

概率图模型理论及应用

Theory and Applications of Probabilistic Graphical Models
(Lesson 4 - mlEstimate)

欧智坚

清华大学电子工程系

Addr: 罗姆楼 6-104

Tel: 62796193

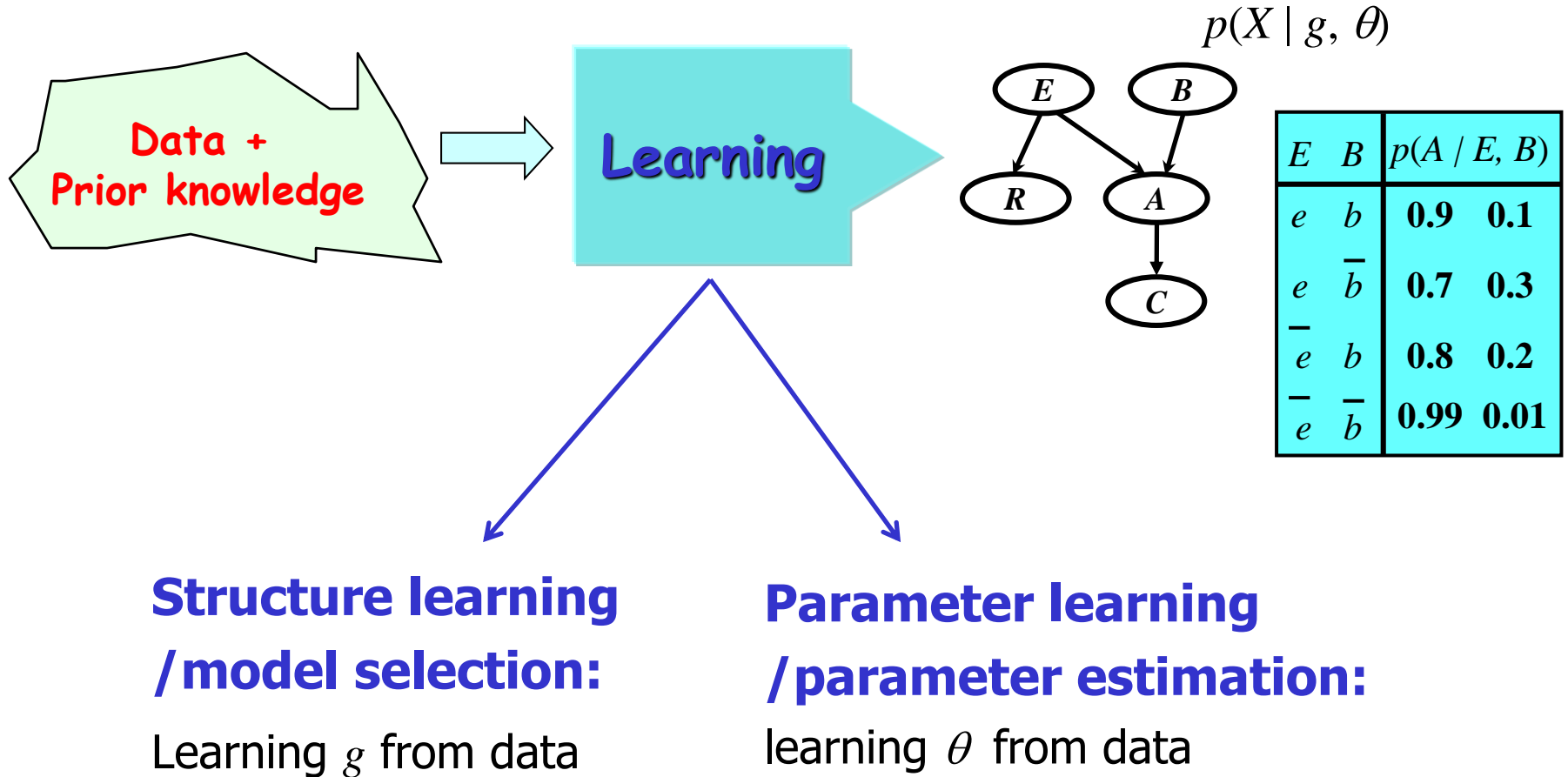
Email: ozj@tsinghua.edu.cn

- 希望课程讲述不要太粗略就好。由于课时有限，希望尽可能学习到概率图模型的精髓。
- 希望在课堂上重点讲解图模型的基本概念和背后思想，对于公式繁冗的推导过程可以只在课上进行一些提示，并推荐一些参考资料供自学，不要占用过多的课堂时间。
- 最好理论部分能在直接有step by step的推导指导，应用部分有hand by hand的工程或作业指导。
- 求老师讲的详细一点，因为发现上课真的听不懂呀。很多概念稍微解释一下吧。
- 感觉头两次课听不太明白，老师讲的有点像期末复习过一遍的感觉。希望老师能降低听课的先修知识门槛，讲的具体一点，多举一些例子。谢谢
- 数学推导可以简略，偏工程应用。
- 可以重点讲一下系统和应用的东西。
- 希望老师能够提前将期末 project 的选题发布，这样同学们能够较为充分地准备。
- 建议：给个5分钟让做笔记的同学自己说说总结或者问题，下面的同学可以提问之类的。

课程章节

- ❖ 第一章 引言 (**1**)
- ❖ 第二章 图模型的表示理论 (**2**)
 - **Semantics (DGM, UGM)**
 - **HMM, CRF**
- ❖ 第三章 图模型的推理理论 (**6**)
 - 精确推理: **variable-elimination, cluster-tree, triangulate**
 - 连续变量: **Kalman**
 - 采样近似: **sampling**
 - 变分近似: **variational**
- ❖ 第四章 图模型的学习理论 (**3**)
 - 参数学习: **maxlikelihoodEstimate, RFLearning, BayesEstimate**
 - 结构学习: **StructureLearning**
- ❖ 第五章 一个综合例子 (**1**)

Learning



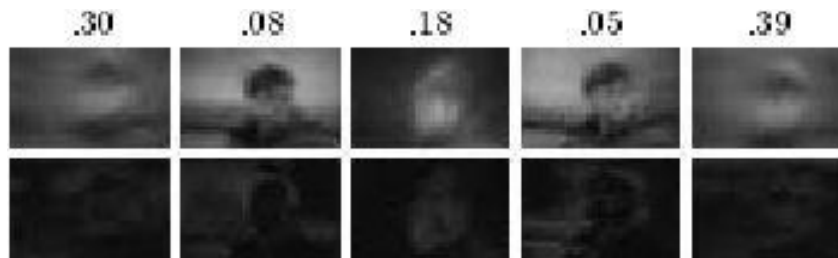
Data

- ❖ 总体分布 $p(x_{1:N} | g, \theta)$ 的一个个样本构成样本集/观测数据 $D = (x[1], \dots, x[M])$
 - 一个样本 $x_{1:N}[m] = (x_1[m], x_2[m], \dots, x_N[m])$
 - 独立同分布采样 (IID) : assume $x[1], \dots, x[M]$ are Independent and Identically Distributed $\sim p(x | g, \theta)$.
- ❖ 目标: 从 $D = (x[1], \dots, x[m], \dots, x[M])$ 中估计出 g, θ

总体分布: $p(x_1, x_2 | g, \theta)$, 其中 $w_k = p(x_1 = k)$ $p(x_2 | x_1 = k) = N(x | \mu_k, \Sigma_k)$

