Introduction to Machine Learning

Brown University CSCI 1950-F, Spring 2012 Prof. Erik Sudderth

Lecture 24: Topic Models & Latent Dirichlet Allocation Monte Carlo Methods

> Many figures courtesy Kevin Murphy's textbook, Machine Learning: A Probabilistic Perspective

Gaussian Mixture Models



Mixture models describe a single, "flat" dataset.

Probabilistic Mixture Models



Mixture models describe a single, "flat" dataset.

Collections of Mixture Models

- Many applications involve multiple "groups" of data
 - Multiple documents in a text corpus
 - Multiple images in a photo repository
 - Multiple users with their own spam filtering decisions
 - Multiple hospitals in a clinical trial
 - Multiple companies in a financial market
- How can we jointly model this data?
- Lumping into single large dataset ignores group differences
- Modeling groups independently can be ineffective, especially when limited data about any one group
- Hierarchical Bayesian models share between groups

Multiple Gaussian Mixtures



Use data to learn set of shared mixture identities, and their frequencies across groups



Each group has its own weight on shared mixture parameters:

 $\pi_j \sim \operatorname{Dir}(\alpha)$

Each observation comes from some mixture component:

$$p(x_{ji} \mid \pi_j, \theta_1, \dots, \theta_K) = \sum_{k=1}^K \pi_{jk} f(x_{ji} \mid \theta_k)$$

Admixture or Topic Models: Multiple Multinomial Mixtures



Latent Dirichlet Allocation (LDA)

Example Data for a Topic Model

Poisoning by ice-cream.

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

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

RNA Editing and the Evolution of Parasites

Larry Simpson and Dmitri A. Maslov

 ${f T}$ he kinetoplastid flagellates, together with their sister group of euglenoids, represent the earliest extant lineage of eukaryotc organisms containing mitochondria (1). Within the kinetoplastids, there are two major groups, the poorly studied bodonidscryptobiids, which consist of both free-living and parasitic cells, and the better known trypanosomatids, which are obligate parasites (2). Perhaps because of the antiquity of the

trypanosomatid lineage, these cells possess several unique genetic fea tures (see accompanying Per-spective by Nilsen)-one of

which is RNA editing of mi-5'-Edited cryptogene tochondrial transcripts. This RNA editing function (3-7) creates open reading frames in "cryptogenes" by insertion (or occasional deletion) of uridine (U) residues at a few specific sites within the cod ing region of an mRNA (5'editing) or at multiple spe-cific sites throughout the mRNA (pan-editing). The Recombination

quenc Ċrithi tral, but there is disagreement on the nacent ture of the primary parasitic host. The "innucle vertebrate first" model (10, 11) states that as an the initial parasitism was in the gut of pretical Cambrian invertebrates. Coevolution of Trype parasite and host would have led to a wide the F distribution of trypanosomatids in insects by th and leeches. In this theory, digenetic life fish r cycles (alternating invertebrate and vertetutes branc

arthr

netic

would

the al

pothe

mitoc



Chaotic Beetles

Charles Godfray and Michael Hassell

curious geometric objects with

SCIENCE • VOL. 275 • 17 IANUARY 1997

Ecologists have known since the pioneering convincing evidence to date of work of May in the mid-1970s (1) that the complex dynamics and chaos population dynamics of animals and plants in a biological population-of the flour beetle, Tribolium can be exceedingly complex. This complexity arises from two sources: The tangled web of interactions that constitute any natural castaneum (see figure). It has proven extremely difcommunity provide a myriad of different ficult to demonstrate complex pathways for species to interact, both di dynamics in populations in the field. By its very nature, a cha-otically fluctuating population rectly and indirectly. And even in isolated populations the nonlinear feedback processes present in all natural populations can result in complex dynamic behavior. Natural will superficially resemble a stable or cyclic population bufpopulations can show persistent oscillatory feted by the normal random permamics and chaos, the latter characterized turbations experienced by all species. Given a long enough by extreme sensitivity to initial conditions. If time series, diagnostic tools such chaotic dynamics were common in nature, then this would have important ramififrom nonlinear mathematics can be used to identify the tellcations for the management and conservation of natural resources. On page 389 of this tale signatures of chaos. In phase issue, Costantino et al. (2) provide the most space, chaotic trajectories come to lie on "strange attractors."

