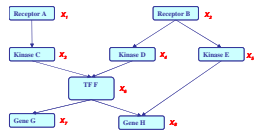


Probabilistic Graphical Models

Learning one-node GM

Eric Xing

Lecture 5, September 23, 2009



Reading:

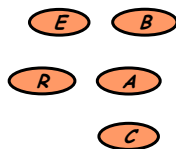
© Eric Xing @ CMU, 2005-2009

1

Learning Graphical Models

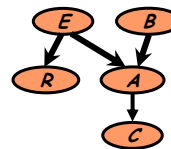
The goal:

Given set of independent samples (*assignments* of random variables), find the **best** (the most likely?) Bayesian Network (both DAG and CPDs)



(B,E,A,C,R)=(T,F,F,T,F)
 (B,E,A,C,R)=(T,F,T,T,F)

 (B,E,A,C,R)=(F,T,T,T,F)



Structural learning

F	B	P(A E,B)	
e	\underline{b}	0.9	0.1
\underline{e}	b	0.2	0.8
\underline{e}	\underline{b}	0.9	0.1
e	\underline{b}	0.01	0.99

Parameter learning

© Eric Xing @ CMU, 2005-2009

2

Learning Graphical Models

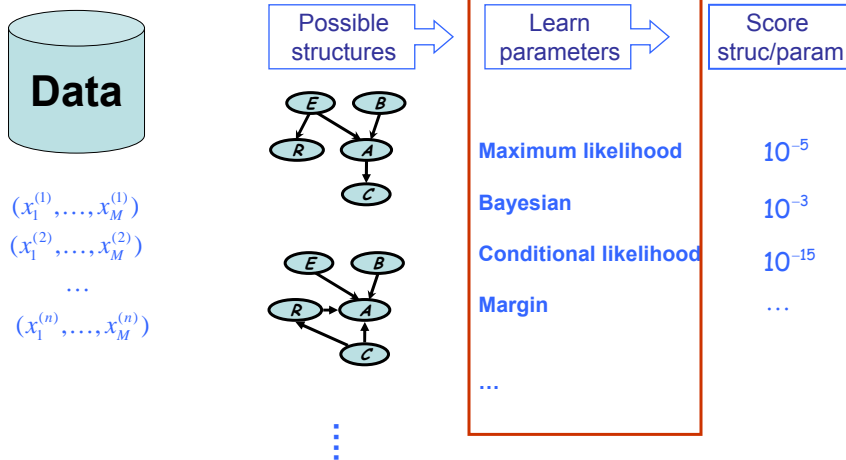


- Scenarios:
 - completely observed GMs
 - directed
 - undirected
 - partially or unobserved GMs
 - directed
 - undirected (an open research topic)
- Estimation principles:
 - Maximal likelihood estimation (MLE)
 - Bayesian estimation
 - Maximal conditional likelihood
 - Maximal "Margin"
 - Maximum entropy
- We use **learning** as a name for the process of **estimating the parameters**, and in some cases, the topology of the network, from data.

© Eric Xing @ CMU, 2005-2009

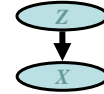
3

Score-based approach



© Eric Xing @ CMU, 2005-2009

4



ML Parameter Est. for completely observed GMs of given structure

- The data:

$$\{(z^{(1)}, x^{(1)}), (z^{(2)}, x^{(2)}), (z^{(3)}, x^{(3)}), \dots, (z^{(N)}, x^{(N)})\}$$



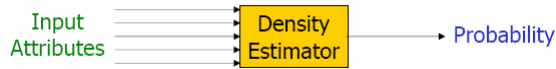
Parameter Learning

- Assume G is known and fixed,
 - from expert design
 - from an intermediate outcome of iterative structure learning
- Goal: estimate from a dataset of N independent, **identically distributed (iid)** training cases $D = \{x^{(1)}, \dots, x^{(N)}\}$.
- In general, each training case $\mathbf{x}^{(n)} = (x_1^{(n)}, \dots, x_M^{(n)})$ is a vector of M values, one per node,
 - the model can be completely observable, i.e., every element in x_n is known (no missing values, no hidden variables),
 - or, partially observable, i.e., $\exists i$, s.t. $x_i^{(n)}$ is not observed.
- **In this lecture we consider learning parameters for a single node.**
 - Often known as “density estimation”

Density Estimation



- A Density Estimator learns a mapping from a set of attributes to a Probability



- Often know as *parameter estimation* if the distribution form is specified
 - Binomial, Gaussian ...
- Four important issues:
 - Nature of the data (iid, correlated, ...)
 - Objective function (MLE, MAP, Margin ...)
 - Algorithm (simple algebra, gradient methods, EM, ...)
 - Evaluation scheme (likelihood on test data, predictability, consistency, ...)