头姿类别 $x_1 \in 1:K$

观测图像 $x_2 \in R^{44*28}$



$x[1]$

	x_1	x_2
	斜脸	
	侧脸	
	斜脸	
	斜脸	
	正脸	

参数: $\theta = \{w_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$

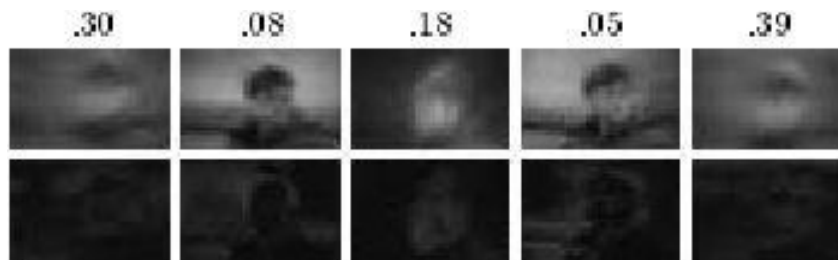
$x[5]$

Data — complete, incomplete

总体分布: $p(x_1, x_2 | g, \theta)$, 其中 $w_k = p(x_1 = k)$ $p(x_2 | x_1 = k) = N(x | \mu_k, \Sigma_k)$

头姿类别 $x_1 \in 1:K$

观测图像 $x_2 \in R^{44 \times 28}$



$x[1]$

x_1	x_2
斜脸	
侧脸	
斜脸	
斜脸	
正脸	

参数: $\theta = \{w_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$

$x[5]$

数据: 400 幅图像及其头姿类别, $D = (x[1], \dots, x[400])$

- 完备数据: 样本中各变量都赋值
- 不完备数据:
latent/hidden variables: 样本中某些变量的取值未知

The learning problem

参数学习

结构学习

	参数学习		结构学习
	Known structure	Bayesian	Unknown structure
Complete data	ML	Bayesian	
Incomplete data	ML	Bayesian	

DGMs, UGMs

Parameter learning

— ML (Known structure, complete data)

对单个分布的参数进行估计？

对一个贝叶斯网络的全体参数进行估计？

最大似然参数估计 (MLE)

给定一个概率分布的参数表达式 (parametric form)

记为 $p_{\theta}(x)$ 或 $p(x | \theta)$

从独立同分布样本集 $D = (x[1], \dots, x[M])$ 中估计出参数 θ ?



- $x \in \{1, 2, \dots, K\}$ is discrete r.v.
- $\theta_k = p(x=k), 1 \leq k \leq K$, is the parameters, $\theta = \{\theta_k | 1 \leq k \leq K\}$
- x is Gauss r.v. $X \sim N(\mu, \Sigma)$
- $\theta = (\mu, \Sigma)$ is the parameters

最大似然参数估计 (MLE)

给定一个概率分布的参数表达式 (parametric form)

记为 $p_{\theta}(x)$ 或 $p(x|\theta)$

从独立同分布样本集 $D = (x[1], \dots, x[M])$ 中估计出参数 θ ?

- 将 θ 视为一个未知常数
- θ 的一个估计/猜测 $\hat{\theta}$: 样本集的一个函数 $(x[1], \dots, x[M])$

样本集 $x[1:M]$ 下, θ 的似然函数 $p(x[1:M]|\theta) = \prod_{m=1}^M p(x[m]|\theta)$

概率分布函数
似然函数

http://en.wikipedia.org/wiki/Likelihood_function

给定参数 λ 下, 随机变量 Y 特定取值 y 的概率 (密度) 值 $p(y|\lambda)$ 视为

给定随机变量 Y 特定取值 y 下, 参数 λ 的似然值

最大似然估计: 使似然函数取最大 $\theta^{ML}(x[1:M]) = \arg \max_{\theta} p(x[1:M]|\theta)$

Multinomial distribution

- $x \in \{1, 2, \dots, K\}$ is discrete r.v.
- $\theta_k = p(x=k)$, $1 \leq k \leq K$, is the parameters, $\theta = \{\theta_k \mid 1 \leq k \leq K\}$
- 观测到独立同分布样本集 $D = (x[1], \dots, x[M])$
- 希望估计 θ ?



$$\text{似然函数 } p(x[1:M] | \theta) = \prod_{m=1}^M p(x[m] | \theta) = \prod_{k=1}^K \theta_k^{N_k}$$

N_k : 在样本集中 $x[m]=k$ 出现的次数

$$\text{最大似然估计 } \theta_k^{ML} = \frac{N_k}{\sum_{l=1}^K N_l}$$

(N_1, \dots, N_K) are sufficient statistics

Sufficient statistics

- ❖ 统计量：样本集 $D = (x[1], \dots, x[M])$ 的某函数
- ❖ Neyman Factorization theorem
一个统计量 $s(D)$, i.e., $s(x[1], \dots, x[M])$ 是充分统计量当且仅当 似然函数可以如下分解：

$$p(D | \theta) = g(\theta, s(D)) \cdot h(D)$$

参数与样本的关联 完全通过充分统计量来体现

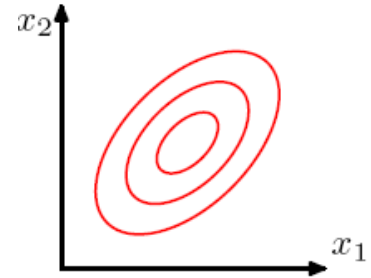
$$p(D | \theta) = \prod_{m=1}^M p(x[m] | \theta) = \prod_{k=1}^K \theta_k^{N_k}$$

N_k : 在样本集中 $x[m]=k$ 出现的次数



Gauss distribution

- x is Gauss r.v. $X \sim N(\mu, \Sigma)$
- $\theta = (\mu, \Sigma)$ is the parameters
- 观测到独立同分布样本集 $D = (x[1], \dots, x[M])$
- 希望估计 θ ?



$$p(x | \theta) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]$$

对数似然函数

$$\log p(x[1:M] | \mu, \Sigma) = \sum_{m=1}^M \log p(x[m] | \mu, \Sigma)$$

$$\max_{\mu, \Sigma} = -\frac{Md}{2} \log(2\pi) - \frac{M}{2} \left\{ \log |\Sigma| + \text{tr}(\Sigma^{-1} \bar{\Sigma}) + (\bar{\mu} - \mu) \Sigma^{-1} (\bar{\mu} - \mu)^T \right\}$$

$$\begin{cases} \frac{\partial L}{\partial \mu} = M \cdot \Sigma^{-1} (\bar{\mu} - \mu) = 0 \\ \frac{\partial L}{\partial \Sigma^{-1}} = -\frac{M}{2} \cdot \left\{ -\Sigma + \bar{\Sigma} + (\bar{\mu} - \mu)(\bar{\mu} - \mu)^T \right\} = 0 \end{cases} \quad \Rightarrow \quad \begin{cases} \bar{\mu} = \frac{1}{M} \sum_{m=1}^M x[m] \\ \bar{\Sigma} = \frac{1}{M} \sum_{m=1}^M (x[m] - \bar{\mu})(x[m] - \bar{\mu})^T \end{cases}$$

Learning parameters for BNs (complete data)

- 考虑贝叶斯网络 $X = \{X_1, X_2, \dots, X_N\}$
 假设：各个条件分布 $p(x_1 | pa_1), p(x_2 | pa_2), \dots, p(x_N | pa_N)$
 有各自表征参数 $\{\theta_1, \theta_2, \dots, \theta_N\} = \theta$