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 fractal structure and hence noninteger dimension. As they

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 estimating the Liapunov expo nent, which is positive for cha-



populations and then compar their predictions with the dy-The flour beetle, Tribo-lium castaneum, exhibits namics in the field. This tech nique has been gaining popu chaotic population dv larity in recent years, helped by namice when the amoun statistical advances in pa in a mathematical model. rameter estimation. Good ex

323

An alternative approach is

parameterize population

- Our data are the pages Science from 1880-2002 (from JSTOR)
- No reliable punctuation, meta-data, or references.
- Note: this is just a subset of JSTOR's archive.

Example Output: 4 Topics

human genome dna genetic genes sequence gene molecular sequencing map information genetics mapping project sequences

evolution evolutionary species organisms life origin biology groups phylogenetic living diversity group new two common

disease host bacteria diseases resistance bacterial new strains control infectious malaria parasite parasites united tuberculosis

computer models information data computers system network systems model parallel methods networks software new simulations

Columns sorted by probability of word given topic.

LDA: Intuition Seeking Life's Bare (Genetic) Necessities

Haemophilus

genome

COLD SPRING HARBOR, NEW YORK— How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms

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

Although the numbers don't match precisely, those predictions

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

SCIENCE • VOL. 272 • 24 MAY 1996

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center

lecular 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.

Every document discusses a mixture of multiple topics.

LDA: Generative Model Seeking Life's Bare (Genetic) Necessities



SCIENCE • VOL. 272 • 24 MAY 1996

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

LDA: Graphical Model



LDA: Graphical Model



- 1 Draw each topic $\beta_i \sim \text{Dir}(\eta)$, for $i \in \{1, \ldots, K\}$.
- 2 For each document:
 - **1** Draw topic proportions $\theta_d \sim \text{Dir}(\alpha)$.
 - **2** For each word:
 - Draw Z_{d,n} ~ Mult(θ_d).
 Draw W_{d,n} ~ Mult(β_{Zd,n}).

Geometry of Topic Models



- Documents are multinomial distributions over some predefined vocabulary of (tens of thousands) of words
- Topics are multinomial distributions on same vocabulary
- Generative model: Each document is (nearly) a convex combination of the topic distributions

LDA: Learning via EM Algorithm?



- M-Step: Maximize likelihood bound with respect to global topic usage distribution α and topic-word distributions β_k
- **E-Step:** Find posterior distribution of document-specific topic frequencies θ_d and word token assignments z_{dn}
- **Problem:** Posterior is intractable for large N
- Variational methods: Create a looser but more tractable loglikelihood bound by constraining form of posterior approx.
- Alternative: Markov Chain Monte Carlo (MCMC)

Uses of Monte Carlo Methods
$$z^{(\ell)} \sim p(z)$$
 $\mathbb{E}[f] = \int f(z)p(z) dz$

- Basic goals: Sampling or estimation of expectations
- Instead of learning by optimization, do simulation
- Given estimated parameters for some statistical model, quantitatively or qualitatively assess accuracy of fit
- Parameter estimation when closed forms unavailable
- Parameter estimation for models with hidden "nuisance" variables (alternative to the EM algorithm)
- General approach to applying computational resources to solve statistical learning problems...



Chaotic dynamical systems generate sequences of pseudo-random numbers approximately distributed uniformly on [0,1]

Monte Carlo Estimators

$$\mathbb{E}_{p}[f(x)] = \int_{\mathcal{X}} f(x)p(x) \, dx \qquad \{x^{(\ell)}\}_{\ell=1}^{L} \quad \begin{array}{l} \text{independent} \\ \text{samples} \end{array}$$
$$\approx \frac{1}{L} \sum_{\ell=1}^{L} f(x^{(\ell)}) = \mathbb{E}_{\tilde{p}}[f(x)] \qquad \qquad \tilde{p}(x) = \frac{1}{L} \sum_{\ell=1}^{L} \delta(x, x^{(\ell)})$$

Good properties if *L* **sufficiently large:**

- Unbiased for any sample size
- Variance inversely proportional to sample size (and independent of dimension of space)
- Weak law of large numbers
- Strong law of large numbers
- **Problem:** Drawing samples from complex distributions...

Alternatives for hard problems:

- Importance sampling
- Markov chain Monte Carlo (MCMC)