Discrete Distributions



- Bernoulli distribution: $\text{Ber}(p)$

$$P(x) = \begin{cases} 1-p & \text{for } x=0 \\ p & \text{for } x=1 \end{cases} \Rightarrow P(x) = p^x (1-p)^{1-x}$$



- Multinomial distribution: $\text{Mult}(1, \theta)$

- Multinomial (indicator) variable:

$$X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \\ X_6 \end{bmatrix}, \quad \text{where} \quad \begin{aligned} X_j &= [0,1], \quad \text{and} \quad \sum_{j=1, \dots, 6} X_j = 1 \\ X_j &= 1 \text{ w.p. } \theta_j, \quad \sum_{j=1, \dots, 6} \theta_j = 1. \end{aligned}$$



$$\begin{aligned} p(x(j)) &= P(\{X_j = 1, \text{ where } j \text{ index the dice-face}\}) \\ &= \theta_j = \theta_A^{x_A} \times \theta_C^{x_C} \times \theta_G^{x_G} \times \theta_T^{x_T} = \prod_k \theta_k^{x_k} = \theta^x \end{aligned}$$

Continuous Distributions

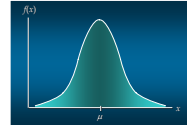
- Uniform Probability Density Function

$$p(x) = \begin{cases} 1/(b-a) & \text{for } a \leq x \leq b \\ 0 & \text{elsewhere} \end{cases}$$



- Normal (Gaussian) Probability Density Function

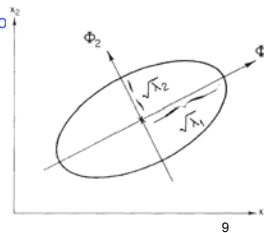
$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



- The distribution is symmetric, and is often illustrated as a bell-shaped curve.
- Two parameters, μ (mean) and σ (standard deviation), determine the location and shape of the distribution.
- The highest point on the normal curve is at the mean, which is also μ .
- The mean can be any numerical value: negative, zero, or positive.

- Multivariate Gaussian

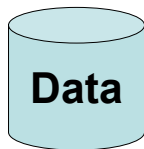
$$p(X; \bar{\mu}, \Sigma) = \frac{1}{(\sqrt{2\pi})^n |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(X - \bar{\mu})^T \Sigma^{-1} (X - \bar{\mu})\right\}$$



© Eric Xing @ CMU, 2005-2009

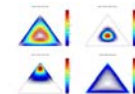
9

Density Estimation Schemes



$(x_1^{(1)}, \dots, x_M^{(1)})$
 $(x_1^{(2)}, \dots, x_M^{(2)})$
 ...
 $(x_1^{(n)}, \dots, x_M^{(n)})$

Learn parameters	Algorithm	Score param
Maximum likelihood	Analytical	10^{-5}
Bayesian	Gradient	10^{-3}
Conditional likelihood	EM	10^{-15}
Margin	Sampling	...
Entropy
...



© Eric Xing @ CMU, 2005-2009

10

Parameter Learning from *iid* Data



- Goal: estimate distribution parameters θ from a dataset of N independent, identically distributed (*iid*), fully observed, training cases

$$\mathcal{D} = \{x^{(1)}, \dots, x^{(n)}\}$$

- Maximum likelihood estimation (MLE)
 1. One of the most common estimators
 2. With iid and full-observability assumption, write $L(\theta)$ as the likelihood of the data:

$$\begin{aligned} L(\theta) &= P(x^{(1)}, x^{(2)}, \dots, x^{(N)}; \theta) \\ &= P(x^{(1)}; \theta)P(x^{(2)}; \theta), \dots, P(x^{(N)}; \theta) \\ &= \prod_{i=1}^N P(x^{(i)}; \theta) \end{aligned}$$

3. pick the setting of parameters most likely to have generated the data we saw:

$$\theta^* = \arg \max_{\theta} L(\theta) = \arg \max_{\theta} \log L(\theta)$$

Example: Bernoulli model



- Data:
 - We observed N *iid* coin tossing: $\mathcal{D}=\{1, 0, 1, \dots, 0\}$



- Representation:

Binary r.v: $x^{(n)} \in \{0,1\}$

- Model: $P(x) = \begin{cases} 1-\theta & \text{for } x=0 \\ \theta & \text{for } x=1 \end{cases} \Rightarrow P(x) = \theta^x (1-\theta)^{1-x}$

- How to write the likelihood of a single observation $x^{(i)}$?

$$P(x^{(i)}) = \theta^{x^{(i)}} (1-\theta)^{1-x^{(i)}}$$

- The likelihood of dataset $\mathcal{D}=\{x_1, \dots, x_N\}$:

$$P(x^{(1)}, x^{(2)}, \dots, x^{(N)} | \theta) = \prod_{i=1}^N P(x^{(i)} | \theta) = \prod_{i=1}^N (\theta^{x^{(i)}} (1-\theta)^{1-x^{(i)}}) = \theta^{\sum_{i=1}^N x^{(i)}} (1-\theta)^{\sum_{i=1}^N (1-x^{(i)})} = \theta^{\#\text{head}} (1-\theta)^{\#\text{tails}}$$

Maximum Likelihood Estimation



- Objective function:

$$\ell(\theta; D) = \log P(D | \theta) = \log \theta^{n_h} (1 - \theta)^{n_t} = n_h \log \theta + (N - n_h) \log(1 - \theta)$$

- We need to maximize this w.r.t. θ
- Take derivatives wrt θ

$$\frac{\partial \ell}{\partial \theta} = \frac{n_h}{\theta} - \frac{N - n_h}{1 - \theta} = 0 \quad \Rightarrow \quad \hat{\theta}_{MLE} = \frac{n_h}{N} \quad \text{or} \quad \hat{\theta}_{MLE} = \frac{1}{N} \sum_i x^{(i)}$$

Frequency as sample mean

- Sufficient statistics

- The counts n_h , where $n_k = \sum_i x^{(i)}$, are sufficient statistics of data D

Overfitting



- Recall that for Bernoulli Distribution, we have

$$\hat{\theta}_{ML}^{head} = \frac{n^{head}}{n^{head} + n^{tail}}$$

- What if we tossed too few times so that we saw zero head?

