

DSC250: Advanced Data Mining

Machine Learning Basics

Zhiting Hu

Lecture 5, Jan 21, 2025

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Outline

- Probability
 - Bayes' rule
 - Exponential family
 - Probabilistic graphical models
 - Entropy, KL divergence, cross entropy
- Functional derivatives
- Practice: MLE vs Maximum entropy

High dimensional distributions

Probabilistic graphical models (PGMs) are about representing probability distributions over random variables

$$p(X) \equiv p(X_1, \dots, X_n)$$

Assume $x_i \in \{0, 1\}^n$

Naively, since there are 2^n possible assignments to X_1, \dots, X_n , can represent this distribution completely using $2^n - 1$ numbers, but quickly becomes intractable for large n

PGMs are methods to represent these distributions more compactly, by exploiting *conditional independence*

Recap: Bayesian networks (directed PGMs)

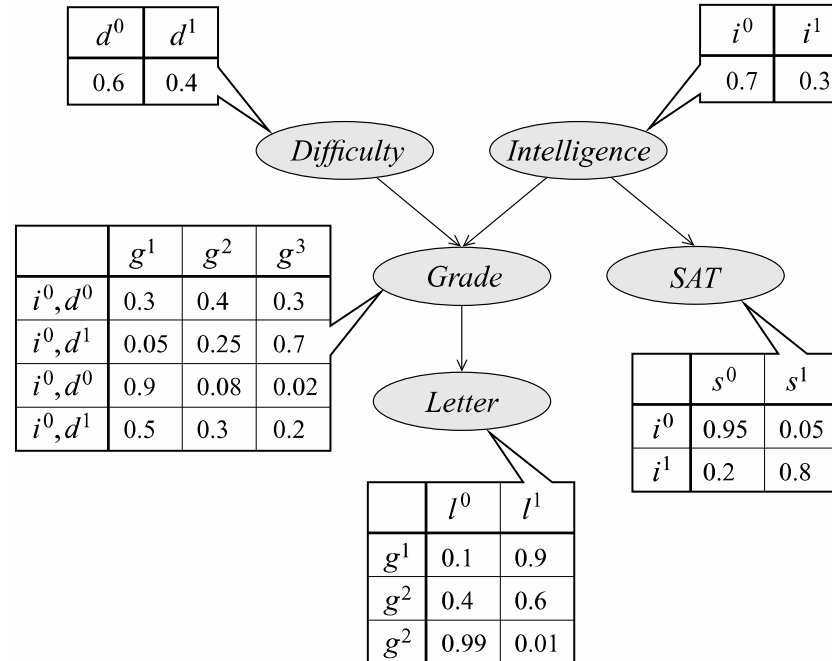
- A **Bayesian network** is specified by a directed *acyclic* graph $G = (V, E)$ with:
 - ① One node $i \in V$ for each random variable X_i
 - ② One conditional probability distribution (CPD) per node, $p(x_i \mid \mathbf{x}_{\text{Pa}(i)})$, specifying the variable's probability conditioned on its parents' values
- Corresponds 1-1 with a particular factorization of the joint distribution:

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i \mid \mathbf{x}_{\text{Pa}(i)})$$

- Powerful framework for designing *algorithms* to perform probability computations

Recap: Example

- Consider the following Bayesian network:



- What is its joint distribution?

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i \mid \mathbf{x}_{\text{Pa}(i)})$$

$$p(d, i, g, s, l) = p(d)p(i)p(g \mid i, d)p(s \mid i)p(l \mid g)$$

Entropy

Entropy

- Shannon entropy $H(p) = - \sum_x p(x) \log p(x)$
 - The average level of "information", "surprise", or "uncertainty" inherent to the variable x 's possible outcomes

KL Divergence

- Kullback-Leibler (KL) divergence: measures the closeness of two distributions $p(\mathbf{x})$ and $q(\mathbf{x})$

$$\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

- a.k.a. Relative entropy
- $\text{KL} \geq 0$ (Jensen's inequality)
- **Questions:**
 - If q is high and p is high in a region, then KL divergence is _____ in this region.
 - If q is high and p is low in a region, then KL divergence is _____ in this region.
 - If q is low in a region, then KL divergence is _____ in this region.

