variables with free parameters, called **variational parameters**.

  - Optimize the variational parameters to make the distribution close (in KL divergence) to the true posterior

- In some settings, variational inference is faster than MCMC.

- It is easier to handle nonconjugate pairs of distributions with variational inference. (This is important in the CTM, DTM, etc.)

## A useful picture (from Wainwright and Jordan, 2008)

## Variational inference (in general)

- Let $x = x_{1:N}$ be observed variables;
  let $z = z_{1:M}$ be the latent variables.

- Our goal is to compute the posterior distribution

$$p(z \mid x) = \frac{p(z, x)}{\int p(z, x) dz}$$

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

## Variational inference

- Introduce a distribution over the latent variables $q_\nu(z)$, parameterized by *variational parameters* $\nu$.

- Use Jensen's inequality to bound the log probability of the observations, (Jordan et al., 1999)

$$
\begin{aligned}
\log p(x) &= \log \int p(z, x) dz \\
&= \log \int p(z, x) \frac{q_\nu(z)}{q_\nu(z)} dz \\
&\geq \mathrm{E}_{q_\nu}[\log p(Z, x)] - \mathrm{E}_{q_\nu}[\log q_\nu(Z)]
\end{aligned}
$$

  (J. McAuliffe calls this the **evidence lower bound**, or ELBO.)

- Optimize the variational parameters to tighten this bound.

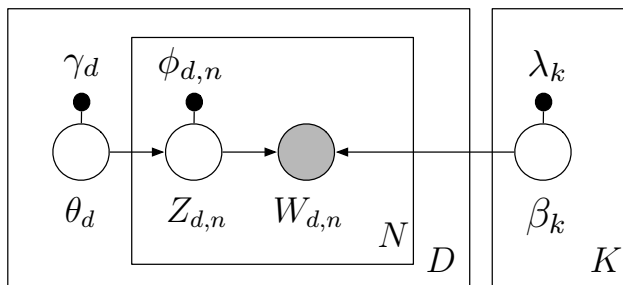- This is the same as finding the member of the family $q_\nu$ that is closest in KL divergence to $p(z \mid x)$.

## Mean-field variational inference

- Complexity is determined by the factorization of $q_\nu$

- In *mean field variational inference $q_\nu$* is fully factored

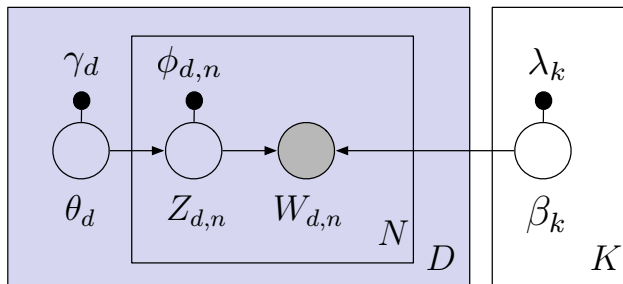$$q_\nu(z) = \prod_{m=1}^{M} q_{\nu_m}(z_m).$$

- Each latent variable is independently governed by its own variational parameter $\nu_m$.

- In the true posterior they can exhibit dependence.
  (Often, this is what makes exact inference difficult.)

# Variational inference for LDA



- The *mean field distribution* places a variational parameter on each hidden variable.

- Optimize these with coordinate ascent, iteratively optimizing each parameter while holding the others fixed.

## Variational inference for LDA



- In the "local step" we iteratively update the parameters for each document, holding the topic parameters fixed.

$$
\begin{aligned}
\gamma^{(t+1)} &= \alpha + \sum_{n=1}^{N} \phi_n^{(t)} \\
\phi_n^{(t+1)} &\propto \exp\{\mathbb{E}_q[\log \theta] + \mathbb{E}_q[\log \beta_{\cdot, w_n}]\}.
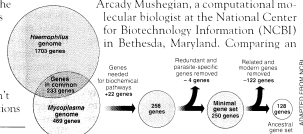\end{aligned}
$$

# Example inference

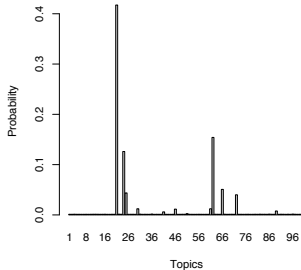## Seeking Life's Bare (Genetic) Necessities

