

Optimization and Probability in Machine Learning

Neal Parikh

<http://nparikh.org>

August 19, 2019

Machine learning

- intersection of computer science and statistics
- algorithms that learn from data and/or improve with experience
- roots in AI, but also used in huge variety of other domains
- draws on a number of mathematical areas, but rests in particular on optimization theory and probability theory
- by framing various tasks that don't appear to involve optimization in that way, can bring the optimization toolbox to bear on a wider range of problems in machine learning

Outline

Regularized loss minimization

Latent variables and the EM algorithm

Variational inference

Conclusion

Regularized loss minimization

optimization-based framework for learning (parameter estimation)

$$\text{minimize } l(w) + \lambda r(w)$$

- $w \in \mathbf{R}^n$ are the **model parameters** or **weights**
- $l : \mathbf{R}^n \rightarrow \mathbf{R}$ is a **loss function** measuring lack of fit on training data
- $r : \mathbf{R}^n \rightarrow \mathbf{R}$ is a **regularizer** measuring model complexity
- $\lambda > 0$ is a **regularization parameter** trading off between the two

if l and r are convex, problem can essentially be solved

Regularized loss minimization

- many models fall under this framework, and may or may not be probabilistic in nature
- ridge regression (Tikhonov regularized linear regression)

$$\text{minimize } \|Xw - y\|_2^2 + \lambda\|w\|_2^2$$

- lasso (sparse linear regression)

$$\text{minimize } \|Xw - y\|_2^2 + \lambda\|w\|_1$$

- support vector machine

$$\text{minimize } \sum_{i=1}^N (1 - y_i w^T x_i)_+ + \lambda\|w\|_2^2$$

Regularized loss minimization

- also subsumes common ways of learning probabilistic models
- maximum likelihood estimation of model $p(x)$

$$\text{minimize} \quad - \sum_{i=1}^N \log p(x_i)$$

- could also model $y | x$ and maximize $\sum_i \log p(y_i | x_i)$
- includes, e.g., linear regression, generalized linear models, ...

- maximum a posteriori (MAP) estimation (finding posterior mode)

$$\text{minimize} \quad - \sum_{i=1}^N \log p(x_i | w) - \log p(w)$$

- $w \sim$ Gaussian: Tikhonov regularization
- $w \sim$ Laplacian: ℓ_1 /lasso regularization

Naive Bayes

- joint model that factorizes as

$$p(\mathbf{x}, y) = p(y) \prod_{i=1}^n p(x_i | y)$$

where x_i are words in an email and $y \in \{\text{spam}, \text{ham}\}$

- need to estimate $p(y = \text{spam})$ and $p(x = \text{bank} | y)$
- maximum likelihood estimates have a trivial closed form solution, e.g.:

$$\hat{p}(y = \text{spam}) = \frac{\# \text{ spam emails}}{\text{total } \# \text{ emails}}$$

$$\hat{p}(x = \text{bank} | \text{spam}) = \frac{\# \text{ times } x = \text{bank in spam mails}}{\# \text{ words in spam emails}}$$

Naive Bayes

- once parameters are fit, want to classify new example \mathbf{x}^{new}
- compute label posterior

$$p(y = \text{spam} | \mathbf{x}^{\text{new}}) = \frac{p(\mathbf{x}^{\text{new}} | \text{spam}) \cdot p(\text{spam})}{p(\mathbf{x}^{\text{new}})}$$

- pick whichever class has higher posterior probability
- in general, called a (probabilistic) **inference** problem

Gaussian discriminant analysis

- assume data comes from generative model

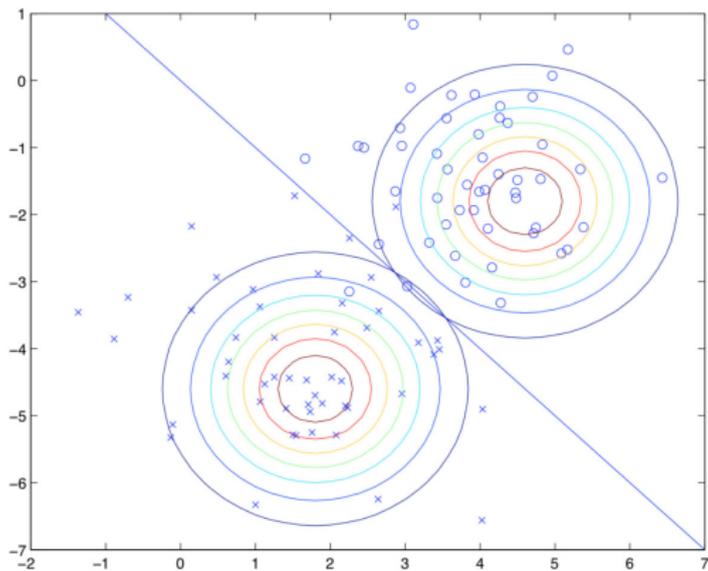
$$\begin{aligned}y &\sim \text{Bernoulli}(\phi) \\x | y = 0 &\sim \text{N}(\mu_0, \Sigma) \\x | y = 1 &\sim \text{N}(\mu_1, \Sigma)\end{aligned}$$

i.e., data comes from one of two Gaussians chosen with a ϕ -coin flip

- estimate $w = (\phi, \mu_k, \Sigma)$ by maximizing

$$\begin{aligned}\ell(w) &= \log \prod_{i=1}^N p(x_i, y_i; w) \\&= \sum_{i=1}^N \log p(x_i | y_i; \mu_0, \mu_1, \Sigma) + \sum_{i=1}^N \log p(y_i; \phi)\end{aligned}$$