KL Divergence

- Kullback-Leibler (KL) divergence: measures the closeness of two distributions $p(\mathbf{x})$ and $q(\mathbf{x})$

$$\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

- a.k.a. Relative entropy
- $\text{KL} \geq 0$ (Jensen's inequality)
- Intuitively:
 - If q is high and p is high, then we are happy (i.e. low KL divergence)
 - If q is high and p is low then we pay a price (i.e. high KL divergence).
 - If q is low then we don't care (i.e. also low KL divergence, regardless of p)
- not a true "distance":
 - not commutative (symmetric) $\text{KL}(p \parallel q) \neq \text{KL}(q \parallel p)$
 - doesn't satisfy triangle inequality

KL Divergence

- Kullback-Leibler (KL) divergence: measures the closeness of two distributions $p(\mathbf{x})$ and $q(\mathbf{x})$

$$\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

- a.k.a. Relative entropy

- Maximum likelihood estimation (MLE) $\min_{\theta} - \mathbb{E}_{\mathbf{x} \sim \tilde{p}(\mathbf{x})} \left[\log p_{\theta}(\mathbf{x}) \right]$

- **Question:** Show that MLE is minimizing the KL divergence between the empirical data distribution and the model distribution

KL Divergence

- Kullback-Leibler (KL) divergence: measures the closeness of two distributions $p(\mathbf{x})$ and $q(\mathbf{x})$

$$\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

- a.k.a. Relative entropy

- Maximum likelihood estimation (MLE) $\min_{\theta} - \mathbb{E}_{\mathbf{x} \sim \tilde{p}(\mathbf{x})} \left[\log p_{\theta}(\mathbf{x}) \right]$

- **Question:** Show that MLE is minimizing the KL divergence between the empirical data distribution and the model distribution

$$\text{KL}(\tilde{p}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) = -\mathbb{E}_{\tilde{p}(\mathbf{x})} [\log p_{\theta}(\mathbf{x})] + H(\tilde{p}(\mathbf{x}))$$



Cross entropy

Key Takeaways

- Probability $p(\mathbf{x})$

- Bayes' rule
$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{y})p(\mathbf{y})}{p(\mathbf{x})}$$
 - prior, posterior

- Exponential family:
 - Gaussian, multinomial, categorical, ...

- Probabilistic graphical models: Bayesian networks

- KL Divergence
 - relation to Cross-entropy

$$\text{KL}(q(\mathbf{x}) || p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

Functional Derivatives (Optional)

Functional derivative

- $\nabla_q - \mathbb{H}(q) = \log q + 1$
- Functional $F(y)$: an operator that takes a function $y(x)$ and returns an output value F
- Functional derivative (aka, variational derivative): relates a change in a Functional $F(y)$ to a change in the function y

Functional derivative

- Recall the conventional derivative $\frac{dy}{dx}$
 - Taylor expansion

$$y(x + \epsilon) = y(x) + \frac{dy}{dx}\epsilon + O(\epsilon^2)$$

- Functional derivative
 - How much a functional $F[y]$ changes when we make a small change $\epsilon\eta(x)$ to the function $y(x)$

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2)$$

- A function $y(x)$ that maximizes (or minimizes) a functional $F[y]$ must satisfy

$$\frac{\delta F}{\delta y(x)} = 0 \text{ for all } x$$

Functional derivative

$$F[y + \epsilon \eta(x)] = F[y(x)] + \epsilon \int \frac{\partial G}{\partial y} \eta(x) dx + O(\epsilon^2)$$

- Consider a functional that is defined by an integral over a function $G(y, x)$

- Ex.1, $-\mathbb{H}(q) = \int q(x) \log q(x) dx$
 - $G = q(x) \log q(x)$

- Consider variations in the function $y(x)$,

$$F[y + \epsilon \eta(x)] = F[y(x)] + \epsilon \int \frac{\partial G}{\partial y} \eta(x) dx + O(\epsilon^2)$$

Functional derivative

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2)$$

- Consider a functional that is defined by an integral over a function $G(y, x)$

- Ex.1, $-\mathbb{H}(q) = \int q(x) \log q(x) dx$
 - $G = q(x) \log q(x)$

- Consider variations in the function $y(x)$,

$$F[y + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\partial G}{\partial y} \eta(x) dx + O(\epsilon^2)$$

Practice: Maximum likelihood vs Maximum Entropy (Optional)