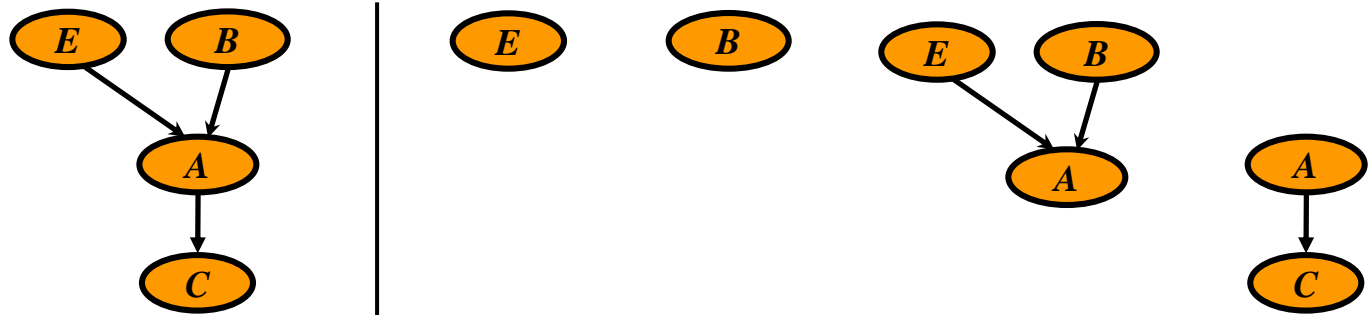
- IID样本集 $D = (x[1], \dots, x[M])$ 下似然函数

$$\max_{\theta} p(D | \theta) = \prod_{m=1}^M p(x[m] | \theta) = \prod_{m=1}^M \prod_{n=1}^N p(x_n[m] | pa_n[m], \theta_n) = \prod_{n=1}^N \max_{\theta_n} \prod_{m=1}^M p(x_n[m] | pa_n[m], \theta_n)$$

对每个条件分布 $p(x_n | pa_n)$ 分别估计其参数 θ_n

$$\hat{\theta}_n = \arg \max_{\theta_n} \prod_{m=1}^M p(x_n[m] | pa_n[m], \theta_n)$$

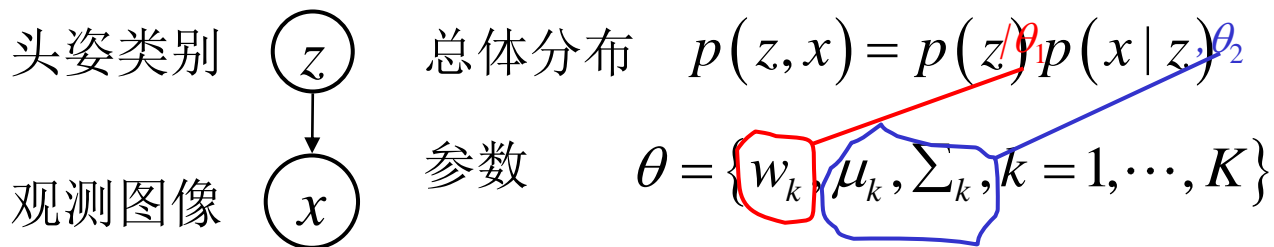
E, B, A, C
$\langle 1, 0, 0, 0 \rangle$
$\langle 1, 1, 1, 1 \rangle$
...
$\langle 1, 0, 1, 1 \rangle$



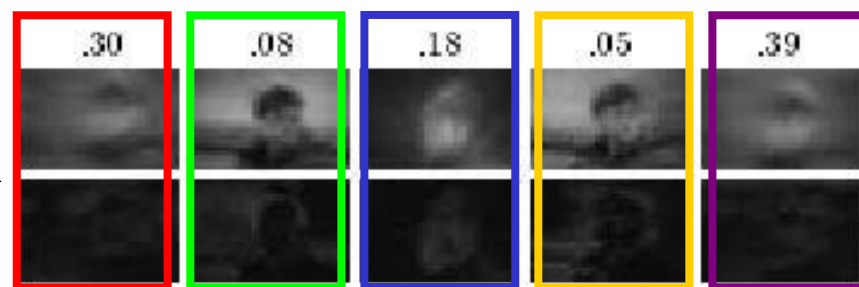
$$\max_{\theta} p(D | \theta) = \prod_{m=1}^M p(E[m] | \theta_1) p(B[m] | \theta_2) p(A[m] | E[m], B[m], \theta_3) p(C[m] | A[m], \theta_4)$$

$$\left. \left\{ \max_{\theta_1} \prod_{m=1}^M p(E[m] | \theta_1), \max_{\theta_2} \prod_{m=1}^M p(B[m] | \theta_2), \max_{\theta_3} \prod_{m=1}^M p(A[m] | E[m], B[m], \theta_3), \max_{\theta_4} \prod_{m=1}^M p(C[m] | A[m], \theta_4) \right\} \right|_{14}$$

高斯混合模型—完备数据



参数估计

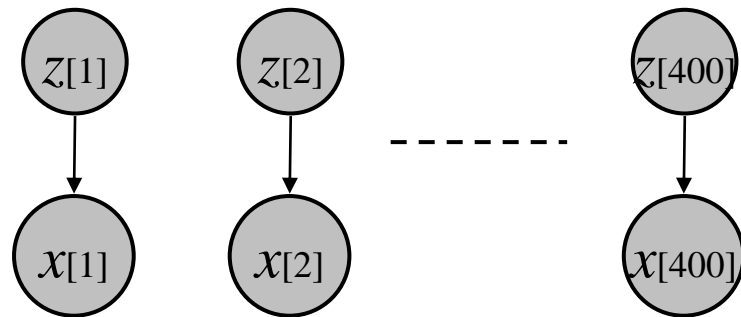


400 幅图像及其头姿类别标注
完备数据 $D = \left(\left(\begin{matrix} x[1] \\ z[1] \end{matrix} \right), \dots, \left(\begin{matrix} x[400] \\ z[400] \end{matrix} \right) \right)$



高斯混合模型—完备数据

$$\theta = \{w_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$$



$$\log p(x[1:M], z[1:M] | \theta)$$

$$= \sum \log p(x[m], z[m] | \theta)$$

$$= \sum_m \left[\log p(z[m] | \theta) + \log p(x[m] | z[m], \theta) \right]$$

$$= \sum_m \log p(z[m] | w_{1:K}) + \sum_m \log N(x[m] | \mu_{z[m]}, \Sigma_{z[m]})$$

$$\sum_{k=1}^K \sum_{m: z[m]=k} \log w_k$$

(A blue arrow points from the first term of the previous equation to this one.)

$$\sum_{k=1}^K \sum_{m: z[m]=k} \log N(x[m] | \mu_k, \Sigma_k)$$

(A green arrow points from the second term of the previous equation to this one.)

离散变量 z 的400个样本下的对数似然值,

对第 k 类, 高斯变量 x 的

N_k 个样本下的对数似然值,

$$\mu_k^{ML}, \Sigma_k^{ML}$$

$$w_k^{ML} = \frac{N_k}{M}$$

$$= \frac{1}{M} \sum_{m=1}^M 1(z[m] = k)$$

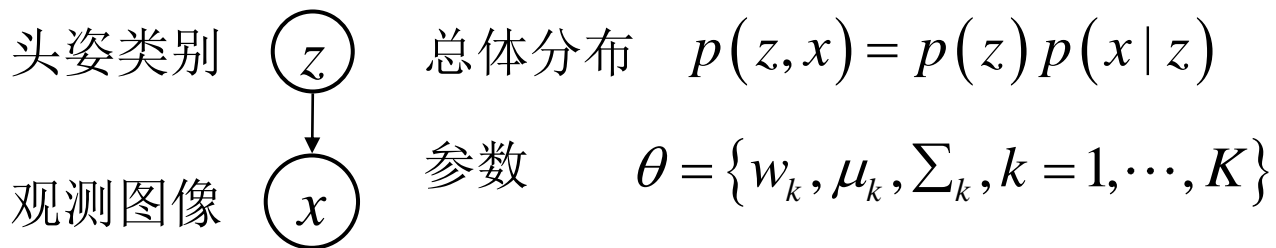
$$\begin{cases} \bar{\mu}_k^{ML} = \frac{1}{N_k} \sum_{m: z[m]=k} x[m] \\ \bar{\Sigma}_k^{ML} = \frac{1}{N_k} \sum_{m: z[m]=k} (x[m] - \bar{\mu}_k^{ML})(x[m] - \bar{\mu}_k^{ML})^T \end{cases}$$