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 complemen-
tary views of the basic genes needed for life.
One research team, using computer analy-
ses 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 esti-
mated 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 hu-
man genome, notes Siv Andersson of Uppsala
University in Sweden, who arrived at the
800 number. But coming up with a consen-
sus 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 mo-
lecular biologist at the National Center
for Biotechnology Information (NCBI)
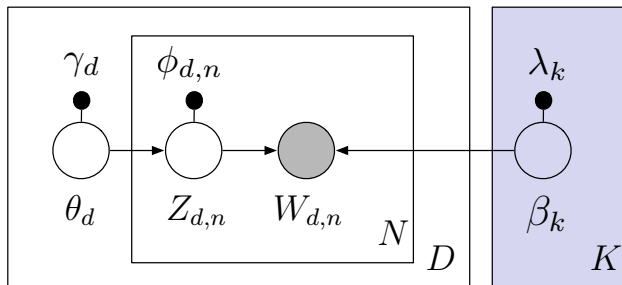in Bethesda, Maryland. Comparing an

* Genome Mapping and Sequenc-
ing, Cold Spring Harbor, New York,
May 8 to 12.

**Stripping down.** Computer analysis yields an esti-
mate of the minimum modern and ancient genomes.

## Variational inference for LDA



- In the "global step" we aggregate the parameters computed from the local step and update the parameters for the topics,

$$\lambda_k = \eta + \sum_d \sum_n w_{d,n} \phi_{d,n}.$$

## Example topic inference

| 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 |

## Variational inference for LDA (sketch)

1: Initialize topics randomly.
2: **repeat**
3:   **for** each document **do**
4:     **repeat**
5:       Update the topic assignment variational parameters.
6:       Update the topic proportions variational parameters.
7:     **until** document objective converges
8:   **end for**
9:   Update the topics from aggregated per-document parameters.
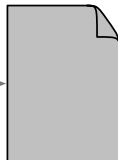10: **until** corpus objective converges.

## Variational inference for LDA

1: Initialize topics $\lambda_{1:K}$ randomly.
2: **while** relative improvement in $\mathcal{L}(\mathbf{w}, \phi, \gamma, \lambda) > \epsilon$ **do**
3:    **for** $d = 1$ to $D$ **do**
4:       Initialize $\gamma_{d,k} = 1$.
5:       **repeat**
6:          Set $\phi_{d,n} \propto \exp\{\mathbb{E}_q[\log\theta_d] + \mathbb{E}_q[\log\beta_{\cdot,w_n}]\}$
7:          Set $\gamma_d = \alpha + \sum_n \phi_{d,n}$
8:       **until** $\frac{1}{K} \sum_k |\text{change in } \gamma_{d,k}| < \epsilon$
9:    **end for**
10:   Set $\lambda_k = \eta + \sum_d \sum_n w_{d,n}\phi_{d,n}$
11: **end while**

## "E step"

```
 1: Initialize topics λ_{1:K} randomly.
 2: while relative improvement in 𝓛(w, φ, γ, λ) > ε do
 3:   for d = 1 to D do
 4:     Initialize γ_{d,k} = 1.
 5:     repeat
 6:       Set φ_{d,n} ∝ exp{𝔼_q[log θ_d] + 𝔼_q[log β_{·,w_n}]}
 7:       Set γ_d = α + ∑_n φ_{d,n}
 8:     until (1/K) ∑_k |change in γ_{d,k}| < ε
 9:   end for
10:   Set λ_k = η + ∑_d ∑_n w_{d,n} φ_{d,n}
11: end while
```

Do variational inference for each document.