Supervised Maximum Likelihood

- Model to be learned $p_{\theta}(\mathbf{x})$
- Observe full data $\mathcal{D} = \{ \mathbf{x}^* \}$
 - i.i.d: independent, identically distributed
- Maximum Likelihood Estimation (MLE)
 - The most classical learning algorithm

$$\min_{\theta} - \mathbb{E}_{\mathbf{x}^* \sim \mathcal{D}} \left[\log p_{\theta}(\mathbf{x}^*) \right]$$

- MLE is closely connected to the Maximum Entropy (MaxEnt) principle

Recap: Exponential Family

- A distribution

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = h(\mathbf{x}) \exp\{\boldsymbol{\theta} \cdot T(\mathbf{x})\} / Z(\boldsymbol{\theta})$$

is an exponential family distribution

- $\boldsymbol{\theta} \in R^d$: natural (canonical) parameter
 - $T(\mathbf{x}) \in R^d$: sufficient statistics, features of data \mathbf{x}
 - $Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}, y} h(\mathbf{x}) \exp\{\boldsymbol{\theta} \cdot T(\mathbf{x})\}$: normalization factor
- Examples: Bernoulli, multinomial, Gaussian, Poisson, gamma,...

Maximum Likelihood for Exponential Family

$m(\mathbf{x})$: the number of times \mathbf{x} is observed in D

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) &= \sum_{\mathbf{x}} m(\mathbf{x}) \log p(\mathbf{x} | \boldsymbol{\theta}) \\ &= \sum_{\mathbf{x}} m(\mathbf{x}) \left(\sum_i \theta_i T_i(\mathbf{x}) - \log Z(\boldsymbol{\theta}) \right) \\ &= \sum_{\mathbf{x}} m(\mathbf{x}) \sum_i \theta_i T_i(\mathbf{x}) - N \log Z(\boldsymbol{\theta})\end{aligned}$$

- Take gradient and set to 0

$$\begin{aligned}\frac{\partial}{\partial \theta_i} \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) &= \sum_{\mathbf{x}} m(\mathbf{x}) T_i(\mathbf{x}) - N \frac{\partial}{\partial \theta_i} \log Z(\boldsymbol{\theta}) \\ &= \sum_{\mathbf{x}} m(\mathbf{x}) T_i(\mathbf{x}) - N \sum_{\mathbf{x}} p(\mathbf{x} | \boldsymbol{\theta}) T_i(\mathbf{x})\end{aligned}$$

$$\Rightarrow \left[\sum_{\mathbf{x}} p(\mathbf{x} | \boldsymbol{\theta}) T_i(\mathbf{x}) \right] = \sum_{\mathbf{x}} \frac{m(\mathbf{x})}{N} T_i(\mathbf{x}) = \left[\sum_{\mathbf{x}} \tilde{p}(\mathbf{x} | \boldsymbol{\theta}) T_i(\mathbf{x}) \right]$$

At MLE, the expectations of the sufficient statistics under the model must match empirical feature average

Maximum Entropy (MaxEnt)

- Given \mathcal{D} , to estimate $p(\mathbf{x})$
- We can approach the problem from an entirely different point of view. Begin with some fixed feature expectations:

$$\sum_{\mathbf{x}} p(\mathbf{x})T_i(\mathbf{x}) = \sum_{\mathbf{x}} \frac{m(\mathbf{x})}{N}T_i(\mathbf{x}) := \alpha_i$$

- There may exist many distributions which satisfy them. Which one should we select?
 - MaxEnt principle: the most uncertain or flexible one, i.e., the one with maximum entropy
- This yields a new optimization problem:
 - This is a variational definition of a distribution!

$$\max_p H(p(\mathbf{x})) = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x})$$

$$\text{s.t. } \sum_{\mathbf{x}} p(\mathbf{x})T_i(\mathbf{x}) = \alpha_i$$

$$\sum_{\mathbf{x}} p(\mathbf{x}) = 1$$

Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(\mathbf{x})} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left(\sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left(\sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(\mathbf{x})} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left(\sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left(\sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

$$\frac{\partial L}{\partial p(\mathbf{x})} = 1 + \log p(\mathbf{x}) - \sum_i \theta_i T_i(\mathbf{x}) - \mu$$

$$p^*(\mathbf{x}) = e^{\mu-1} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\}$$