Parameter learning

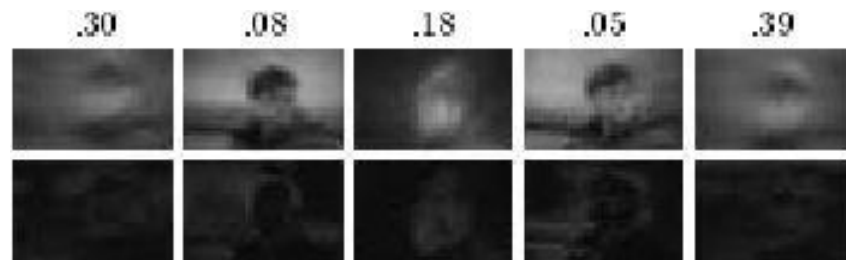
— ML (Known structure, incomplete data)

Expectation-Maximization 算法

高斯混合模型—不完备数据



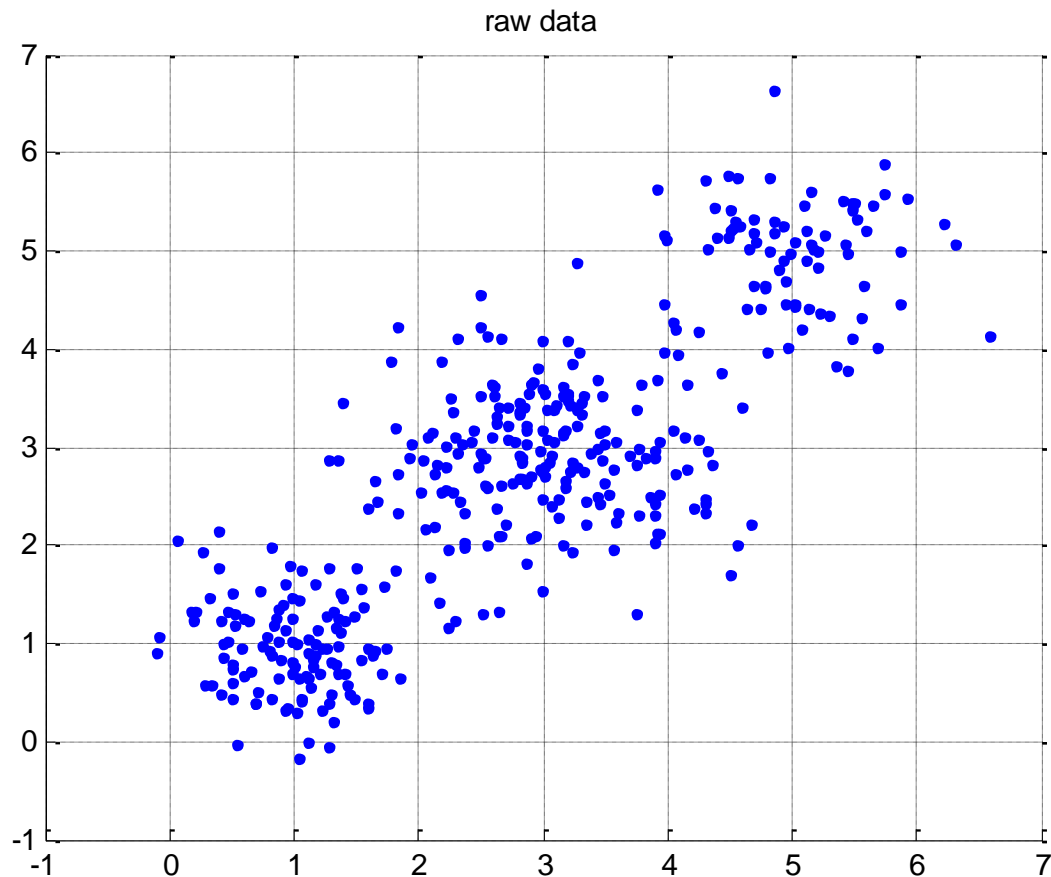
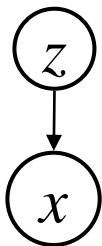
参数估计



数据：400 幅图像 $D = (x[1], \dots, x[400])$

不完备数据

Homework3_em



```
function [weight, meanvec, stdvec] = EmEstimate(x, iternum)
% x is the input observation, D-dim vectors * N
% iternum is the given number for EM iterations
```

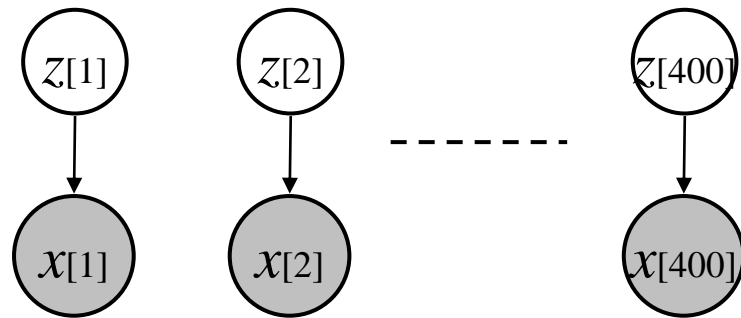
高斯混合模型—不完备数据

$$\theta = \{w_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$$

$$\log p(x[1:M] | \theta)$$

$$= \sum_m \log p(x[m] | \theta)$$

$$\stackrel{\text{m}}{\text{a}} \sum_m \log \left(\sum_k w_k N(x[m] | \mu_k, \Sigma_k) \right)$$



$$\frac{\partial}{\partial w_k} \log p(x[1:M] | \theta) = \sum_m \frac{\partial}{\partial w_k} \log \left(\sum_k w_k N(x[m] | \mu_k, \Sigma_k) \right)$$

$$= \sum_m \frac{1}{\sum_k w_k N(x[m] | \mu_k, \Sigma_k)} \frac{\partial}{\partial w_k} \left(\sum_k w_k N(x[m] | \mu_k, \Sigma_k) \right)$$

$$\frac{\partial}{\partial \mu_k} \log p(x[1:M] | \theta) =$$

$$\frac{\partial}{\partial \Sigma_k} \log p(x[1:M] | \theta) =$$

需要求解

$$\theta = \{w_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$$

联立非线性方程!

EM一般讨论

❖ 记 x 为全体观测值，记 z 为全体隐变量

■ 联合分布： $p(x, z | \theta)$

$$\theta^{ML} = \arg \max_{\theta} \log p(x | \theta)$$

$$= \log \sum_z p(x, z | \theta)$$

❖ $\log p(\overset{z}{x} | \theta)$ is called the incomplete log-likelihood.

■ 联合分布 $p(x, z | \theta)$ 的分解表示得不到利用

❖ $\log p(x, z | \theta)$ is called the complete log-likelihood.

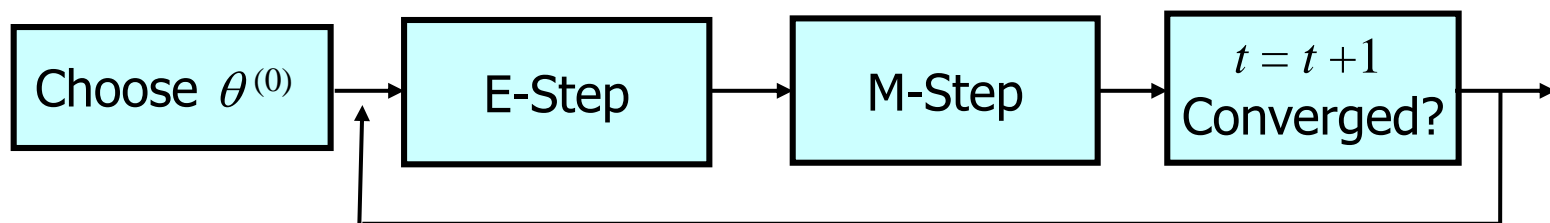
■ 利用分解

EM算法描述

- ❖ z 是隐变量，完备对数似然函数 $\log p(x, z | \theta)$ 是 z 的一个函数
- ❖ 使用 完备对数似然函数的期望？
- ❖ 完备对数似然函数 $\log p(x, z | \theta)$ 在条件分布 $p(z | \theta^{old}, x)$ 下的期望

$$Q(\theta | \theta^{old}) = E[\log p(x, z | \theta) | \theta^{old}, x] = \sum_z p(z | \theta^{old}, x) \log p(x, z | \theta)$$