Gaussian discriminant analysis



Maximum likelihood estimation

maximum likelihood estimates of parameters given by

$$\begin{aligned}\hat{\phi} &= \frac{1}{N} \sum_{i=1}^N [y_i = 1] \\ \hat{\mu}_k &= \frac{\sum_{i=1}^N [y_i = k] x_i}{\sum_{i=1}^N [y_i = k]} \\ \hat{\Sigma} &= \frac{1}{N} \sum_{i=1}^N (x_i - \mu_{y_i})(x_i - \mu_{y_i})^T\end{aligned}$$

very natural interpretations:

- $\hat{\phi}$ is empirical proportion of positive label in \mathcal{D}
- $\hat{\mu}_k$ is empirical average of x_i with label k
- $\hat{\Sigma}$ is empirical covariance, with variance measured to relevant mean

Regularized loss minimization

- subsumes a very wide class of machine learning models
- generative/discriminative, probabilistic/non-probabilistic, ...
- sometimes problem can just be solved in closed form
- sometimes need sophisticated optimization algorithms like L-BFGS, accelerated proximal gradient method, ...
- calculations, or subsequent model use, may involve probability operations

Outline

Regularized loss minimization

Latent variables and the EM algorithm

Variational inference

Conclusion

Mixture of Gaussians

- probabilistic model for clustering / density estimation
- consider data $\mathcal{D} = \{x_1, \dots, x_N\}$
- generative model

$$\begin{aligned}z &\sim \text{Multinomial}(\phi) \\x | z = k &\sim \text{N}(\mu_k, \Sigma_k)\end{aligned}$$

- *i.e.*, each x_i generated by sampling a **unobserved** (hidden, latent) $z_i \in [K]$ and then drawing x_i from the corresponding Gaussian
- presence of these latent variables is the key new wrinkle
- model parameters are ϕ, μ_k, Σ_k

Maximum likelihood estimation

- model parameters are ϕ, μ_k, Σ_k
- as usual, write down likelihood for $w = (\phi, \mu_k, \Sigma_k)$

$$\begin{aligned}\ell(w) &= \sum_{i=1}^N \log p(x_i; w) \\ &= \sum_{i=1}^N \log \sum_{z_i=1}^K p(x_i | z_i) p(z_i)\end{aligned}$$

- this function is *nonconvex* due to sum over values of z_i
- can no longer easily solve the relevant optimization problem

Maximum likelihood estimation

- if z_i were known, problem is easy and becomes

$$\ell(w) = \sum_{i=1}^N \log p(x_i | z_i) + \sum_{i=1}^N \log p(z_i)$$

- maximizing with respect to ϕ, μ, Σ gives

$$\phi_j = \frac{1}{N} \sum_{i=1}^N [z_i = j], \quad \mu_j = \frac{\sum_{i=1}^N [z_i = j] x_i}{\sum_{i=1}^N [z_i = j]}$$

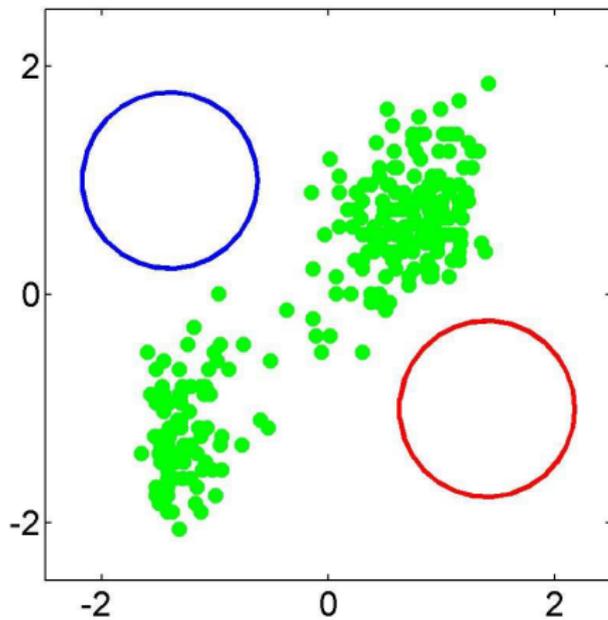
similar expression for Σ

- *i.e.*, if z_i were known, nearly identical to maximum likelihood estimates in GDA (with z_i 's as class labels)