$$Z(\boldsymbol{\theta}) = e^{\mu-1} = \sum_{\mathbf{x}} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\} \quad \left(\text{since } \sum_{\mathbf{x}} p^*(\mathbf{x}) = 1 \right)$$

$$p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(\mathbf{x})} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left(\sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left(\sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

$$\frac{\partial L}{\partial p(\mathbf{x})} = 1 + \log p(\mathbf{x}) - \sum_i \theta_i T_i(\mathbf{x}) - \mu$$

$$p^*(\mathbf{x}) = e^{\mu-1} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\}$$

$$Z(\boldsymbol{\theta}) = e^{\mu-1} = \sum_{\mathbf{x}} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\} \quad \left(\text{since } \sum_{\mathbf{x}} p^*(\mathbf{x}) = 1 \right)$$

$$p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

- So feature constraints + MaxEnt \Rightarrow **exponential family**.
- Problem is strictly convex w.r.t. $p(\mathbf{x})$, so solution is unique.

Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(\mathbf{x})} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left(\sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left(\sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

$$p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

plug $p(\mathbf{x} | \boldsymbol{\theta})$ back into L , and since $\sum_{\mathbf{x}} \frac{m(\mathbf{x})}{N} T_i(\mathbf{x}) := \alpha_i$:

$$\max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \sum_{\mathbf{x}} m(\mathbf{x}) \sum_i \theta_i T_i(\mathbf{x}) - N \log Z(\boldsymbol{\theta})$$

- Recovers precisely the MLE problem of exponential family

- So feature constraints + MaxEnt \Rightarrow **exponential family**.
- Problem is strictly convex w.r.t. $p(\mathbf{x})$, so solution is unique.

Constraints from Data

- We have seen a case of **convex duality**:
 - In one case, we assume exponential family and show that Maximum Likelihood implies model expectations must match empirical expectations.
 - In the other case, we assume model expectations must match empirical feature counts and show that MaxEnt implies exponential family distribution.

A more general MaxEnt problem

$$\min_p \text{KL}(p(\mathbf{x}) \| h(\mathbf{x}))$$

$$\stackrel{\text{def}}{=} \sum_{\mathbf{x}} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{h(\mathbf{x})} = -\text{H}(p) - \sum_{\mathbf{x}} p(\mathbf{x}) \log h(\mathbf{x})$$

$$\text{s.t.} \quad \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) = \alpha_i$$

$$\sum_{\mathbf{x}} p(\mathbf{x}) = 1$$

$$\Rightarrow p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(\mathbf{x}) \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

Summary

- Maximum entropy is dual to maximum likelihood of exponential family distributions
- This provides an alternative view of the problem of fitting a model into data:
 - The data instances in the training set are treated as constraints, and the learning problem is treated as a constrained optimization problem.
 - We'll revisit this optimization-theoretic view of learning repeatedly in the future!

$$\begin{aligned} \max_p \quad & H(p(\mathbf{x})) = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) \\ \text{s.t.} \quad & \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) = \alpha_i \\ & \sum_{\mathbf{x}} p(\mathbf{x}) = 1 \end{aligned}$$

Key Takeaways

- Probability
 - Bayes' rule
 - Exponential family
 - Probabilistic graphical models: Bayesian networks
 - KL divergence
- Functional derivative
- Convex duality between MLE and MaxEnt (optional)

Topic Models

Outline

- Representations of Text and Topics
- Topic Model v1: Multinomial Mixture Model
- Topic Model v2: Probabilistic Latent Semantic Analysis (pLSA)
- Topic Model v3: Latent Dirichlet Allocation (LDA)

Slides adapted from:

- Y. Sun, CS 247: Advanced Data Mining
- M. Gormley, 10-701 Introduction to Machine Learning

Motivation

Suppose you're given a massive corpora and asked to carry out the following tasks

- **Organize** the documents into **thematic categories**
- **Describe** the evolution of those categories **over time**
- Enable a domain expert to **analyze and understand** the content
- Find **relationships** between the categories
- Understand how **authorship** influences the content



Motivation

Suppose you're given a massive corpora and asked to carry out the following tasks

- **Organize** the documents into **thematic categories**
- **Describe** the evolution of those categories **over time**
- Enable a domain expert to **analyze and understand** the content
- Find **relationships** between the categories
- Understand how **authorship** influences the content

Topic Modeling:

A method of (usually unsupervised) discovery of latent or hidden structure in a corpus