- ❖ 求解 $\theta^* = \arg \max_{\theta} Q(\theta | \theta^{old})$
成立 $\log p(x | \theta^{old}) \leq \log p(x | \theta^*)$
- ❖ EM算法是一个迭代过程



Jensen Inequality : for convex \cap function f

$$E[f(U)] \leq f(E[U])$$

EM算法证明

$$\log p(x, z | \theta) = \log p(x | \theta) + \log p(z | \theta, x), \quad \forall \theta \text{ applying } E[\dots | \theta^{(old)}, x]$$

$$E[\log p(x, z | \theta) | \theta^{(old)}, x] = \log p(x | \theta) + E[\log p(z | \theta, x) | \theta^{(old)}, x], \quad \forall \theta$$

$$E[\log p(x, z | \theta^{(old)}) | \theta^{(old)}, x] = \log p(x | \theta^{(old)}) + E[\log p(z | \theta^{(old)}, x) | \theta^{(old)}, x]$$

$$\begin{pmatrix} E[\log p(x, z | \theta) | \theta^{(old)}, x] \\ -E[\log p(x, z | \theta^{(old)}) | \theta^{(old)}, x] \end{pmatrix} = \begin{pmatrix} \log p(x | \theta) \\ -\log p(x | \theta^{(old)}) \end{pmatrix} + E\left[\log \frac{p(z | \theta, x)}{p(z | \theta^{(old)}, x)} \mid \theta^{(old)}, x\right]$$

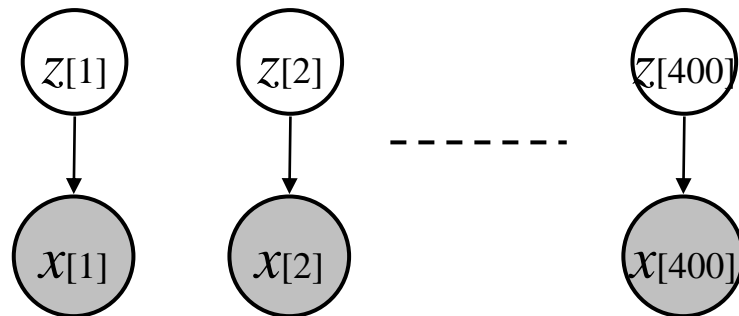
$$\begin{aligned} Q(\theta | \theta^{(old)}) &= E[\log p(x, z | \theta) | \theta^{(old)}, x] \leq \log E\left[\frac{p(z | \theta, x)}{p(z | \theta^{(old)}, x)} \mid \theta^{(old)}, x\right] \\ &= \log \sum_z \frac{p(z | \theta, x)}{p(z | \theta^{(old)}, x)} p(z | \theta^{(old)}, x) \end{aligned}$$

EM Example: Learning with GMM

$$Q(\theta | \theta^{(old)}) = E[\log p(x, z | \theta) | \theta^{(old)}, x]$$

❖ 给定不完备数据 $(x[1], \dots, x[M])$

$$E[\log p(x[1:M], z[1:M] | \theta) | \theta^{(old)}, x[1:M]]$$



$$= \sum_m E[\log p(x[m], z[m] | \theta) | \theta^{(old)}, x[1:M]]$$

$$= \sum_m \sum_{z[m]} p(z[m] | \theta^{(old)}, x[m]) \log p(x[m], z[m] | \theta)$$

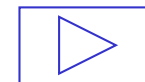
$$= \sum_m \sum_{z[m]} p(z[m] | \theta^{(old)}, x[m]) \{ \log p(x[m] | \theta, z[m]) + \log p(z[m] | \theta) \}$$

$$= \sum_m \sum_k p(z[m] = k | \theta^{(old)}, x[m]) \{ \log p(x[m] | \theta, z[m] = k) + \log p(z[m] = k | \theta) \}$$

$$\max_{\{w_k, \mu_k, \Sigma_k, k=1:K\}} \left\{ \sum_k \sum_m \gamma_m(k) \log N(x[m] | \mu_k, \Sigma_k) + \sum_k \sum_m \gamma_m(k) \log w_k \right\}$$

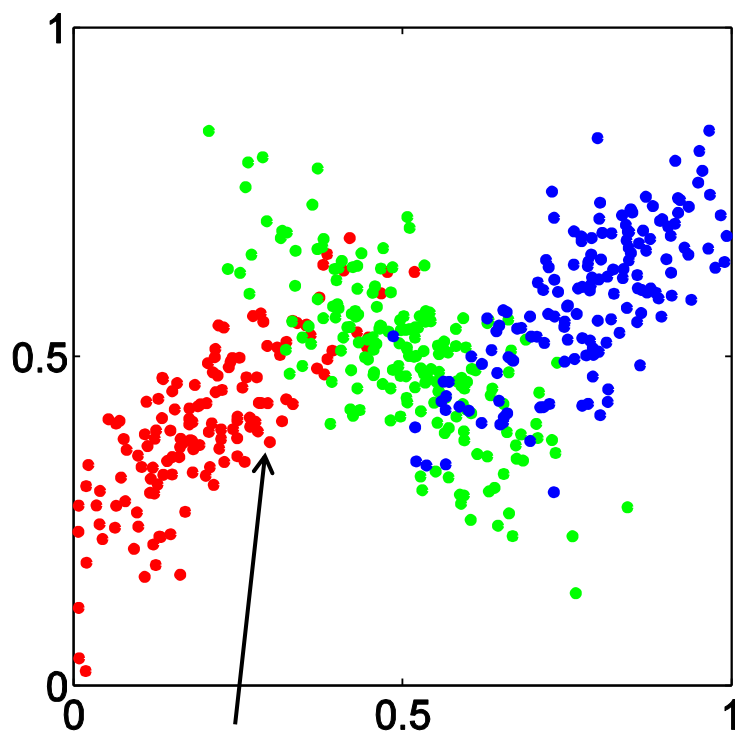
subject to: $\sum_k w_k = 1$

Posterior Probabilities

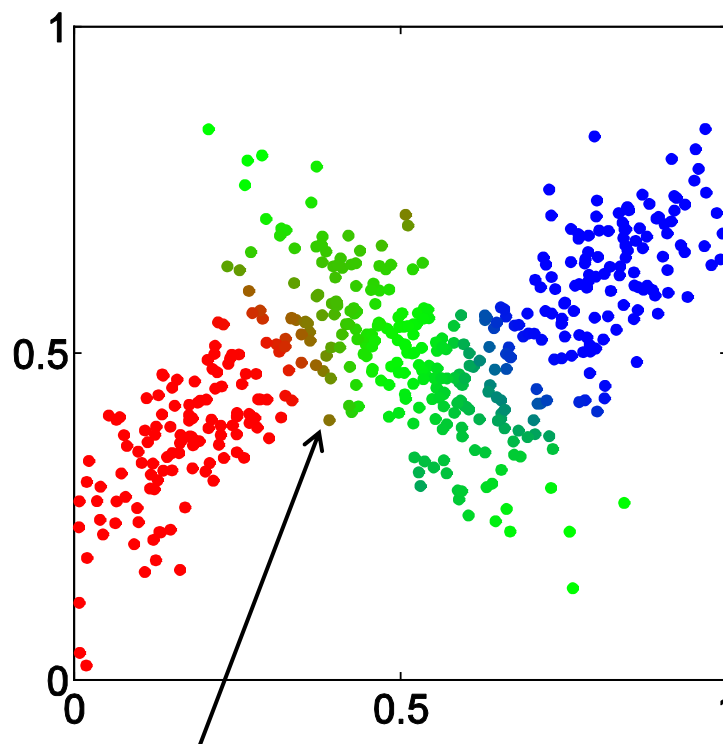


$\frac{m}{a}$ 不完备数据下目标函数 = $\sum_{k=1}^K \sum_{m=1}^M \gamma_m(k) \log N(x[m] | \mu_k, \Sigma_k) + \frac{m}{X} \sum_{k=1}^K \sum_{m=1}^M \gamma_m(k) \log w_k$

