# Bayesian models for missing data



### Missing data

### EM algorithm





### **Missing data**











- **Definition:** No data value is stored for the variable in an observation.
- Reasons for missing data
  - ♦ inaccessible
  - omitted
  - unavailable attributes
  - high cost



#### Notations

$$X = (X_{obs} + X_{mis})$$

$$X = \begin{bmatrix} X_{11} & \cdots & X_{1p} \\ \vdots & \ddots & \vdots \\ X_{n1} & \cdots & X_{np} \end{bmatrix} = \begin{bmatrix} X_1, X_2, \dots X_p \end{bmatrix}$$

$$M_{ij} = \begin{cases} 1 \text{ as } X_{ij} \text{ is observed} \\ 0 \text{ as } X_{ij} \text{ is missing} \end{cases}$$



#### Types of missing data

- MCAR: missing completely at random p(M|X)=p(M)
- MAR: missing at random
   p(M|X)=p(M|X<sub>obs</sub>)
- MNAR: missing not at random





#### A simple exercise

- 1. A sudden traffic violation happened, some of those surveyed submitted questionnaires in a hurry with some missing responses. (MCAR)
- 2. People with bad jobs seldom answer their incomes.(MAR)
- People with low incomes are less likely to report their incomes. (MNAR)





#### Solutions of missing data

• **Deletion(** listwise and pairwise )









#### **Multiple imputations**



#### Inference with multiple imputations







x <sub>i</sub>	4	6	8	9	11	13	16	18	20	25
y <sub>i</sub>	4	20	10		14		16	15	18	22

$$\widehat{Y_R} = \frac{\sum y_i}{\sum x_i} \overline{X} = \widehat{R}\overline{X}$$

#### Step 1

<i>xi</i>	4	6	8	9	11	13	16	18	20	25
$y_{i}^{(1)}$	4	20	10	10	14	14	16	15	18	22
$y_{i}^{(2)}$	4	20	10	16	14	14	16	15	18	22

Step 2

	$\widehat{\overline{Y_R}}$	v <sub>i</sub>
1	13.38	2.965
2	13.57	3.157

Step 3

$$\widehat{Y_R} = \frac{13.38 + 13.57}{2} = 13.48$$

$$v = \frac{1}{2}(2.965 + 3.157) = 3.061$$

$$B = \frac{(13.38 - 13.48)^2 + (13.57 - 13.48)^2}{2 - 1} = 0.0181$$

$$T = 3.061 + (1 + \frac{1}{2}) * 0.0181 = 3.0882$$

### **EM algorithm**











- The EM algorithm formalizes an intuitive idea for obtaining parameter estimates when some of the data are missing:
  - 1. replace missing values by estimated values,
  - 2. estimate parameters.
  - 3. Repeat

Step (1) use estimated parameter values as true values, and

step (2) use estimated values as "observed" values, iterating until convergence.

#### • EM is a method to find $\theta_{ML}$ where

$$\theta_{ML} = \arg \max_{\theta \in \Omega} L(\Theta)$$

- $= \underset{\theta \in \Omega}{\operatorname{arg}} \underset{\theta \in \Omega}{\operatorname{max}} \quad \log P(X \mid \Theta) \Rightarrow \underset{\theta \in \Omega}{\operatorname{arg}} \underset{\theta \in \Omega}{\operatorname{max}} \quad \log P(Z \mid \Theta)$
- $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})$ 
  - ◆Z: complete data ("augmented data")
  - ◆X: observed data ("incomplete" data)
  - ♦Y: hidden data ("missing" data)
  - $\boldsymbol{\blacklozenge} \Theta$ : a parameter vector.



#### The inner loop of the EM algorithm

E-step: Compute expectation of  $(\theta, \theta^t : \text{old}, \text{ new distribution parameters})$ 

$$Q(\theta, \theta^t) = E_{\theta} \{ \log p(z; \theta) \mid x \}$$

M-step: Find  $\theta$  that maximizes Q

$$\theta^{t+1}$$
=arg max Q( $\theta, \theta^t$ ),for all  $\theta$ 

 In particular, when the distribution of the complete-data vector belongs to the exponential family,and the log-likelihood is linear in the sufficient statistic for θ,the E-step reduces to computing the expectation of the complete-data sufficient statistic givien the observed data.