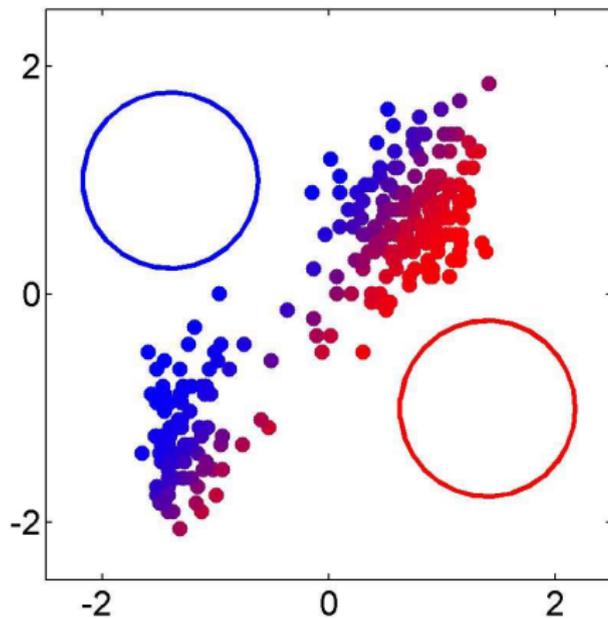
EM algorithm

- **idea:** iteratively guess the z_i and then use formulas above:
 - ① E-step (probability): compute $\rho_{ij} = p(z_i = j | x_i; \theta, \mu, \Sigma)$
 - ② M-step (optimization): use formulas above with ρ_{ij} in place of $[z_i = j]$
- E-step is an inference task: compute posterior probability of z_i 's, given data and current setting of parameters; 'soft guesses' for values of z_i
- M-step is 'regular' maximum likelihood estimation, but there is uncertainty around the value of the z_i and that's incorporated in estimates
- yields a 'soft' version of k -means for this model