$\frac{m}{a}$ 完备数据下目标函数 = $\sum_{k=1}^K \sum_{m=1}^M 1(z[m] = k) \log N(x[m] | \mu_k, \Sigma_k) + \frac{m}{X} \sum_{k=1}^K \sum_{m=1}^M 1(z[m] = k) \log w_k$

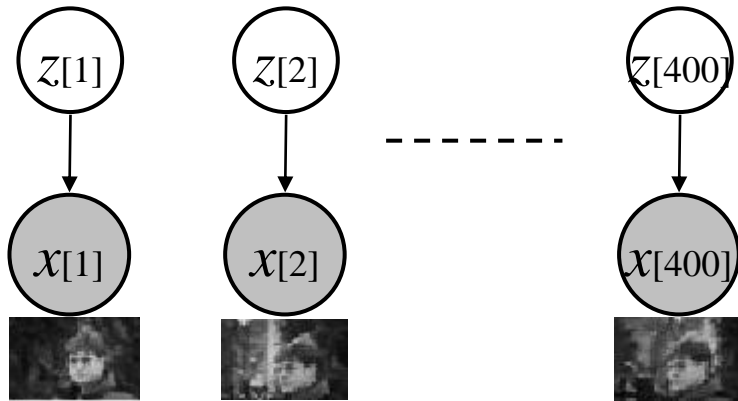
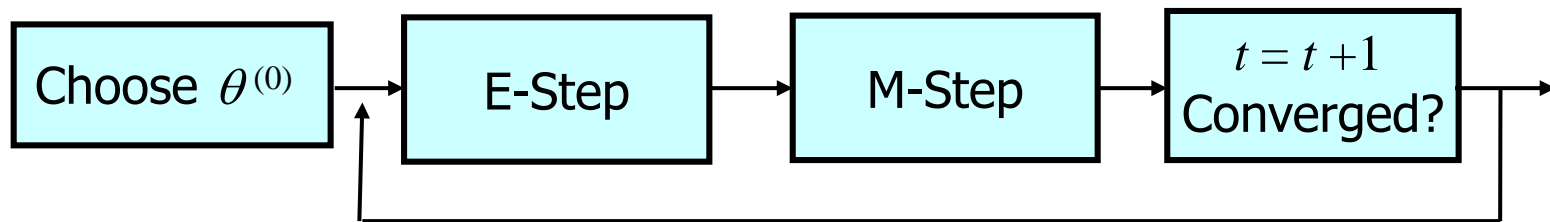


Hard assignment



Soft assignment

EM Example: Learning with GMM



$$\gamma_m(k) = p(z[m] = k | \theta^{(old)}, x[m])$$

$$= \frac{w_k^{(old)} N(x[m] | \mu_k^{(old)}, \Sigma_k^{(old)})}{\sum_k w_k^{(old)} N(x[m] | \mu_k^{(old)}, \Sigma_k^{(old)})}$$

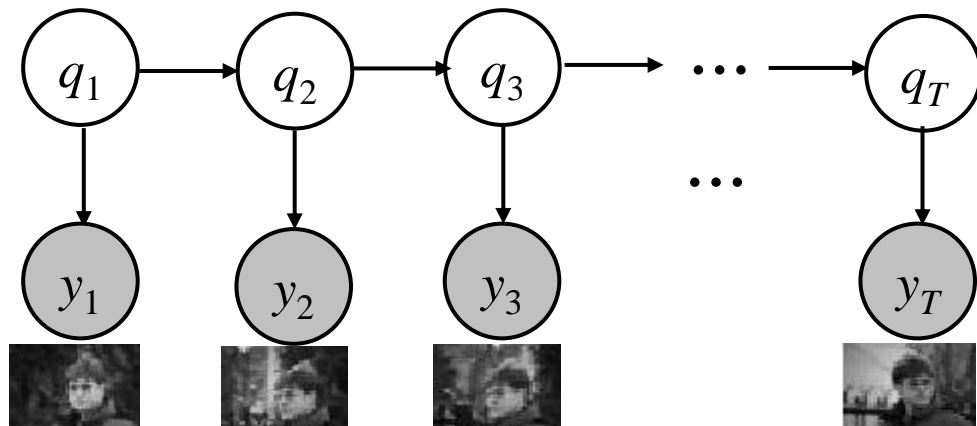
$$\mu_k^* = \frac{\sum_m \gamma_m(k) x[m]}{\sum_m \gamma_m(k)} \quad \bar{\mu}_k^{ML} = \frac{\sum_{m=1}^M 1(z[m] = k) x[m]}{\sum_{m=1}^M 1(z[m] = k)}$$

$$\Sigma_k^* = \frac{\sum_m \gamma_m(k) (x[m] - \mu_k^*)(x[m] - \mu_k^*)^T}{\sum_m \gamma_m(k)}$$

$$= \frac{\sum_m \gamma_m(k) x[m] x[m]^T}{\sum_m \gamma_m(k)} - \mu_k^* (\mu_k^*)^T$$

$$w_k^* = \frac{\sum_m \gamma_m(k)}{M}$$

EM Example: Learning with HMM



$$\lambda = (\pi, A, B) \quad \max_{\lambda} \log p(y_{1:T} | \lambda) ?$$

$$\lambda^* = \arg \max_{\lambda} E \left[\log p(y_{1:T}, q_{1:T} | \lambda) \mid \lambda^{(t)}, y_{1:T} \right]$$

$$= \arg \max_{\lambda} \left\{ \begin{array}{l} E \left[\log p(q_1 | \lambda) \mid \lambda^{(t)}, y_{1:T} \right] \\ + \sum_{t=2}^T E \left[\log p(q_t | \lambda, q_{t-1}) \mid \lambda^{(t)}, y_{1:T} \right] \\ + \sum_{t=1}^T E \left[\log p(y_t | \lambda, q_t) \mid \lambda^{(t)}, y_{1:T} \right] \end{array} \right\} \quad \leftarrow \begin{array}{l} \log p(y_{1:T}, q_{1:T} | \lambda) \\ = \log p(q_1 | \lambda) \\ + \sum_{t=2}^T \log p(q_t | \lambda, q_{t-1}) \\ + \sum_{t=1}^T \log p(y_t | \lambda, q_t) \end{array}$$

From EM to SA

❖ 记 x 为全体观测值, 记 z 为全体隐变量

■ 联合分布: $p(x, z | \theta)$

$$\theta^{ML} = \arg \max_{\theta} \log p(x | \theta)$$

$$Q(\theta | \theta^{(old)}) = E[\log p(x, z | \theta) | \theta^{(old)}, x] = \sum_z p(z | \theta^{(old)}, x) \log p(x, z | \theta)$$

$$\text{Fisher Equality: } \frac{\partial \log p(x | \theta)}{\partial \theta} = E_{p(z|x, \theta)} \left[\frac{\partial \log p(x, z | \theta)}{\partial \theta} \right]$$

$$\therefore E_{p(z|x, \theta)} \left[\frac{\partial \log p(z | x, \theta)}{\partial \theta} \right] = 0$$

Problem: The objective is to find a solution θ to $E_{Y \sim f(\cdot; \theta)}[H(Y; \theta)] = \alpha$, where $\theta \in R^d$, noisy observation $H(Y; \theta) \in R^d$.

- Gu & Kong, A stochastic approximation algorithm with Markov chain Monte-Carlo method for incomplete data estimation problems, PNAS 1998.
- Delyon, Lavielle, and Moulines. "Convergence of a stochastic approximation version of the EM algorithm." Annals of statistics, 1999.