We have $\hat{\theta}_{ML}^{head} = 0$, and we will predict that the probability of seeing a head next is zero!!!

- Frequentist vs. Bayesian estimate

Bayesian Parameter Estimation



- Treat the distribution parameters θ also as a *random variable*
- The *a posteriori* distribution of θ after seeing the data is:

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)} = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta)p(\theta)d\theta}$$

This is Bayes Rule

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}$$

Bayes, Thomas (1763) An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, 53:370-418



The prior $p(\cdot)$ encodes our prior knowledge about the domain

© Eric Xing @ CMU, 2005-2009

15

Frequentist Parameter Estimation



Two people with different priors $p(\theta)$ will end up with different estimates $p(\theta|D)$.

- Frequentists dislike this “subjectivity”.
- Frequentists think of the parameter as a **fixed, unknown constant**, not a random variable.
- Hence they have to come up with different “objective” **estimators** (ways of computing from data), instead of using Bayes’ rule.
 - These estimators have different properties, such as being “unbiased”, “minimum variance”, etc.
 - The **maximum likelihood estimator**, is one such estimator, which is simple and has good statistical properties.

© Eric Xing @ CMU, 2005-2009

16

Discussion



θ or $p(\theta)$, this is the problem!

Maximum Likelihood Estimation



- The log-likelihood is monotonically related to the likelihood:

$$\ell(\theta; D) = \log p(D | \theta) = \sum_{n=1}^N \log p(x^{(n)} | \theta)$$

- The Idea underlying maximum likelihood estimation (MLE): pick the setting of parameters most likely to have generated the data we saw:

$$\hat{\theta}_{ML} = \arg \max_{\theta} \ell(\theta; D)$$

- Problem of MLE:

- Overfitting: means that "some of the relationships that appear statistically significant are actually just noise. It occurs when the complexity of the statistical model is too great for the amount of data that you have"
- Often the MLE **overfits** the training data, so it is common to maximize a **regularized log-likelihood** instead:
$$\hat{\theta} = \arg \max_{\theta} \ell(\theta; D) - c(\theta)$$
- Insufficient training data can lead to spurious estimator (e.g., certain possible values are not observed due to data sparsity), so it is common to **smooth** the estimated parameter

Being a pragmatic frequentist



- Maximum *a posteriori* (MAP) estimation:

$$\hat{\theta}_{MAP} = \arg \max_{\theta} p(\theta | D) = \arg \max_{\theta} \mathcal{L}(\theta; D) + \log p(\theta)$$

- Smoothing with pseudo-counts
 - Recall that for Binomial Distribution, we have

$$\hat{\theta}_{MLE}^{head} = \frac{n^{head}}{n^{head} + n^{tail}}$$

- What if we tossed too few times so that we saw zero head? We have $\hat{\theta}_{MLE}^{head} = 0$, and we will predict that the probability of seeing a head next is zero!!!

- The rescue:
$$\hat{\theta}_{MLE}^{head} = \frac{n^{head} + n'}{n^{head} + n^{tail} + n'}$$
 But are we still objective?

- Where n' is known as the pseudo- (imaginary) count

© Eric Xing @ CMU, 2005-2009

19

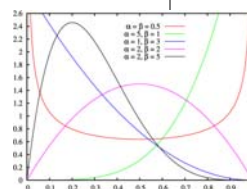
Bayesian estimation for Bernoulli



- Beta distribution:

$$P(\theta; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} = B(\alpha, \beta) \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

- When x is discrete $\Gamma(x+1) = x\Gamma(x) = x!$



- Posterior distribution of θ :

$$P(\theta | x^{(1)}, \dots, x^{(N)}) = \frac{p(x^{(1)}, \dots, x^{(N)} | \theta) p(\theta)}{p(x^{(1)}, \dots, x^{(N)})} \propto \theta^{n_h} (1-\theta)^{n_t} \times \theta^{\alpha-1} (1-\theta)^{\beta-1} = \theta^{n_h + \alpha - 1} (1-\theta)^{n_t + \beta - 1}$$

- Notice the isomorphism of the posterior to the prior,
- such a prior is called a **conjugate prior**
- α and β are hyperparameters (parameters of the prior) and correspond to the number of “virtual” heads/tails (pseudo counts)

© Eric Xing @ CMU, 2005-2009

20

Bayesian estimation for Bernoulli, con'd



- Posterior distribution of θ :

$$P(\theta | x^{(1)}, \dots, x^{(N)}) = \frac{P(x^{(1)}, \dots, x^{(N)} | \theta) P(\theta)}{P(x^{(1)}, \dots, x^{(N)})} \propto \theta^{n_h} (1-\theta)^{n_t} \times \theta^{\alpha-1} (1-\theta)^{\beta-1} = \theta^{n_h+\alpha-1} (1-\theta)^{n_t+\beta-1}$$

- Maximum *a posteriori* (MAP) estimation:

$$\theta_{MAP} = \arg \max_{\theta} \log P(\theta | x^{(1)}, \dots, x^{(N)})$$

- Posterior mean estimation:

$$\theta_{Bayes} = \int \theta p(\theta | D) d\theta = C \int \theta \times \theta^{n_h+\alpha-1} (1-\theta)^{n_t+\beta-1} d\theta = \frac{n_h + \alpha}{N + \alpha + \beta}$$

Bata parameters
can be understood
as pseudo-counts

- Prior strength: $A = \alpha + \beta$

- A can be interperated as the size of an imaginary data set from which we obtain the **pseudo-counts**