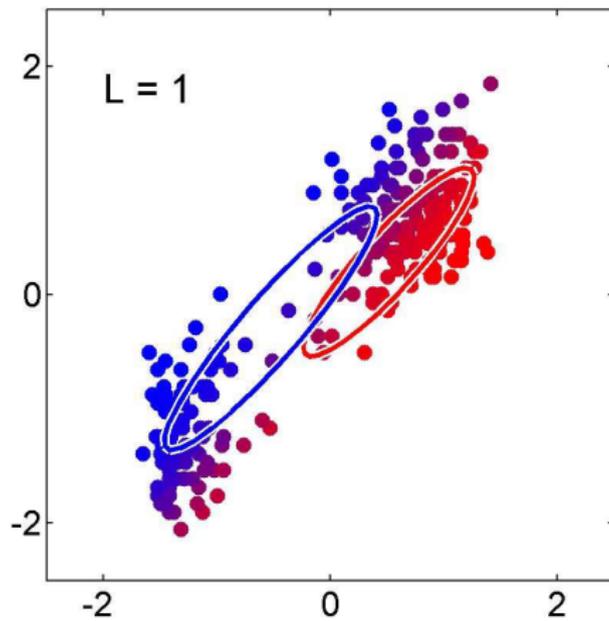
Gaussian mixture model



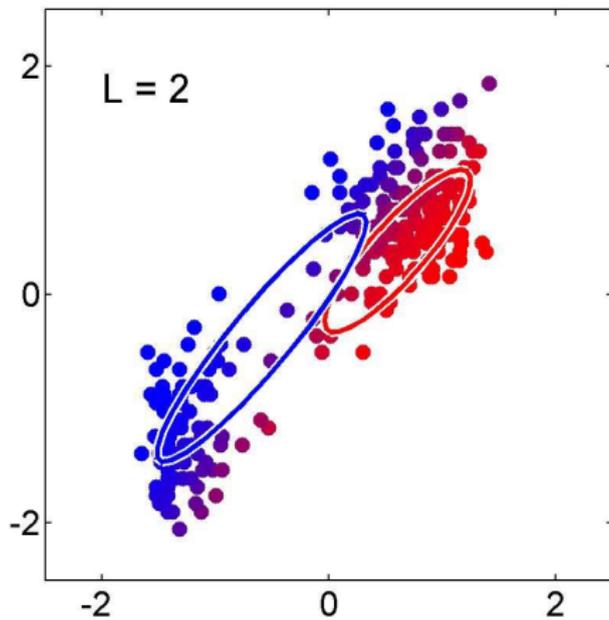
Gaussian mixture model



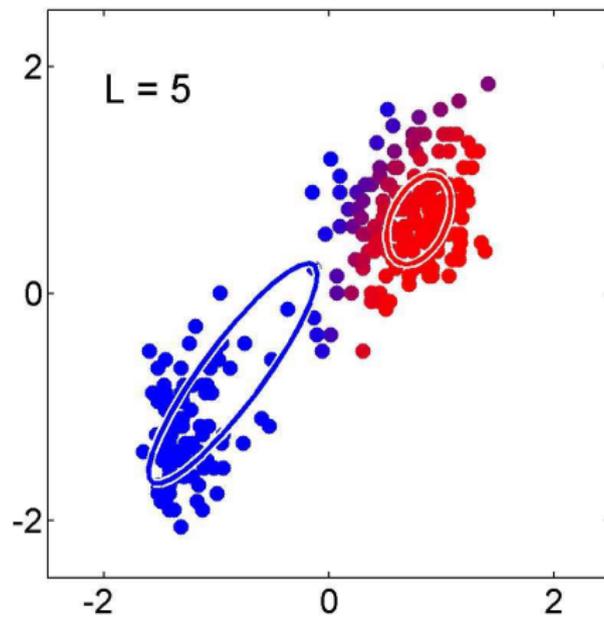
Gaussian mixture model



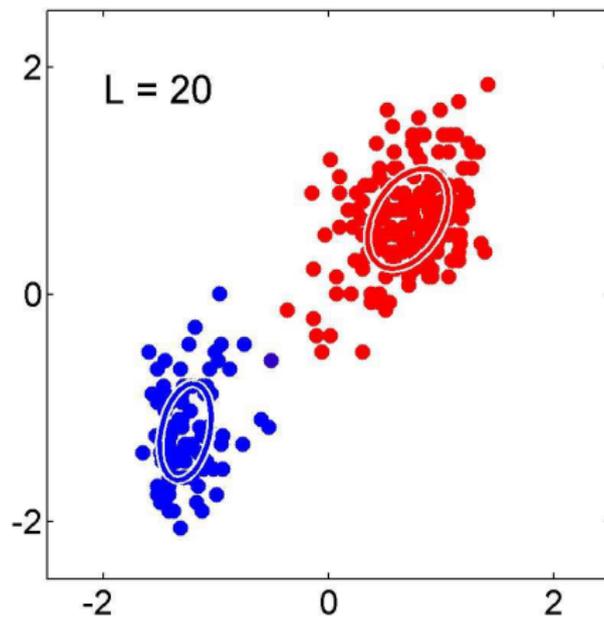
Gaussian mixture model



Gaussian mixture model



Gaussian mixture model



EM algorithm

- in general, EM algorithm is standard approach to maximum likelihood estimation with latent variable models
- data $\mathcal{D} = \{x_1, \dots, x_N\}$
- want to fit model $p(x, z)$ with z hidden
- likelihood is given by

$$\ell(w) = \sum_{i=1}^N \log p(x_i; w) = \sum_{i=1}^N \log \sum_z p(x_i, z; w)$$

- often the case that maximum likelihood estimation of x would be easy if z were known, so alternate the two steps

EM algorithm

- iteratively lower bound ℓ , then maximize that lower bound
- for each i , let q_i be a distribution over z 's

$$\begin{aligned}\sum_{i=1}^N \log p(x_i) &= \sum_{i=1}^N \log \sum_{z_i} p(x_i, z_i) \\ &= \sum_{i=1}^N \log \sum_{z_i} q_i(z_i) \frac{p(x_i, z_i)}{q_i(z_i)} \\ &\geq \sum_{i=1}^N \sum_{z_i} q_i(z_i) \log \frac{p(x_i, z_i)}{q_i(z_i)}\end{aligned}$$

by Jensen's inequality

EM algorithm

- previous formula gives lower bound for *any* q_i ; ideally, have the lower bound be tight (inequality holds with equality) for current value of w
- can show that this is the case when $q_i(z_i) = p(z_i | x_i; w)$
- E-step: lower bound ℓ via computing $p(z | x)$
- M-step: maximize this lower bound

$$E_q[\log p(z, x; w)]$$

with respect to w

EM algorithm

- previous motivation is as a ‘majorization-minimization’ algorithm
- can also be viewed as coordinate ascent on

$$F(q, w) = \sum_{i=1}^N \sum_{z_i} q_i(z_i) \log \frac{p(x_i, z_i; w)}{q_i(z_i)}$$

- E-step: maximization with respect to q
- M-step: maximization with respect to w
- (note: can also be modified for MAP estimation)
- this perspective suggests/justifies many variations

Outline

Regularized loss minimization

Latent variables and the EM algorithm

Variational inference

Conclusion

Inference

comes up in several places:

- using a trained model, e.g., to predict out-of-sample outcomes
- E-step of EM for MLE in partially observed model
- Bayesian learning (work with joint model $p(x, z, w)$ with w random)

trivial in very simple cases, but expensive or intractable in complex models

(note: in some models, like hidden Markov models or Kalman filters, inference can be carried out exactly but requires use of an algorithm)

Variational inference

- let x be observed and z be hidden
- interested in computing the posterior distribution

$$p(z | x) = \frac{p(z, x)}{p(x)} = \frac{p(z, x)}{\int_z p(z, x)}$$

- denominator ('evidence') is hard to compute and makes this difficult
- main idea is to pick family \mathcal{Q} of distributions over the latent variables indexed by *variational parameters*

$$q(z | \nu)$$

and set ν *via optimization* to make q close to $p(z | x)$

- *i.e.*, turn a probability problem into an optimization problem

Variational inference

- recall following lower bound from EM

$$\log p(x) \geq L(q) = \mathbb{E}_q[\log p(x, z)] - \mathbb{E}_q[\log q(z)]$$

where RHS is called the **evidence lower bound**

- here, choose a parametrized family of distributions for q such that these expectations are computable, then maximize lower bound L with respect to these 'variational parameters'
- can show that

$$\text{KL}(q(z) \parallel p(z \mid x)) + L(q) = \log p(x)$$

- minimizing KL divergence is equivalent to maximizing L , plus obtain a lower bound on $\log p(x)$