Parameter learning

— ML (Known structure, complete data)

无向图

Parameter learning for UGMs (complete data)

- 贝叶斯网络:

全体参数的对数似然函数 = 各个节点处参数的对数似然函数之和

- 无向图模型:

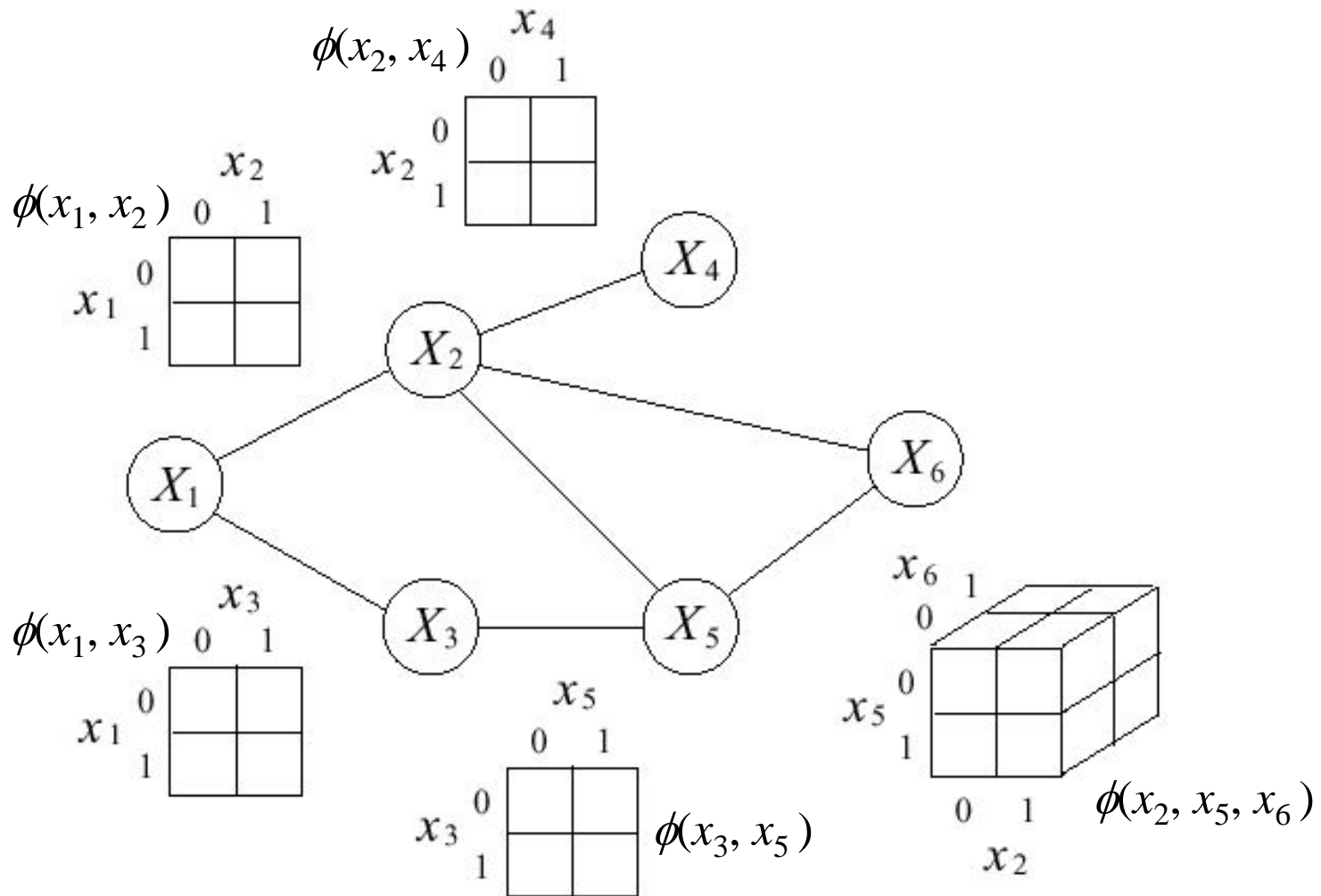
归一化常数 Z 是所有参数的函数

$$\log p(x) = \sum_{c \in \mathcal{C}} \log \phi_c(x_c) - \log Z \quad Z = \sum_x \prod_{c \in \mathcal{C}} \phi_c(x_c)$$

- 即使是完备数据情形下，无向图模型参数估计需要推理计算（求边缘分布）
- 写出目标函数，直接运用梯度下降法等最优化方法

$$p(x_1, x_2, x_3, x_4, x_5, x_6) = \frac{1}{Z} \phi(x_1, x_2) \phi(x_1, x_3) \phi(x_2, x_4) \phi(x_3, x_5) \phi(x_2, x_5, x_6)$$

Example



无向图的对数似然函数：离散情形

❖ 给定IID样本集 $D = (x[1], \dots, x[M])$

对数似然函数 $\log p(D | \theta) = \sum_m \log p(x[m] | \theta)$

$$= \sum_m \sum_x 1(x[m] = x) \log p(x | \theta)$$

样本集中特定 x 出现次数

$$\text{count}(x) = \sum_m 1(x[m] = x)$$

$$= \sum_x \sum_m 1(x[m] = x) \log p(x | \theta)$$

$$= \sum_x \text{Count}(x) \log \left(\frac{1}{Z} \prod_{c \in \mathcal{C}} \phi_c(x_c) \right)$$

样本集中特定 x_c 出现次数

$$\text{count}(x_c) = \sum_{x \setminus x_c} \text{count}(x)$$

$$= \sum_x \sum_c \text{Count}(x) \log \phi_c(x_c) - M \log Z$$

簇变量取特定值的次数(**clique count**)
是(离散)无向图模型的充分统计量

$$= \sum_c \sum_{x_c} \text{Count}(x_c) \log \phi_c(x_c) - M \log Z$$

对数似然函数求导

对数似然函数 $\log p(D|\theta) = \sum_c \sum_{x_c} \text{Count}(x_c) \log \phi_c(x_c) - M \log Z$

第一项求导 $\frac{\partial l_1}{\partial \phi_c(x_c)} = \frac{\text{Count}(x_c)}{\phi_c(x_c)}$

第二项求导 $\frac{\partial \log Z}{\partial \phi_c(x_c)} = \frac{1}{Z} \frac{\partial Z}{\partial \phi_c(x_c)} = \frac{1}{Z} \frac{\partial \sum_y \prod_d \phi_d(y_d)}{\partial \phi_c(x_c)}$ $Z = \sum_x \prod_{c \in \mathcal{C}} \phi_c(x_c)$

$$= \frac{1}{Z} \sum_y 1(y_d, x_c) \frac{\partial \prod_d \phi_d(y_d)}{\partial \phi_c(x_c)} \quad \text{只留下 } y_d = x_c \text{ 的项}$$

$$= \frac{1}{Z} \sum_y 1(y_d, x_c) \frac{1}{\phi_c(x_c)} \prod_d \phi_d(y_d)$$

$$= \frac{1}{\phi_c(x_c)} \sum_y 1(y_d, x_c) \frac{1}{Z} \prod_d \phi_d(y_d) = \frac{p(x_c)}{\phi_c(x_c)}$$

必要条件：簇上边缘分布...