#### **Mixture Model Training**

Let the complete-data vector  $y = (y_1, \dots, y_n)^T$  be a random sample from N(u, $\sigma^2$ ). Suppose  $y_i$ ,  $i = 1, \dots, m$  are observed and  $i = m + 1, \dots, n$  are missing.

$$f(y;u,\sigma^2) = (\frac{1}{2\pi\sigma^2})^{n/2} \exp\{-\frac{1}{2}\sum_{i=1}^n \frac{(y_i - u)^2}{\sigma^2}\}$$

 $(\sum y_i, \sum y_i^2)$  are sufficient statistics for  $\theta = (\mu, \sigma^2)$ 



#### At the t<sup>th</sup> iteration

#### For the E-step, compute

• 
$$E_{\theta}(\sum y_i | y_{obs}) = \sum_{i=1}^m y_i + (n-m)u^{(t)}$$

•  $E_{\theta}(\sum y_i^2 | y_{obs}) = \sum_{i=1}^m y_i^2 + (n-m)(u^{(t)^2} - \sigma^{(t)^2})$ For the M-step

• 
$$\widehat{u} = \frac{1}{n} \sum_{i=1}^{n} y_i$$

• 
$$\hat{\sigma}^2 = (\frac{1}{n} \sum_{i=1}^n y_i)^2 - \frac{1}{n} \sum_{i=1}^n y_i^2$$



#### An example of EM

#### • Let events be "grade in a class"

Get's an A	P(A) = 1/2	a=number 0f A's
Get's a B	P(B)=u	b=number 0f B's
Get's a C	P(C)=2u	c=number 0f C's
Get's a D	P(D) = 1/2 - 3u	d=number 0f D's

Suppose we know the number 0f (A's+B's) =h number 0f C's=c number 0f D's=d



#### Same Problem with Hidden Information

Someone tells us that Number of High grades (A's + B's) = hNumber of C's = c= dNumber of D's

REMEMBER  $P(A) = \frac{1}{2}$  $P(B) = \mu$  $P(C) = 2\mu$  $P(D) = \frac{1}{2} - 3\mu$ 

What is the max. like estimate of  $\mu$  now?

We can answer this question circularly:

If we know the value of  $\mu$  we could compute the Expected value of *a* and *b* Since the ratio a:b should be the same as the ratio  $\frac{1}{2}$ :  $\mu$   $a = \frac{\frac{1}{2}}{\frac{1}{2} + \mu}h$   $b = \frac{\mu}{\frac{1}{2} + \mu}h$ 

#### MAXIMIZATION

EXPECTATION

If we know the expected values of *a* and *b* we could compute the maximum likelihood value of µ

 $\mu = \frac{b+c}{6(b+c+d)}$ 



### E.M. for our Trivial Problem

We begin with a guess for  $\mu$ We iterate between EXPECTATION and MAXIMALIZATION to improve our estimates of  $\mu$  and *a* and *b*. REMEMBER  $P(A) = \frac{1}{2}$   $P(B) = \mu$   $P(C) = 2\mu$  $P(D) = \frac{1}{2} - 3\mu$ 

Define  $\mu(t)$  the estimate of  $\mu$  on the t'th iteration b(t) the estimate of b on t'th iteration  $\mu(0) = initial guess$  $b(t) = \frac{\mu(t)h}{\frac{1}{2} + \mu(t)} = \mathbf{E}[b \mid \mu(t)]$  $\mu(t+1) = \frac{b(t) + c}{6(b(t) + c + d)}$ M-step = max like est of  $\mu$  given b(t)Continue iterating until converged. Good news: Converging to local optimum is assured.

Bad news: I said "local" optimum.



### E.M. Convergence

- Convergence proof based on fact that Prob(data | μ) must increase or remain same between each iteration [NOT OBVIOUS]
- But it can never exceed 1 [OBVIOUS]

So it must therefore converge [OBVIOUS]



## Thank you!