- Applied primarily to text corpora, but **techniques are more general**
- Provides a **modeling toolbox**
- Has prompted the exploration of a variety of new **inference methods** to accommodate **large-scale datasets**

Topic Modeling: Examples

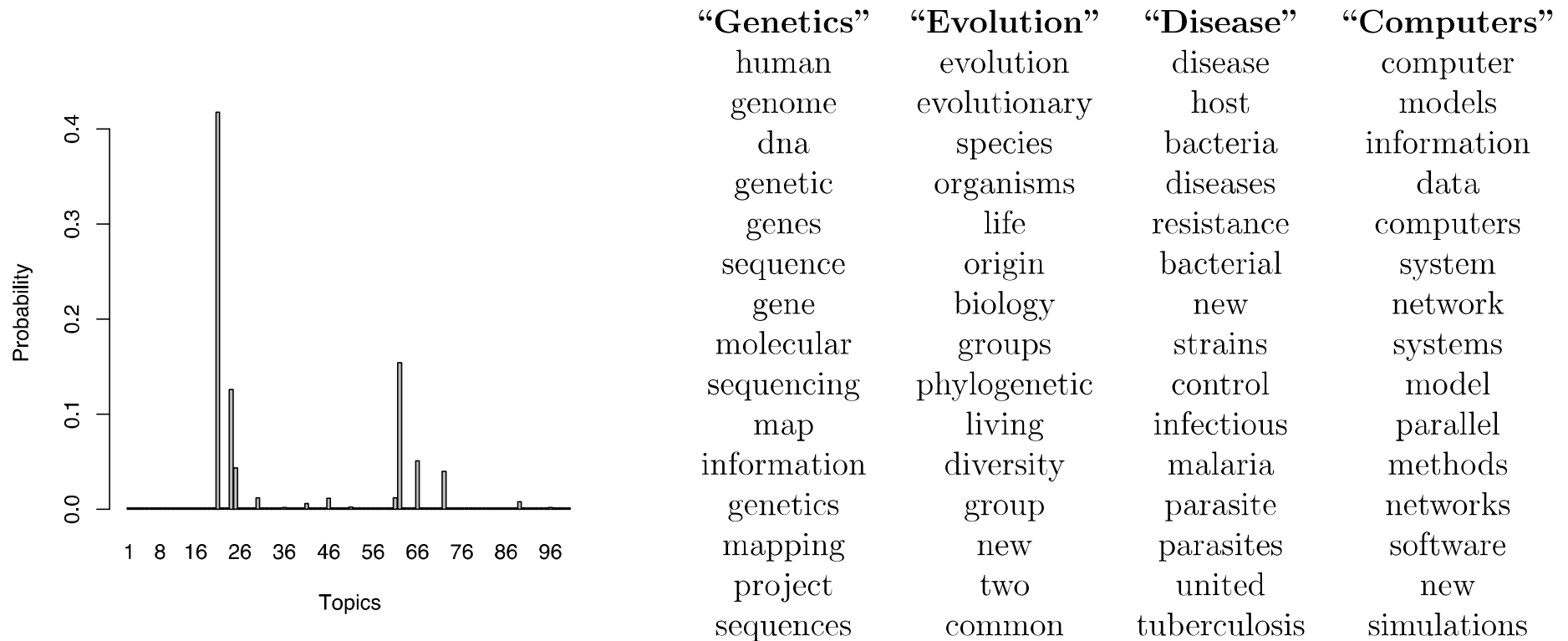


Figure from (Blei, 2011), shows topics and top words learned automatically from reading 17,000 Science articles

Topic Modeling: Examples

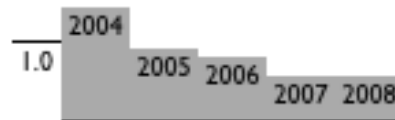
Dirichlet multinomial regression (DMR) topic model & CML (Mimno & McCallum, 2008)

Topic 0 [0.152]



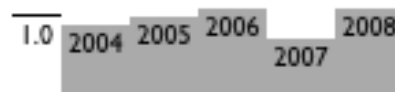
problem, optimization, problems, convex, convex optimization, linear, semidefinite programming, formulation, sets, constraints, proposed, margin, maximum margin, optimization problem, linear programming, programming, procedure, method, cutting plane, solutions