© Eric Xing @ CMU, 2005-2009

21

Effect of Prior Strength



- Suppose we have a uniform prior ($\alpha = \beta = 1/2$), and we observe $\vec{n} = (n_h = 2, n_t = 8)$

- Weak prior $A = 2$. Posterior prediction:

$$p(x = h | n_h = 2, n_t = 8, \bar{\alpha} = \bar{\alpha} \times 2) = \frac{1+2}{2+10} = 0.25$$

- Strong prior $A = 20$. Posterior prediction:

$$p(x = h | n_h = 2, n_t = 8, \bar{\alpha} = \bar{\alpha} \times 20) = \frac{10+2}{20+10} = 0.40$$

- However, if we have enough data, it washes away the prior. e.g., $\vec{n} = (n_h = 200, n_t = 800)$. Then the estimates under weak and strong prior are $\frac{1+200}{2+1000}$ and $\frac{10+200}{20+1000}$, respectively, both of which are close to 0.2

© Eric Xing @ CMU, 2005-2009

22

Being a “subjective” Bayesian



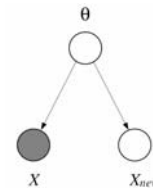
- The prior $p(\cdot)$ encodes our prior knowledge about the domain
 - therefore Bayesian estimation has been criticized for being "subjective"
- Empirical Bayes – fit prior from "training" data

How estimators should be used?



- $\hat{\theta}_{MAP}$ is not Bayesian (even though it uses a prior) since it is a point estimate.
- Consider predicting the future. A sensible way is to combine predictions based on all possible values of θ , weighted by their posterior probability, this is what a **Bayesian** will do:

$$\begin{aligned}
 p(x_{new} | \mathbf{x}) &= \int p(x_{new}, \theta | \mathbf{x}) d\theta \\
 &= \int p(x_{new} | \theta, \mathbf{x}) p(\theta | \mathbf{x}) d\theta \\
 &= \int p(x_{new} | \theta) p(\theta | \mathbf{x}) d\theta
 \end{aligned}$$



- A **frequentist** will typically use a “plug-in” estimator such as ML/MAP:

$$p(x_{new} | \mathbf{x}) = p(x_{new} | \hat{\theta}_{ML}), \quad \text{or, } p(x_{new} | \mathbf{x}) = p(x_{new} | \hat{\theta}_{MAP})$$

- The Bayesian estimate will collapse to MAP for concentrated posterior

Frequentist vs. Bayesian



- This is a “theological” war.
- Advantages of Bayesian approach:
 - Mathematically elegant.
 - Works well when amount of data is much less than number of parameters (e.g., one-shot learning).
 - Easy to do incremental (sequential) learning.
 - Can be used for model selection (max likelihood will always pick the most complex model).
- Advantages of frequentist approach:
 - Mathematically/ computationally simpler.
 - "objective", unbiased, invariant to reparameterization
- As $|D| \rightarrow \infty$, the two approaches become the same:

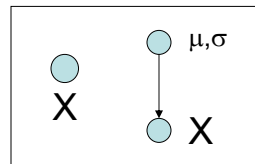
$$p(\theta | D) \rightarrow \delta(\theta, \hat{\theta}_{ML})$$

Simplest GMs: the building blocks



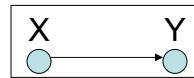
Density estimation

Parametric and nonparametric methods



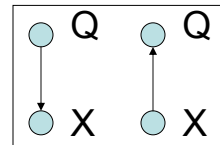
Regression

Linear, conditional mixture, nonparametric



Classification

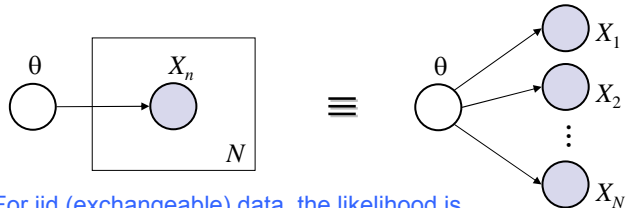
Generative and discriminative approach



Plates



- A plate is a “macro” that allows subgraphs to be replicated



- For iid (exchangeable) data, the likelihood is

$$p(D|\theta) = \prod_n p(x_n|\theta)$$

- We can represent this as a Bayes net with N nodes.
 - The rules of plates are simple: repeat every structure in a box a number of times given by the integer in the corner of the box (e.g. N), updating the plate index variable (e.g. n) as you go.
 - Duplicate every arrow going into the plate and every arrow leaving the plate by connecting the arrows to each copy of the structure.

© Eric Xing @ CMU, 2005-2009

27

Discrete Distributions



- Bernoulli distribution: $\text{Ber}(p)$

$$p(x) = \begin{cases} 1-p & \text{for } x=0 \\ p & \text{for } x=1 \end{cases} \Rightarrow p(x) = p^x(1-p)^{1-x}$$



- Multinomial distribution: $\text{Mult}(1, \theta)$

- Multinomial (indicator) variable:

$$X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \\ X_6 \end{bmatrix}, \quad \text{where} \quad \begin{aligned} X_j &= [0,1], \quad \text{and} \quad \sum_{j=1, \dots, 6} X_j = 1 \\ X_j &= 1 \text{ w.p. } \theta_j, \quad \sum_{j=1, \dots, 6} \theta_j = 1. \end{aligned}$$



$$\begin{aligned} p(x(j)) &= p(\{X_j = 1, \text{ where } j \text{ index the dice-face}\}) \\ &= \theta_A^{x_A} \times \theta_C^{x_C} \times \theta_G^{x_G} \times \theta_T^{x_T} = \prod_k \theta_k^{x_k} = \theta^x \end{aligned}$$