对数似然函数求导 $\frac{\partial \log p(D|\theta)}{\partial \phi_c(x_c)} = \frac{\text{Count}(x_c)}{\phi_c(x_c)} - M \frac{p(x_c)}{\phi_c(x_c)}$

在最大似然参数估计 θ^{ML} 下 $p^{ML}(x_c) = \frac{\text{Count}(x_c)}{M} = q_c(x_c)$

取 θ^{ML} 时簇上的边缘分布 = 簇上的经验分布

Model marginal = empirical marginal

只是指出：取得 θ^{ML} 时，簇上的边缘分布所需满足的必要条件
具体怎么求出 θ^{ML} ？

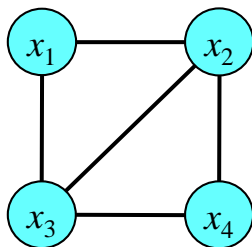
课程章节

- ❖ 第一章 引言 (**1**)
- ❖ 第二章 图模型的表示理论 (**2**)
 - **Semantics (DGM, UGM)**
 - **HMM, CRF**
- ❖ 第三章 图模型的推理理论 (**6**)
 - 精确推理: **variable-elimination, cluster-tree, triangulate**
 - 连续变量: **Kalman**
 - 采样近似: **sampling**
 - 变分近似: **variational**
- ❖ 第四章 图模型的学习理论 (**3**)
 - 参数学习: **maxlikelihoodEstimate, RFLearning, BayesEstimate**
 - 结构学习: **StructureLearning**
- ❖ 第五章 一个综合例子 (**1**)

MLE for undirected graph

Tabular?	Decomposable?	Max-clique?	Method
√	√	√	Direct
√	—	—	IPF Iterative Proportional Fitting
×	—	—	GIS Generalized Iterative Scaling
×	—	—	梯度下降法

局部势函数 是否 全部定义在最大簇上？



$$p(x_1, x_2, x_3, x_4) = \frac{1}{Z} \phi(x_1, x_2, x_3) \phi(x_2, x_3, x_4)$$

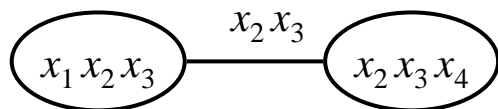
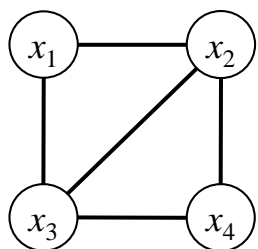
$$p(x_1, x_2, x_3, x_4) = \frac{1}{Z} \phi(x_1, x_2) \phi(x_1, x_3) \phi(x_2, x_3, x_4)$$

MLE for decomposable undirected graph

❖ 定理：一个可分解无向图 \mathcal{G} 上的联合分布可以表示成：

$$p(x) = \frac{\prod_{c \in \mathcal{C}} p_c(x_c)}{\prod_{s \in \mathcal{S}} p_s(x_s)}$$

■ \mathcal{C} 是最大簇集合， \mathcal{S} 是图 \mathcal{G} 的连接树上的隔离子集合



$$\begin{aligned} p(x_1, x_2, x_3, x_4) &= \frac{1}{Z} \phi(x_1, x_2, x_3) \phi(x_2, x_3, x_4) \\ &= \frac{p(x_1, x_2, x_3) p(x_2, x_3, x_4)}{p(x_2, x_3)} \end{aligned}$$

在最大似然参数估计 θ^{ML} 下 $p^{ML}(x) = \frac{\prod_{c \in \mathcal{C}} q_c(x_c)}{\prod_{s \in \mathcal{S}} q_s(x_s)} = \frac{1}{Z} \prod_{c \in \mathcal{C}} \phi_c^{ML}(x_c)$ 取 $\phi_c^{ML}(x_c) = \frac{q_c(x_c)}{\dots}$

$$p^{ML}(x_1, x_2, x_3, x_4) = \frac{q(x_1, x_2, x_3) q(x_2, x_3, x_4)}{q(x_2, x_3)} = \frac{1}{Z} \phi^{ML}(x_1, x_2, x_3) \phi^{ML}(x_2, x_3, x_4)$$

$$\phi^{ML}(x_1, x_2, x_3) = \quad \phi^{ML}(x_2, x_3, x_4) =$$

IPF (Iterative Proportional Fitting)

对数似然函数求导
$$\frac{\partial \log p(D|\theta)}{\partial \phi_c(x_c)} = M \frac{q_c(x_c)}{\phi_c(x_c)} - M \frac{p(x_c)}{\phi_c(x_c)}$$

最大似然参数估计 θ^{ML} 应满足
$$\frac{q(x_c)}{\phi_c(x_c)} = \frac{p(x_c)}{\phi_c(x_c)}$$

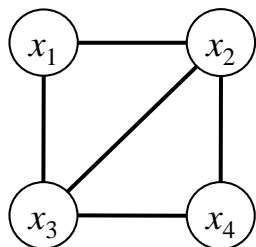
$p(x_c)$ 是 $\{\phi_c(x_c)\}$ 的函数

直接求解 $\{\phi_c(x_c)\}$ 很困难, $\{\phi_c(x_c)\}$ 是非线性方程的解

IPF: 将 $\phi_c(x_c)$ 视为下面不动点方程的不动点, 迭代求解:

$$\phi_c(x_c) = \frac{q(x_c)}{p(x_c)} \phi_c(x_c) \quad \phi_c^{(t+1)}(x_c) = \frac{q(x_c)}{p^{(t)}(x_c)} \phi_c^{(t)}(x_c)$$

推理计算

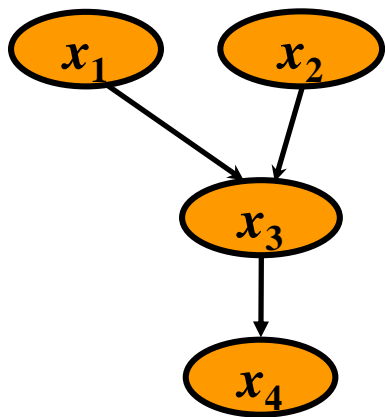


$$p(x_1, x_2, x_3, x_4) = \frac{1}{Z} \phi_{1,2}(x_1, x_2) \phi_{1,3}(x_1, x_3) \phi_{2,3,4}(x_2, x_3, x_4)$$

$$\phi_{1,2}^{(t+1)}(x_1, x_2) = \frac{q(x_1, x_2)}{\frac{1}{Z^{(t)}} \sum_{x_3, x_4} \phi_{1,2}^{(t)}(x_1, x_2) \phi_{1,2}^{(t)}(x_1, x_3) \phi_{2,3,4}^{(t)}(x_2, x_3, x_4)} \phi_{1,2}^{(t)}(x_1, x_2)$$

Example: Multinomial Bayes net

- 假设变量 X_n 有 K_n 个不同可能取值
- 结点 x_n 的条件分布 $p(x_n | pa_n)$ 含有一系列多元分布。对父结点集 pa_n 的每个可能取值组合 i ，有一个多元分布 $p(x_n | pa_n = i)$
- $\theta = \{\theta_n | n = 1, \dots, N\}$ $\theta_n = \{\theta_{n,i} | i = 1, \dots\}$ $\theta_{n,i} = \{\theta_{n,i,k} | k = 1, \dots, K_n\}$



$$p(x_3 | pa_3 = (0,0))$$

$$p(x_3 | pa_3 = (0,1))$$

$$p(x_3 | pa_3 = (1,0))$$

$$p(x_3 | pa_3 = (1,1))$$

$\theta_{n,i,k}$ \square $p(x_n = k | pa_n = i)$

结点 \swarrow \searrow

x_n 的第 k 个可能取值
 $1 \leq k \leq K_n$

pa_n 的第 i 个可能取值
 $1 \leq i \leq \prod_{x_l \in pa_n} K_l$

Example: MLE for multinomial Bayes net

- 似然函数 $p(x_n[1:M] | \theta_n) \square \prod_{m=1}^M p(x_n[m] | pa_n[m], \theta_n)$
$$= \prod_i \prod_k (\theta_{n,i,k})^{N_{n,i,k}}$$

- 充分统计量 $N_{n,i,k} \square \sum_{m=1}^M 1(pa_n[m] = i, x_n[m] = k)$

- 最大似然估计:
$$\hat{\theta}_{n,i,k} = \frac{N_{n,i,k}}{\sum_{l=1}^{K_n} N_{n,i,l}}$$