Topic 54 [0.051]



decision trees, trees, tree, decision tree, decision, tree ensemble, junction tree, decision tree learners, leaf nodes, arithmetic circuits, ensembles modts, skewing, ensembles, anytime induction decision trees, trees trees, random forests, objective decision trees, tree learners, trees grove, candidate split

Topic 99 [0.066]



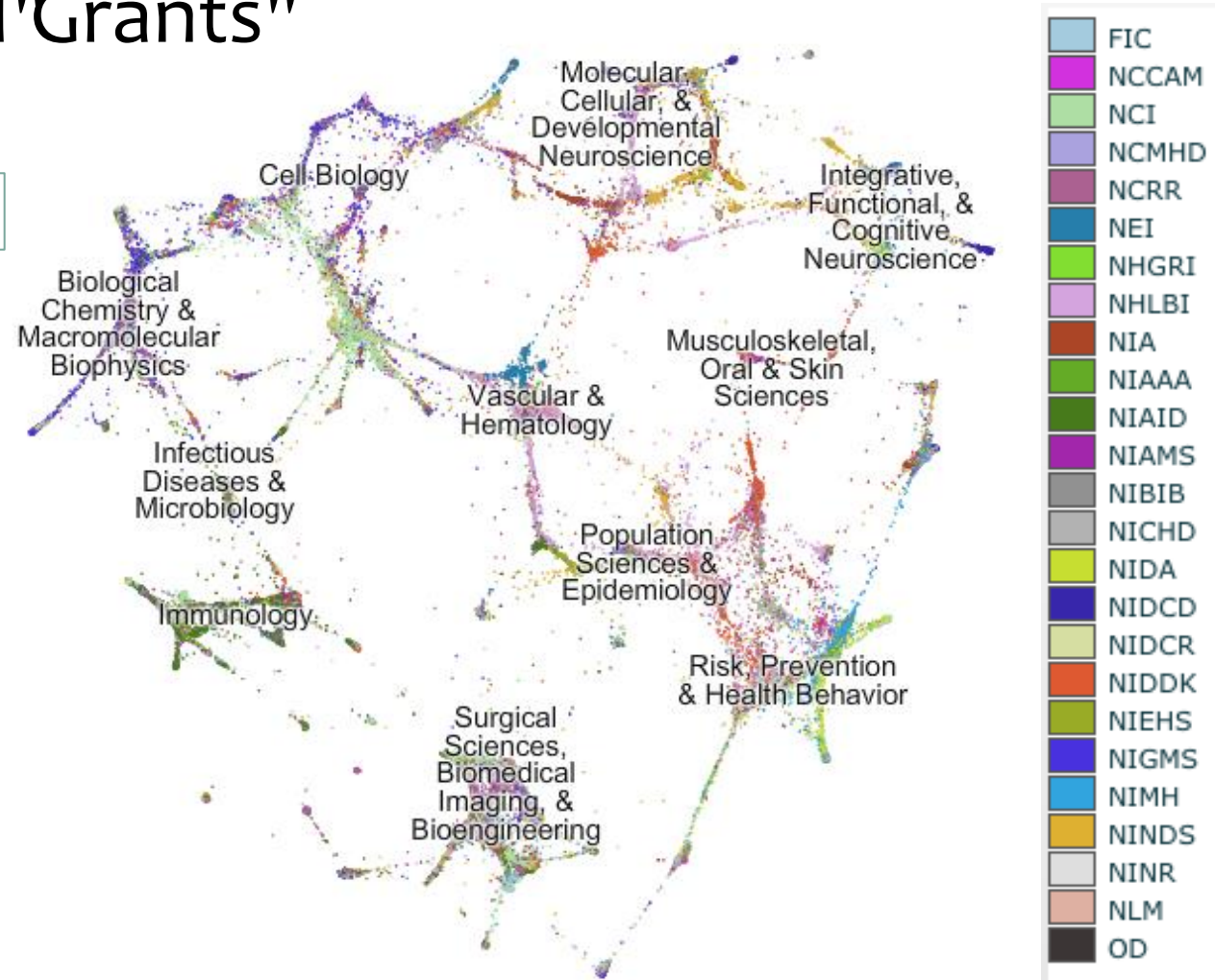
inference, approximate inference, exact inference, markov chain, models, approximate, gibbs sampling, variational, bayesian, variational inference, variational bayesian, approximation, sampling, methods, exact, bayesian inference, dynamic bayesian, process, mcmc, efficient