© Eric Xing @ CMU, 2005-2009

28

Discrete Distributions

- Multinomial distribution: $\text{Mult}(n, \theta)$

- Count variable:

$$n = \begin{bmatrix} n_1 \\ \vdots \\ n_K \end{bmatrix}, \quad \text{where } \sum_j n_j = N$$

$$p(n) = \frac{N!}{n_1! n_2! \dots n_K!} \theta_1^{n_1} \theta_2^{n_2} \dots \theta_K^{n_K} = \frac{N!}{n_1! n_2! \dots n_K!} \theta^n$$

"Arts"	"Hedgers"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FEAR	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
DEAL	FEDERAL	CAMERAS	FOUR
MUSICAL	YEAR	WORK	PUBLIC
REST	SPENDING	PARENTS	TEACHER
ACTION	NEW	SLEET	BENNETT
FIRST	FEAR	FAMILY	MANGAT
YORK	PLAN	WELFARE	NAMBY
OPERA	MONEY	MEN	STATE
THEATRE	PROGRAMS	PERCENT	ELEMENTARY
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	STATE

The Wilson Foundation hopes "Education" will give \$1.2 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Julliard School. "Our best bet" that we had a real opportunity to make a mark on the future of the performing arts with their grant, and we every bit as important as our traditional areas of support. In health, medical research, culture and the social sciences." Home Foundation. President Bush. A House and Senate in recognizing the grant. Lincoln Center's share will be \$500,000 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 Julliard School, whose music and the performing arts are made, will get \$300,000. The Home Foundation, a leading supporter of the Lincoln Center Commission Corporate Fund, will make its total annual \$100,000 donation, too.

Example: multinomial model

- Data:
 - We observed N iid die rolls (K -sided): $\mathcal{D} = \{5, 1, K, \dots, 3\}$

- Representation:

Unit basis vectors: $x_n = \begin{bmatrix} x_{n,1} \\ x_{n,2} \\ \vdots \\ x_{n,K} \end{bmatrix}$, where $x_{n,k} \in \{0,1\}$, and $\sum_{k=1}^K x_{n,k} = 1$

- Model:

$x_{n,k} = 1$ w.p. θ_k , and $\sum_{k \in \{1, \dots, K\}} \theta_k = 1$

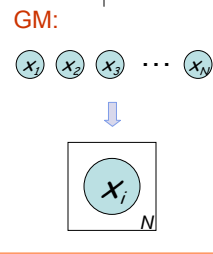
- How to write the likelihood of a single observation x_n ?

$$P(x_i) = P(\{x_{n,k} = 1, \text{ where } k \text{ index the die-side of the } n\text{th roll}\})$$

$$= \theta_k = \theta_1^{x_{n,1}} \times \theta_2^{x_{n,2}} \times \dots \times \theta_K^{x_{n,K}} = \prod_{k=1}^K \theta_k^{x_{n,k}}$$

- The likelihood of dataset $\mathcal{D} = \{x_1, \dots, x_N\}$:

$$P(x_1, x_2, \dots, x_N | \theta) = \prod_{n=1}^N P(x_n | \theta) = \prod_{n=1}^N \left(\prod_{k=1}^K \theta_k^{x_{n,k}} \right) = \prod_k \theta_k^{\sum_{n=1}^N x_{n,k}} = \prod_k \theta_k^{n_k}$$



MLE: constrained optimization with Lagrange multipliers



- Objective function:

$$\ell(\theta; \mathcal{D}) = \log \mathcal{P}(\mathcal{D} | \theta) = \log \prod_k \theta_k^{n_k} = \sum_k n_k \log \theta_k$$

- We need to maximize this subject to the constrain $\sum_{k=1}^K \theta_k = 1$

- Constrained cost function with a Lagrange multiplier

$$\bar{\ell} = \sum_k n_k \log \theta_k + \lambda \left(1 - \sum_{k=1}^K \theta_k \right)$$

- Take derivatives wrt θ_k

$$\frac{\partial \bar{\ell}}{\partial \theta_k} = \frac{n_k}{\theta_k} - \lambda = 0$$

$$n_k = \lambda \theta_k \Rightarrow \sum_k n_k = N = \lambda \sum_k \theta_k = \lambda$$

$$\Rightarrow \hat{\theta}_{k,MLE} = \frac{n_k}{N} \quad \text{or} \quad \hat{\theta}_{k,MLE} = \frac{1}{N} \sum_n x_{n,k}$$

Frequency as sample mean

- Sufficient statistics

- The counts, $\vec{n} = (n_1, \dots, n_K)$, $n_k = \sum_n x_{n,k}$, are **sufficient statistics** of data \mathcal{D}

© Eric Xing @ CMU, 2005-2009

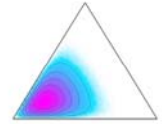
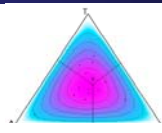
31

Bayesian estimation:



- Dirichlet distribution:

$$P(\theta) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k - 1} = C(\alpha) \prod_k \theta_k^{\alpha_k - 1}$$



- Posterior distribution of θ :