## "M step"

1: Initialize topics $\lambda_{1:K}$ randomly.
2: **while** relative improvement in $\mathcal{L}(\mathbf{w}, \phi, \gamma, \lambda) > \epsilon$ **do**
3:   **for** $d = 1$ to $D$ **do**
4:     Initialize $\gamma_{d,k} = 1$.
5:     **repeat**
6:       Set $\phi_{d,n} \propto \exp\{\mathbb{E}_q[\log \theta_d] + \mathbb{E}_q[\log \beta_{\cdot, w_n}]\}$
7:       Set $\gamma_d = \alpha + \sum_n \phi_{d,n}$
8:     **until** $\frac{1}{K} \sum_k |$ change in $\gamma_{d,k}| < \epsilon$
9:   **end for**
10:   Set $\lambda_k = \eta + \sum_d \sum_n w_{d,n} \phi_{d,n}$
11: **end while**

Update the posterior estimates of the topics based on the "E step."

## Online inference for LDA (with M. Hoffman and F. Bach)



*Sample one document*     *Analyze it*     *Update the model*

- Our goal is to use this (and related) models for analyzing massive collections of millions of documents.

- But, in the first step of batch inference we estimate the posterior for *every document* based on randomly initialized topics.

## Online inference for LDA (with M. Hoffman and F. Bach)



*Sample one document*    *Analyze it*    *Update the model*

- Online variational inference is much more efficient.

- It allows us to easily analyze millions of documents.

- It lets us develop topic models on streaming collections.

## Online inference for LDA



1. Randomly pick a document.
2. Perform local variational inference with the current topics.
3. Form "fake" topics, treating the sampled document as though it were the only document in the collection.
4. Update the topics to be a weighted average of the fake topics and current topics.

## Online variational inference for LDA (sketch)

1: Define an appropriate sequence of weights.
2: Initialize topics randomly.
3: **for** ever **do**
4:   Choose a random document $d$.
5:   **repeat**
6:     Update the topic assignment variational parameters.
7:     Update the topic proportions variational parameters.
8:   **until** document objective converges
9:   Compute topics as though $d$ is the only document.
10:  Set the topics to a weighted average of the current topics and
     the topics from step 9.
11: **end for**

## On-line variational inference for LDA

1: Define $\rho_t \triangleq (\tau_0 + t)^{-\kappa}$
2: Initialize $\lambda$ randomly.
3: **for** $t = 0$ to $\infty$ **do**
4:    Choose a random document $w_t$
5:    Initialize $\gamma_{tk} = 1$. (The constant 1 is arbitrary.)
6:    **repeat**
7:       Set $\phi_{t,n} \propto \exp\{\mathbb{E}_q[\log \theta_t] + \mathbb{E}_q[\log \beta_{.,w_n}]\}$
8:       Set $\gamma_t = \alpha + \sum_n \phi_{t,n}$
9:    **until** $\frac{1}{K} \sum_k |\text{change in } \gamma_{t,k}| < \epsilon$
10:    Compute $\tilde{\lambda}_k = \eta + D \sum_n w_{t,n} \phi_{t,n}$
11:    Set $\lambda_k = (1 - \rho_t)\lambda_k + \rho_t \tilde{\lambda}_k$.
12: **end for**

# Analyzing 3.3M articles from Wikipedia



| Documents analyzed | 2048 | 4096 | 8192 | 12288 | 16384 | 32768 | 49152 | 65536 |
|---|---|---|---|---|---|---|---|---|
| **Top eight words** | systems road made service announced national west language | systems health communication service billion language care road | service systems health companies market communication company billion | service systems companies health market communication company billion | service systems companies business company billion health industry | companies systems business company industry market billion | business service companies industry company management systems services | business service companies industry services company management public | business industry service companies services company management public |

# Why does this work?

**A STOCHASTIC APPROXIMATION METHOD**[1]

By Herbert Robbins and Sutton Monro

*University of North Carolina*

**1. Summary.** Let $M(x)$ denote the expected value at level $x$ of the response to a certain experiment. $M(x)$ is assumed to be a monotone function of $x$ but is unknown to the experimenter, and it is desired to find the solution $x = \theta$ of the equation $M(x) = \alpha$, where $\alpha$ is a given constant. We give a method for making successive experiments at levels $x_1$, $x_2$, $\cdots$ in such a way that $x_n$ will tend to $\theta$ in probability.

- Why waste time with the real gradient, when a cheaper noisy estimate of the gradient will do (Robbins and Monro, 1951)?

- Idea: Follow a noisy estimate of the gradient with a step-size.

- By decreasing the step-size according to a certain schedule, we guarantee convergence to a local optimum.