Mean field family

- in mean field variational inference, assume that the family factorizes

$$q(z_1, \dots, z_m) = \prod_{j=1}^m q(z_j)$$

i.e., all variables are independent

- typically, this family does not contain the true posterior because the hidden variables are dependent (and these dependencies are what make the posterior difficult to work with)
- in 'coordinate ascent variational inference', iteratively optimize each variational distribution while holding others fixed
- computations end up being simple when relevant parts of original model are exponential family distributions

Additional topics

- **variational EM**

- use variational inference to compute approximate posterior $p(z | x)$ in E-step, *i.e.*, do inexact maximization of F with respect to q

- **variational Bayes**

- parameters w are random and model is $p(x, z, w)$
- use lower bound

$$\log p(x) = \iint p(x, z, w) dz dw \geq \mathbb{E}[\log p(x, z, w)] - \mathbb{E}[\log q(z, w)]$$

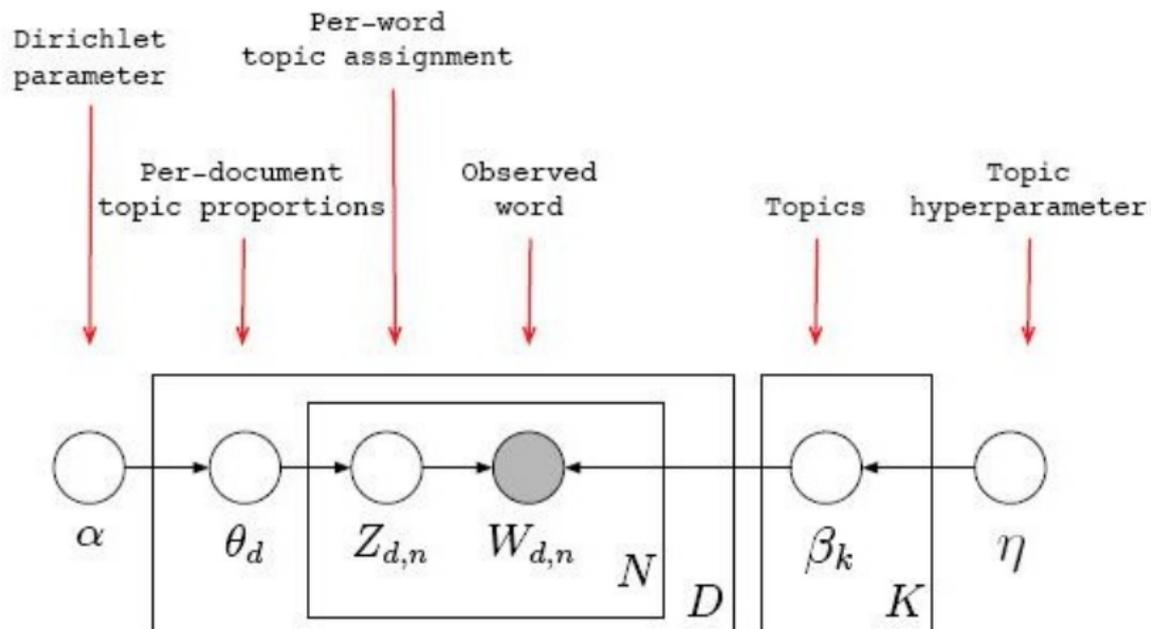
with factorized approximation $q(z, w) = q(z)q(w)$ and do alternating maximization w.r.t. z, w

- yields EM-like algorithm sometimes called ‘variational Bayesian EM’
- **stochastic variational inference**: use stochastic optimization to scale optimization carried out in variational inference

Latent Dirichlet Allocation

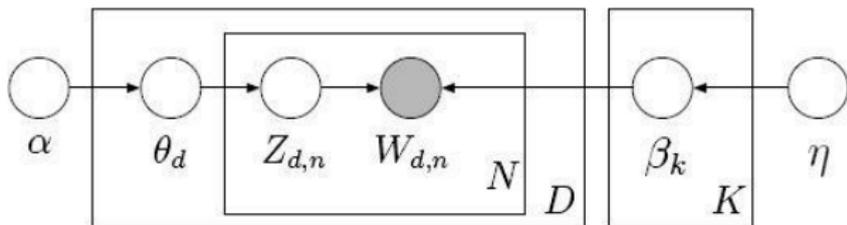
- **words** are multinomial random variables w
- **documents** are sequences of N words $\mathbf{w} = (w_1, \dots, w_N)$
- **topics** are (multinomial) distributions over words
- model document as a random mixture θ over K latent topics
- from an *unlabeled* collection of documents, infer
 - per-word topic assignments in each document
 - per-document topic proportions
 - per-corpus topic distributions

Latent Dirichlet Allocation



Topic models

models for discovering thematic structure in document collections



joint distribution of topic mixture θ , topic distributions \mathbf{z} , and words \mathbf{w} is:

$$p(\theta, \mathbf{z}, \mathbf{w} \mid \alpha, \beta) = p(\theta \mid \alpha) \prod_{n=1}^N p(z_n \mid \theta) p(w_n \mid z_n, \beta)$$

goal: fit parameters and compute posterior $p(\theta, \mathbf{z} \mid \mathbf{w}, \alpha, \beta)$