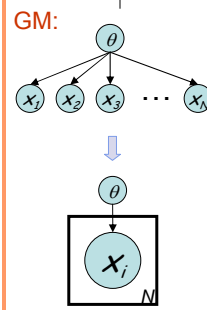
$$P(\theta | x_1, \dots, x_N) = \frac{p(x_1, \dots, x_N | \theta) p(\theta)}{p(x_1, \dots, x_N)} \propto \prod_k \theta_k^{n_k} \prod_k \theta_k^{\alpha_k - 1} = \prod_k \theta_k^{\alpha_k + n_k - 1}$$

- Notice the isomorphism of the posterior to the prior,
- such a prior is called a **conjugate prior**

- Posterior mean estimation:

$$\theta_k = \int \theta_k p(\theta | \mathcal{D}) d\theta = C \int \theta_k \prod_k \theta_k^{\alpha_k + n_k - 1} d\theta = \frac{n_k + \alpha_k}{N + |\alpha|}$$

Dirichlet parameters can be understood as pseudo-counts



© Eric Xing @ CMU, 2005-2009

32

More on Dirichlet Prior:



- Where is the normalize constant $\mathcal{C}(\alpha)$ come from?

$$\frac{1}{C(\alpha)} = \int \dots \int \theta_1^{\alpha_1-1} \dots \theta_K^{\alpha_K-1} d\theta_1 \dots d\theta_K = \frac{\prod_k \Gamma(\alpha_k)}{\Gamma(\sum_k \alpha_k)}$$

- Integration by parts
- $\Gamma(\alpha)$ is the gamma function: $\Gamma(\alpha) = \int_0^\infty t^{\alpha-1} e^{-t} dt$
- For integers, $\Gamma(n+1) = n!$
- Marginal likelihood:

$$p(\{x_1, \dots, x_N\} | \bar{\alpha}) = p(\bar{n} | \bar{\alpha}) = \int p(\bar{n} | \bar{\theta}) p(\bar{\theta} | \bar{\alpha}) d\bar{\theta} = \frac{C(\bar{\alpha})}{C(\bar{n} + \bar{\alpha})}$$

- Posterior in closed-form:

$$P(\bar{\theta} | \{x_1, \dots, x_N\}, \bar{\alpha}) = \frac{p(\bar{n} | \bar{\theta}) p(\bar{\theta} | \bar{\alpha})}{p(\bar{n} | \bar{\alpha})} = C(\bar{n} + \bar{\alpha}) \prod_k \theta_k^{\alpha_k + n_k - 1} = \text{Dir}(\bar{n} + \bar{\alpha})$$

- Posterior predictive rate:

$$p(x_{N+1} = i | \{x_1, \dots, x_N\}, \bar{\alpha}) = \int C(\bar{n} + \bar{\alpha}) \prod_k \theta_k^{\alpha_k + n_k - 1} \times \theta_i^{\alpha_i + n_i} d\bar{\theta} = \frac{C(\bar{n} + \bar{\alpha})}{C(\bar{n} + \bar{\alpha} + x_N)} = \frac{n_i + \alpha_i}{|\bar{n}| + |\bar{\alpha}|}$$

© Eric Xing @ CMU, 2005-2009

33

Sequential Bayesian updating



- Start with Dirichlet prior $\mathcal{P}(\bar{\theta} | \bar{\alpha}) = \text{Dir}(\bar{\theta} : \bar{\alpha})$
- Observe \mathcal{N}' samples with sufficient statistics \bar{n}' . Posterior becomes:

$$\mathcal{P}(\bar{\theta} | \bar{\alpha}, \bar{n}') = \text{Dir}(\bar{\theta} : \bar{\alpha} + \bar{n}')$$

- Observe another \mathcal{N}'' samples with sufficient statistics \bar{n}'' . Posterior becomes:

$$\mathcal{P}(\bar{\theta} | \bar{\alpha}, \bar{n}', \bar{n}'') = \text{Dir}(\bar{\theta} : \bar{\alpha} + \bar{n}' + \bar{n}'')$$

- So sequentially absorbing data in any order is equivalent to batch update.

© Eric Xing @ CMU, 2005-2009

34

Effect of Prior Strength



- Let $N = |\bar{n}| = \sum_k n_k$ be the number of observed samples
- Let $A = |\bar{\alpha}| = \sum_k \alpha_k$ be the number of "pseudo observations"
---- the strength of the prior
- Let $\bar{\alpha}' = |\bar{\alpha}| / A$ denote the prior means
- Then posterior mean is a convex combination of the prior mean and the MLE:

$$\begin{aligned}
 p(x_{N+1} = i | \{x_1, \dots, x_N\}, \bar{\alpha}) &= \frac{n_i + \alpha_i}{|\bar{n}| + |\bar{\alpha}|} = \frac{n_i + \alpha_i}{N + A} \\
 &= \frac{A}{N + A} \frac{\alpha_i}{A} + \frac{N}{N + A} \frac{n_i}{N} \\
 &= \lambda \alpha_i' + (1 - \lambda) \hat{\theta}_{k, MLE}
 \end{aligned}$$

where $\lambda = \frac{A}{N + A}$.