- See Hoffman et al. (2010) and Sato (2001).

## Online inference is promising, in general

- Stochastic variational methods are a general way to approximate the posterior for massive/streaming data.

- No need to process the whole data set in advance; can easily link to web APIs and other data sources

- Powerful algorithm for topic modeling, and can be adapted hierarchical models for many types of data.

- Software and papers: www.cs.princeton.edu/∼blei/

# Latent Dirichlet allocation (flashback)



- This joint defines a posterior.

- From a collection of documents, infer
  - Per-word topic assignment $z_{d,n}$
  - Per-document topic proportions $\theta_d$
  - Per-corpus topic distributions $\beta_k$

- Then use posterior expectations to perform the task at hand, e.g., information retrieval, document similarity, exploration, ...

# Latent Dirichlet allocation (flashback)



Approximate posterior inference algorithms

- 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)
- Online variational inference (Hoffman et al., 2010)

Also see Mukherjee and Blei (2009) and Asuncion et al. (2009).

**Discussion**

## This tutorial

- What are topic models?
- What kinds of things can they do?
- How do I compute with a topic model?
- What are some unsanswered questions in this field?
- How can I learn more?

# Introduction to topic modeling



*Topics*      *Documents*      *Topic proportions and assignments*

- LDA assumes that there are *K* topics shared by the collection.
- Each document exhibits the topics with different proportions.
- Each word is drawn from one topic.
- We discover the structure that best explain a corpus.

# Extensions of LDA



- Topic models can be adapted to many settings

- Bayesian nonparametric topic models let the corpus determine the number of topics (or more complicated topic structure).

# Posterior inference



- Posterior inference is the central computational problem.

- We discussed three algorithms
  - MCMC based on collapsed Gibbs sampling
  - Mean-field variational inference
  - Online variational inference

## Some open issues

- **Model interpretation and model checking**
  Which model should I choose for which task?
  (Chang et al. 2009, Ramadge et al. 2009, Newman et al. 2010, Mimno and Blei 2011, Mimno et al. 2011)

- **Incorporating corpus, discourse, or linguistic structure**
  How can our knowledge of language help us build and use exploratory models of text?

- **Interfaces and "downstream" applications of topic modeling**
  What can I do with an annotated corpus? How can I incorporate latent variables into a user interface?

- **Theoretical understanding of approximate inference**
  What do we know about variational inference from either the statistical or learning perspective?

# Interpretation I: Human studies of topic models



(see Chang et al. 2009 and Newman et al. 2010)

# Interpretation II: Labelled LDA on JSTOR



### Tax Innovation in the States: Capitalizing on Political Opportunity

Frances Stokes Berry;William D. Berry. *American Journal of Political Science* (1992), pp. 715-742

Journal Disciplines:

- **Political Science**

Pie chart labels: Political Science, Statistics, Finance, Economics, Business, Public Policy and Administration

### Chaos and Nonlinear Forecastability in Economics and Finance

Blake LeBaron. *Philosophical Transactions: Physical Sciences and Engineering* (1994), pp. 397-404

Journal Disciplines:

- **Mathematics**
- **Biological Sciences**
- **General Science**

Pie chart labels: Statistics, Mathematics, Business, Economics, Other, Finance

### Reply: Theory Is Not a Social Dilemma

Gerald Marwell;Pamela Oliver. *Social Psychology Quarterly* (1994), pp. 373

Journal Disciplines:

- **Psychology**
- **Sociology**

Pie chart labels: Political Science, Philosophy, Library Science, History of Science and Technology, Feminist and Women's Studies, Other, Sociology

(see Ramadge et al. 2009 and **Ramadge et al. 2011**)

# Interptetation III: Mutual information discrepancy



(see Mimno and Blei 2011)

## Topic modeling resources

- The topic modeling mailing list is a good discussion group.
- Bibliography: http://www.cs.princeton.edi/~mimno/
- Software and papers: http://www.cs.princeton.edu/~blei/

# If you remember one picture...



**Assumptions**

**Data**

**Inference algorithm**

$x_2$

$x_3$

**Discovered structure**

"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."
(J. Tukey, *The Future of Data Analysis*, 1962)