A 100 topic model of Science 1980-2000

sound speech acoustic language sounds	quantum laser light optical electron	brain memory human visual cognitive	computer data information problem computers	ice climate ocean sea temperature
stars universe galaxies astronomers star	research national science new funding	materials organic molecules molecular polymer	fossil species evolution birds evolutionary	volcanic years fig deposits rocks

Topic proportions in documents

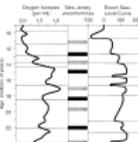
pelagic Group (Benthic) of the University of Miami, coral-reef scientist at Leg 185.

Dr. Inge Leh, 60, swapped up with her mostly oceanographic results over the last year for paleoceanology supports the Eocene curve. "It was (close to) my first job," she tells. "In other places, especially when you go back beyond 12 million years ago, we have to fit together pieces that aren't there." For her, as well as other scientists who study marine sediment cores from recent data from offshore Brazil. There, Vitor Alves and Gordon Hubbert of the University, using well data provided by the Brazilian company Petrobras, tracked sea-level changes that correlate very well with the Florida data. Their core, The researchers believe has zero noise and Eocene sea level, Brazil, gives that the mean of sea-level Eocene curve is not about 22 years old. "We do not believe that either."

The double-barreled documentation of the curve, likely to reveal all, includes, though, Andrew Miller of the University of Toronto, for example, reports a steadily opposite. "I don't think this is good science at all. There are no more reports to Eocene curve, but the margin of error indicates as large as they could be, meaning only," he says. Indeed, Bhat has shown good correlations between the Eocene curve and modern preserved sets of events.

André's point is well taken," says Mike. Matching a sea level change from one site to the Eocene curve is highly subjective, he notes, so there has been a tendency to make mistakes when one site. But, he says, "we're making the timing." At one point, it's reasonable to say that these changes are correlated and therefore they are closely related. Mike adds that "Whether Mike is scientifically correct or not is difficult for me to answer for one of these things, especially when we talk to him, it seems like working."

A mysterious mechanism
Even if the Eocene curve is a highly resolved global indication of sea level, it likely tracks another phenomenon, which is driving sea level change. Researchers have presented data on the mechanism and timing of sea level rise. But the Eocene curve indicates that the mechanism and timing of sea level rise for the entire time and tells us a significant difference from the Florida data. "And while researchers have been able to link the Eocene curve to sea level during the recent past, the link is poor or unclear time. To improve past sea level, researchers apply the Eocene curve to the Eocene of offshore sediments. As global sea levels



Sea changes From about a sea level rise, 100 million years ago, the Eocene curve (solid line) and Florida data (dashed line) and Brazil data (dotted line) show sea level changes in meters.

at the expense of water or stable on the ocean's surface, changes the marine composition of minerals and the carbonate skeletons of marine organisms.

Now the Leg 185 group has combined these changes to compare isotopes with their New Jersey sea level changes, and with the Eocene curve, back to 50 million years ago. And a paper in press by George Miller and James Browning of Rutgers connects the link between isotopic changes and the Eocene curve over 100 million years ago. Although evidence of isotopic data also shows signs of sea level rise on land, ranging up to 48 million years ago. But before that, while the world was experiencing the warmest heat wave of the past 65 million years, both groups had

but the correlation falls apart, leaving two weaknesses to the sea level changes.

In the evidence to report global change, sea level curves over correlation. Barbara Stoll and David Schrag of Princeton University, have used isotopes preserved in carbonate to track the expansion of continental margins and sea level during the period of relative warmth 80 million to 130 million years ago, when oxygen isotope records are available. When filling sea level expansion evidence to leading by study, some the amount of expansion in the world ocean increases. In work presented in his talk, a member of the American Geophysical Union, the researchers found that oceanic expansion included in a few hundred thousand years, suggesting rapid sea level drops of 15 to 50 meters, and the deep oceanic water temperature in the Eocene curve. Stoll and Schrag also link to global expansion, suggesting that sea level rise has temporarily gone large enough to lower sea level, or processes that also, give signs in the low record of lake, high latitude climate.

If climate didn't drive sea level up, and James what? The geologists believe, plans have been suggested. Kenneth has one proposal that marine isotopes might have done the job in world time, by changing isotopic ratios, the other's little evidence for such forces. "The link is having problems" with the Eocene curve in other times, concludes Hag. "because the mechanism is still unclear." Geologists may not be willing to accept Eocene's job, but he hasn't yet been accepted as a mechanism.

—Richard A. Kerr

GENOME MEETING

Seeking Life's Bare (Genetic) Necessities

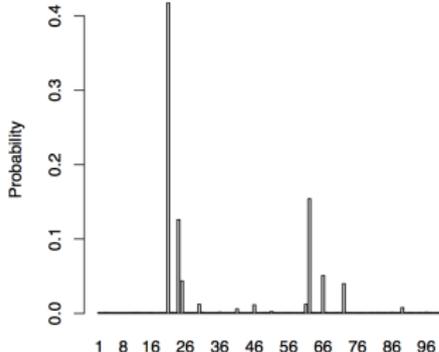
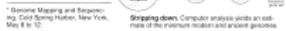
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genomic meeting here, 700 scientists presented with radically different approaches to answer this question.