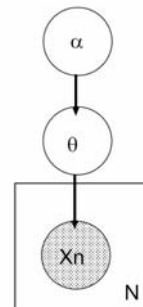
© Eric Xing @ CMU, 2005-2009

35

Hierarchical Bayesian Models



- θ are the parameters for the likelihood $p(x | \theta)$
- α are the parameters for the prior $p(\theta | \alpha)$.
- We can have hyper-hyper-parameters, etc.
- We stop when the choice of hyper-parameters makes no difference to the marginal likelihood; typically make hyper-parameters constants.
- Where do we get the prior?
 - Intelligent guesses
 - Empirical Bayes (Type-II maximum likelihood)
 - computing point estimates of α :

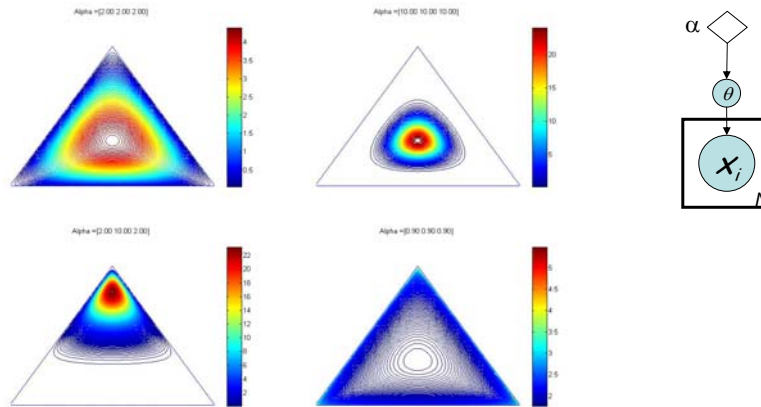


$$\hat{\alpha}_{MLE} = \arg \max_{\bar{\alpha}} p(\bar{n} | \bar{\alpha})$$

© Eric Xing @ CMU, 2005-2009

36

Limitation of Dirichlet Prior:



© Eric Xing @ CMU, 2005-2009

37

The Logistic Normal Prior

$$\theta \sim LN_K(\mu, \Sigma)$$

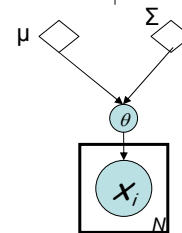
$$\gamma \sim N_{K-1}(\mu, \Sigma) \quad \gamma_K = \mathbf{0}$$

$$\theta_i = \exp\left\{ \gamma_i - \log\left(\mathbf{1} + \sum_{i=1}^{K-1} e^{\gamma_i} \right) \right\}$$

$$C(\gamma) = \log\left(\mathbf{1} + \sum_{i=1}^{K-1} e^{\gamma_i} \right)$$

- Log Partition Function
- Normalization Constant

Problem

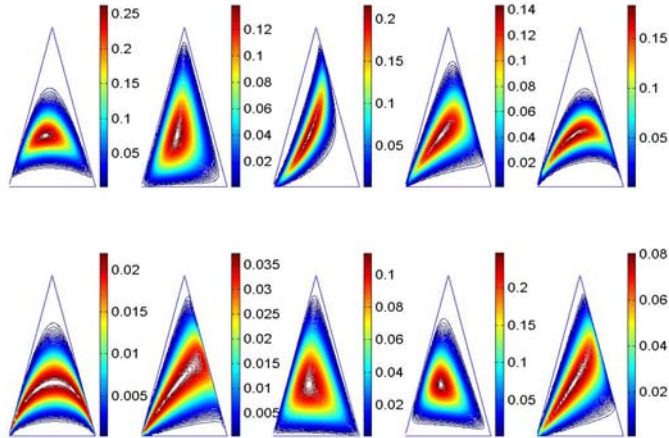


- Pro: co-variance structure
- Con: non-conjugate (we will discuss how to solve this later)

© Eric Xing @ CMU, 2005-2009

38

Logistic Normal Densities



Logistic Normal

© Eric Xing @ CMU, 2005-2009

39

Example: univariate-Gaussian



- Data:
 - We observed N iid real samples:
 $D = \{-0.1, 10, 1, -5.2, \dots, 3\}$

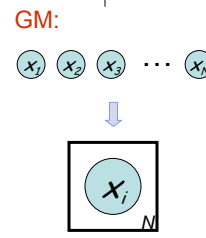
- Model: $P(x) = (2\pi\sigma^2)^{-1/2} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

- Log likelihood:

$$\ell(\theta; D) = \log P(D | \theta) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2} \sum_{n=1}^N \frac{(x_n - \mu)^2}{\sigma^2}$$

- MLE: take derivative and set to zero:

$$\begin{aligned} \frac{\partial \ell}{\partial \mu} &= (1/\sigma^2) \sum_n (x_n - \mu) & \Rightarrow & \mu_{MLE} = \frac{1}{N} \sum_n (x_n) \\ \frac{\partial \ell}{\partial \sigma^2} &= -\frac{N}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_n (x_n - \mu)^2 & & \sigma_{MLE}^2 = \frac{1}{N} \sum_n (x_n - \mu_{ML})^2 \end{aligned}$$



© Eric Xing @ CMU, 2005-2009

40

MLE for a multivariate-Gaussian



- It can be shown that the MLE for μ and Σ is

$$\mu_{MLE} = \frac{1}{N} \sum_n (x_n)$$

$$\Sigma_{MLE} = \frac{1}{N} \sum_n (x_n - \mu_{ML})(x_n - \mu_{ML})^T = \frac{1}{N} S$$

where the scatter matrix is

$$S = \sum_n (x_n - \mu_{ML})(x_n - \mu_{ML})^T = \left(\sum_n x_n x_n^T \right) - N \mu_{ML} \mu_{ML}^T$$

$$x_n = \begin{pmatrix} x_{n,1} \\ x_{n,2} \\ \vdots \\ x_{n,K} \end{pmatrix}$$

$$X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix}$$

- The sufficient statistics are $\sum_n x_n$ and $\sum_n x_n x_n^T$.
- Note that $X^T X = \sum_n x_n x_n^T$ may not be full rank (eg. if $N < D$), in which case Σ_{ML} is not invertible