<http://www.cs.umass.edu/~mimno/icml100.html>

Topic Modeling: Examples

- Map of NIH Grants

(Talley et al., 2011)

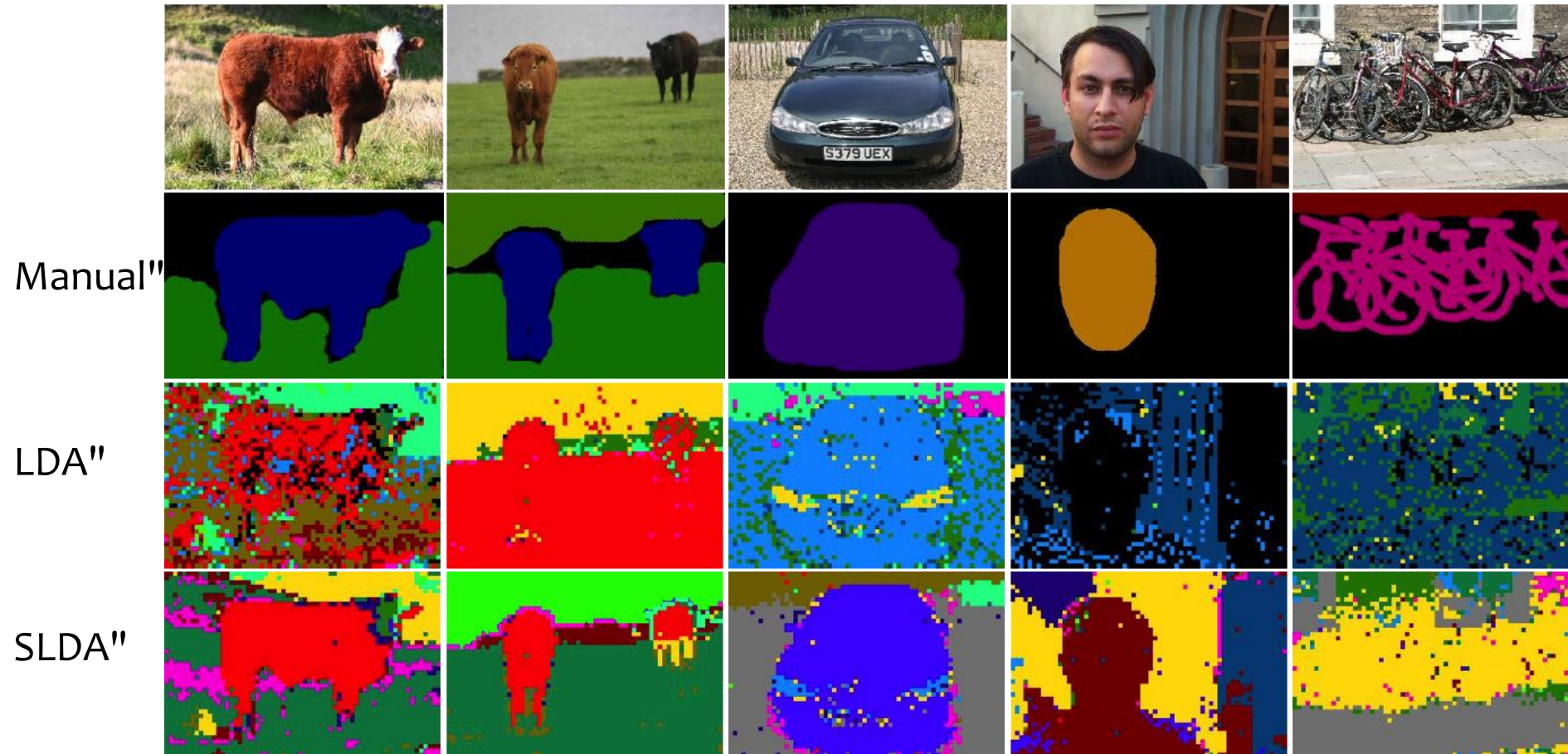


<https://app.nihmaps.org/>

Other Applications of Topic Models

- Spatial "LDA"

(Wang & Grimson, 2007)



Other Applications of Topic Models

- Word Sense Induction

(Brody & Lapata, 2009)