One research team, using computer analysis to compare known genes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms may have needed only 120 genes. The other research group argued that a single genetic and molecular tool that has appeared, 80% of the time, is essential for that organism. "It's not about the number, but the quality of the genes," says one researcher.

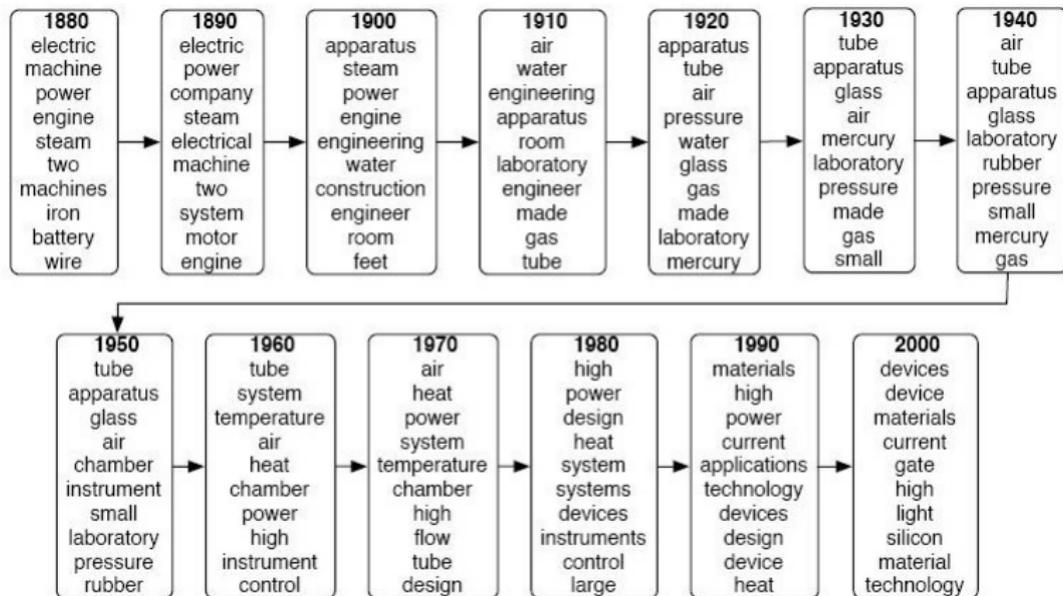
"We're not all that far apart," especially in comparison to the 24,000 genes in the human genome, says the biologist at Cornell University in Ithaca, who arrived at the 80 number. But coming up with common answers may be more than just a matter of numbers, particularly in cases that more genes may be complexly regulated and controlled. "It may be a matter of organizing our really organized genome," explains Anand Mahalingam, a computer scientist 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.

Shrapping down Computer analysis yields an estimate of the minimum number and ancient genomes.