© Eric Xing @ CMU, 2005-2009

41

Bayesian parameter estimation for a Gaussian



- There are various reasons to pursue a Bayesian approach
 - We would like to update our estimates sequentially over time.
 - We may have prior knowledge about the expected magnitude of the parameters.
 - The MLE for Σ may not be full rank if we don't have enough data.
- We will restrict our attention to conjugate priors.
- We will consider various cases, in order of increasing complexity:
 - Known σ , unknown μ
 - Known μ , unknown σ
 - Unknown μ and σ

© Eric Xing @ CMU, 2005-2009

42

Bayesian estimation: unknown μ , known σ

- Normal Prior:

$$P(\mu) = (2\pi\tau^2)^{-1/2} \exp\left\{-\frac{(\mu - \mu_0)^2}{2\tau^2}\right\}$$

- Joint probability:

$$P(\mathbf{x}, \mu) = (2\pi\sigma^2)^{-N/2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2\right\} \\ \times (2\pi\tau^2)^{-1/2} \exp\left\{-\frac{(\mu - \mu_0)^2}{2\tau^2}\right\}$$

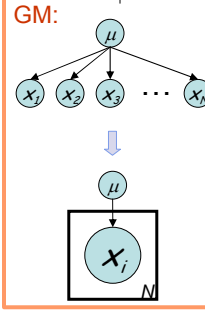
- Posterior:

$$P(\mu | \mathbf{x}) = (2\pi\tilde{\sigma}^2)^{-1/2} \exp\left\{-\frac{(\mu - \tilde{\mu})^2}{2\tilde{\sigma}^2}\right\}$$

where $\tilde{\mu} = \frac{N/\sigma^2}{N/\sigma^2 + 1/\tau^2} \bar{x} + \frac{1/\tau^2}{N/\sigma^2 + 1/\tau^2} \mu_0$, and $\tilde{\sigma}^2 = \left(\frac{N}{\sigma^2} + \frac{1}{\tau^2}\right)^{-1}$

© Eric Xing @ CMU, 2005-2009

43



Bayesian estimation: unknown μ , known σ

$$\mu_N = \frac{N/\sigma^2}{N/\sigma^2 + 1/\sigma_0^2} \bar{x} + \frac{1/\sigma_0^2}{N/\sigma^2 + 1/\sigma_0^2} \mu_0, \quad \tilde{\sigma}^2 = \left(\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}\right)^{-1}$$

- The posterior mean is a convex combination of the prior and the MLE, with weights proportional to the relative noise levels.
- The precision of the posterior $1/\tilde{\sigma}^2$ is the precision of the prior $1/\sigma_0^2$ plus one contribution of data precision $1/\sigma^2$ for each observed data point.

- Sequentially updating the mean

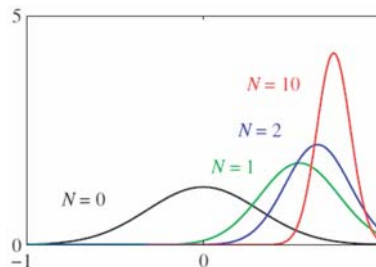
- $\mu^* = 0.8$ (unknown), $(\sigma^2)^* = 0.1$ (known)

- Effect of single data point

$$\mu_1 = \mu_0 + (x - \mu_0) \frac{\sigma_0^2}{\sigma^2 + \sigma_0^2} = x - (x - \mu_0) \frac{\sigma_0^2}{\sigma^2 + \sigma_0^2}$$

- Uninformative (vague/ flat) prior, $\sigma_0^2 \rightarrow \infty$

$$\mu_N \rightarrow \mu_0$$



© Eric Xing @ CMU, 2005-2009

44

Other scenarios



- Known μ , unknown $\lambda = 1/\sigma_2$
 - The conjugate prior for λ is a **Gamma** with shape a_0 and rate (inverse scale) b_0

$$p(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$$

- The conjugate prior for σ^2 is **Inverse-Gamma**

$$IG(\sigma^2|a, b) = \frac{1}{\Gamma(a)} b^a (\sigma^2)^{-(a+1)} \exp(-b/(\sigma^2))$$

- Unknown μ and unknown σ_2

- The conjugate prior is

Normal-Inverse-Gamma

$$\begin{aligned} P(\mu, \sigma^2) &= P(\mu|\sigma^2)P(\sigma^2) \\ &= \mathcal{N}(\mu|m, \sigma^2 V) IG(\sigma^2|a, b) \end{aligned}$$

- Semi conjugate prior

- Multivariate case:

- The conjugate prior is

Normal-Inverse-Wishart

$$\begin{aligned} P(\mu, \Sigma) &= P(\mu|\Sigma)P(\Sigma) \\ &= \mathcal{N}(\mu|\mu_0, \frac{1}{\kappa_0}\Sigma) \mathcal{IW}(\Sigma|\Lambda_0^{-1}, \nu_0) \end{aligned}$$

© Eric Xing @ CMU, 2005-2009

45

Summary



- Learning scenarios:
 - Data
 - Objective function
 - Frequentist and Bayesian
- Learning single-node GM – density estimation
 - Typical discrete distribution
 - Typical continuous distribution
 - Conjugate priors

© Eric Xing @ CMU, 2005-2009

46