Senses of <i>drug</i> (WSJ)
1. U.S., administration, federal, against, war, dealer
2. patient, people, problem, doctor, company, abuse
3. company, million, sale, maker, stock, inc.
4. administration, food, company, approval, FDA

Senses of <i>drug</i> (BNC)
1. patient, treatment, effect, anti-inflammatory
2. alcohol, treatment, patient, therapy, addiction
3. patient, new, find, effect, choice, study
4. test, alcohol, patient, abuse, people, crime
5. trafficking, trafficker, charge, use, problem
6. abuse, against, problem, treatment, alcohol
7. people, wonder, find, prescription, drink, addict
8. company, dealer, police, enforcement, patient

- Selectional Preference

(Ritter et al., 2010)

Topic <i>t</i>	Arg1	Relations which assign highest probability to <i>t</i>	Arg2
18	The residue - The mixture - The reaction mixture - The solution - the mixture - the reaction mixture - the residue - The reaction - the solution - The filtrate - the reaction - The product - The crude product - The pellet - The organic layer - Thereto - This solution - The resulting solution - Next - The organic phase - The resulting mixture - C.)	was treated with, is treated with, was poured into, was extracted with, was purified by, was diluted with, was filtered through, is dissolved in, is washed with	EtOAc - CH2Cl2 - H2O - CH.sub.2Cl.sub.2 - H.sub.2O - water - MeOH - NaHCO3 - Et2O - NHCl - CHCl.sub.3 - NHCl - drop-wise - CH2Cl.sub.2 - Celite - Et.sub.2O - Cl.sub.2 - NaOH - AcOEt - CH2C12 - the mixture - saturated NaHCO3 - SiO2 - H2O - N hydrochloric acid - NHCl - preparative HPLC - to0 C

Text Data

- Word/term
- Document
 - A sequence of words
- Corpus
 - A collection of documents



Represent a Document

- Most common way: Bag-of-Words
 - Ignore the order of words
 - keep the count

c1: *Human machine interface* for Lab ABC computer applications
c2: A *survey* of user opinion of computer system response time
c3: The *EPS* user interface management system
c4: System and human system engineering testing of *EPS*
c5: Relation of user-perceived *response time* to error measurement

m1: The generation of random, binary, unordered *trees*
m2: The intersection *graph* of paths in *trees*
m3: *Graph minors* IV: Widths of *trees* and well-quasi-ordering
m4: *Graph minors*: A *survey*



	c1	c2	c3	c4	c5	m1	m2	m3	m4
<i>human</i>	1	0	0	1	0	0	0	0	0
<i>interface</i>	1	0	1	0	0	0	0	0	0
<i>computer</i>	1	1	0	0	0	0	0	0	0
<i>user</i>	0	1	1	0	1	0	0	0	0
<i>system</i>	0	1	1	2	0	0	0	0	0
<i>response</i>	0	1	0	0	1	0	0	0	0
<i>time</i>	0	1	0	0	1	0	0	0	0
<i>EPS</i>	0	0	1	1	0	0	0	0	0
<i>survey</i>	0	1	0	0	0	0	0	0	1
<i>trees</i>	0	0	0	0	0	1	1	1	0
<i>graph</i>	0	0	0	0	0	0	1	1	1
<i>minors</i>	0	0	0	0	0	0	0	1	1

Vector space model

Represent a Document

- Represent the doc as a vector where each entry corresponds to a different word and the number at that entry corresponds to how many times that word was present in the document (or some function of it)
 - Number of words is huge
 - Select and use a smaller set of words that are of interest
 - E.g. uninteresting words: 'and', 'the' 'at', 'is', etc. These are called stop-words
 - Stemming: remove endings. E.g. 'learn', 'learning', 'learnable', 'learned' could be substituted by the single stem 'learn'
 - Other simplifications can also be invented and used
 - The set of different remaining words is called dictionary or vocabulary. Fix an ordering of the terms in the dictionary so that you can operate them by their index.
 - Can be extended to bi-gram, tri-gram, or so

Limitations of Bag-of-Words

- Dimensionality
 - High dimensionality
- Sparseness
 - Most of the entries are zero
- Shallow representation
 - The vector representation does not capture semantic relations between words

Ex: "Tom loves Kate."

Topic Models

- Topic modeling
 - Get topics automatically from a corpus
 - Assign documents to topics automatically
- Most frequently used topic models
 - pLSA
 - LDA

“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

Questions?