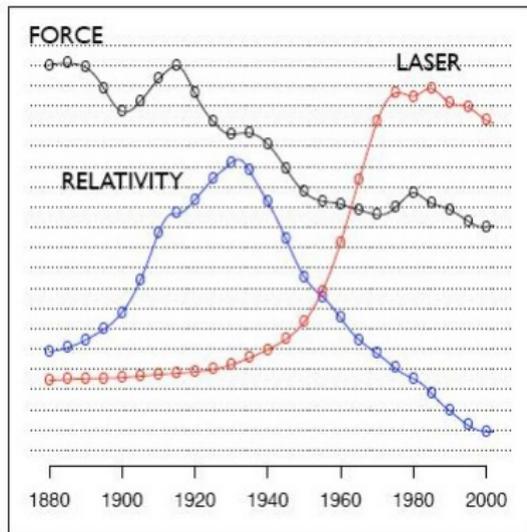


Model Evolution of Topics over Time

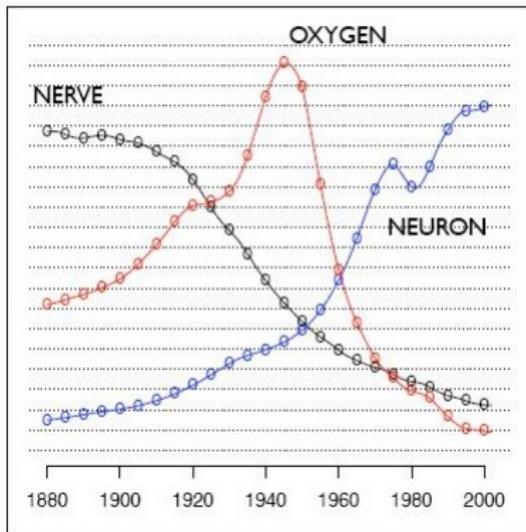


Visualizing Trends Within Topics

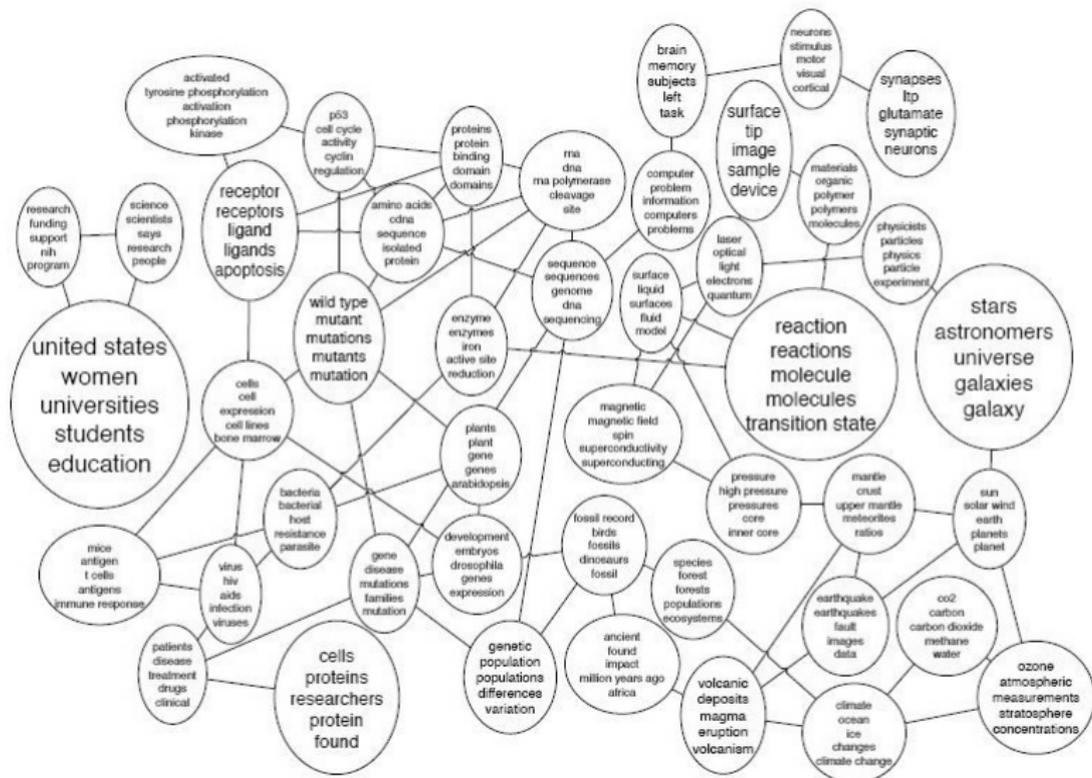
"Theoretical Physics"



"Neuroscience"

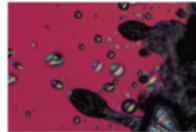
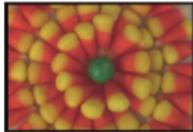


Model Connections Between Topics



Matching Words and Pictures

Candy



Sunset



People
& Fish



Matching Words and Pictures



True caption
market people
Corr-LDA
people market pattern textile display



True caption
scotland water
Corr-LDA
scotland water flowers hills tree



True caption
bridge sky water
Corr-LDA
sky water buildings people mountain



True caption
sky tree water
Corr-LDA
tree water sky people buildings



True caption
birds tree
Corr-LDA
birds nest leaves branch tree



True caption
fish reefs water
Corr-LDA
fish water ocean tree coral



True caption
mountain sky tree water
Corr-LDA
sky water tree mountain people



True caption
clouds jet plane
Corr-LDA
sky plane jet mountain clouds

Variational EM for LDA

- computing evidence $p(\mathbf{w} | \alpha, \beta)$ is intractable, but evidence lower bound

$$\mathbb{E}_q[\log p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta)] - \mathbb{E}_q[\log q(\theta, \mathbf{z} | \gamma, \phi)]$$

gives lower bound on $\log p(\mathbf{w} | \alpha, \beta)$

- plug in form of p and use family of approximate posteriors given by

$$q(\theta, \mathbf{z} | \gamma, \phi) = q(\theta | \gamma) \prod_{n=1}^N q(z_n | \phi_n)$$

where γ is a variational Dirichlet parameter and ϕ_n are variational multinomial parameters

- variational E-step: maximize lower bound with respect to γ and ϕ_n via alternating maximization (both simple closed form expressions)
- M-step: maximize with respect to hyperparameters α (simple numerical method), β (closed form)

Outline

Regularized loss minimization

Latent variables and the EM algorithm

Variational inference

Conclusion

Conclusions

- interactions of optimization and probability in machine learning
- using probabilistic structure can ease optimization
- framing probabilistic computations in variational form can help bring full optimization toolbox to bear on wider range of problems
- can lead to fast and scalable algorithms that enable working with very complex probabilistic models on huge datasets, beyond the reach of other methods
- many analogies between probabilistic methods (e.g., Gibbs sampling) and optimization-based methods (e.g., coordinate ascent)

Acknowledgements

- David Blei & collaborators
- Michael Jordan & collaborators
- Andrew Ng (diagrams)