

Learning

Incomplete Data

Overview

Incomplete Data

- Multiple settings:
 - Hidden variables
 - Missing values
- Challenges
 - Foundational is the learning task well defined?
 - Computational how can we learn with incomplete data?

Why latent variables?

• Model sparsity



Why latent variables?

• Discovering clusters in data

Treating Missing Data

```
Sample sequence: H,T,?,?,H,?,H
```

 Case I: A coin is tossed on a table, occasionally it drops and measurements are not taken

H + H + H

Case II: A coin is tossed, but sometimes tails are not reported

#++++

We need to consider the missing data mechanism

- X = {X₁,...,X_n} are random variables
 Sometimes missing
- **O** = {O₁,...,O_n} are *observability variables* Always observed
- Y = {Y₁,...,Y_n} new random variables
 - $-\operatorname{Val}(Y_i) = \operatorname{Val}(Y_i) \cup \{?\}$
 - Always observed
 - Y_i is a deterministic function of X_i and O_i :

$$Y_i = \begin{cases} X_i & O_{X_i} = o^1 \\ ? & O_{X_i} = o^0 \end{cases}$$

Case I (random missing values)



(a) Random missing values

Case II

(deliberate missing values)



(b) Deliberate missing values

- When can we ignore the missing data mechanism and focus only on the likelihood?
- Missing at Completely at Random (MAR)

A missing data model P_{missing} is missing completely at random (MCAR) if $P_{\text{missing}} \models (X \perp O_X)$



- When can we ignore the missing data mechanism and focus only on the likelihood?
- Missing at Random (MAR)

We say that a missing data model P_{missing} is missing at random (MAR) if for all observations y with $P_{\text{missing}}(y) > 0$, and for all $x_{\text{hidden}}^{y} \in \text{Val}(X_{\text{hidden}}^{y})$, we have that

$$P_{\text{missing}} \vDash \left(o_X \perp x_{\text{hidden}}^{\mathbf{y}} \mid x_{obs}^{\mathbf{y}} \right)$$

where o_X are the specific values of the observation variables given Y.



- When can we ignore the missing data mechanism and focus only on the likelihood?
- Missing at Random (MAR)

We say that a missing data model P_{missing} is missing at random (MAR) if for all observations y with $P_{\text{missing}}(y) > 0$, and for all $x_{\text{hidden}}^{y} \in \text{Val}(X_{\text{hidden}}^{y})$, we have that

$$P_{\text{missing}} \vDash \left(o_X \perp x_{\text{hidden}}^{y} \mid x_{obs}^{y} \right)$$

where o_X are the specific values of the observation variables given **Y**.

Identifiability

• Likelihood can have multiple global maxima



- Example:
 - We can rename the values of the hidden variable H
 - If H has two values, likelihood has two global maxima
- With many hidden variables, there can be an exponential number of global maxima
- Multiple local and global maxima can also occur with missing data (not only hidden variables)

Likelihood for Complete Data

P(X)

X

 X^1

 θ_{x1}

X0

 θ_{x0}

Input Data:



Likelihood decomposes by variables

Likelihood decomposes within CPDs

Likelihood:

 $L(D:\theta) = P(x[1], y[1]) P(x[2], y[2]) P(x[3], y[3])$ P(Y|X)**Y**⁰ X V^1 $= P(x^{0}, y^{0})P(x^{0}, y^{1})P(x^{1}, y^{0})$ **X**0 $\theta_{y0|x0}$ $\theta_{v1|x0}$ **X**1 $\theta_{y0|x1}$ $\theta_{y1|x1}$ $=\theta_{x^0}\theta_{y^0|x^0}\theta_{x^0}\theta_{y^1|x^0}\theta_{x^1}\theta_{y^0|x^1} =$ $\theta_{x^0}^{M[x=0]} (1-\theta_{x^0})^{M[x=1]} \times$ Daphne Koller $\theta_{y^0|x^0}^{M[y=0,x=0]} (1-\theta_{y^0|x^0})^{M[y=1,x=0]} \times$ $\theta_{y^0|x^1}^{M[y=0,x=1]} (1-\theta_{y^0|x^1})^{M[y=1,x=1]}$

Likelihood for Incomplete Data



Multimodal Likelihood



Parameter Correlations

- Total of 8 data points
- Some X's unobserved





Summary

- Incomplete data arises often in practice
- Raises multiple challenges & issues:
 - The mechanism for missingness
 - Identifiability
 - Complexity of likelihood function

Likelihood with Incomplete Data





- Follow gradient of likelihood w.r.t. parameters •
- Line search & conjugate gradient methods for fast • convergence

Gradient Ascent

• Theorem:

Let \mathcal{G} be a Bayesian network structure over \mathcal{X} , and let $\mathcal{D} = \{o[1], \dots, o[M]\}$ be a partially observable data set. Let X be a variable and U its parents in \mathcal{G} . Then

$$\frac{\partial \ell(\boldsymbol{\theta}:\mathcal{D})}{\partial P(x \mid \boldsymbol{u})} = \frac{1}{P(x \mid \boldsymbol{u})} \sum_{m=1}^{M} P(x, \boldsymbol{u} \mid \boldsymbol{o}[m], \boldsymbol{\theta}).$$

 Requires computing P(X_i, U_i|d[m], θ) for all i, m

• Can be done with clique-tree algorithm, since X_i , U_i are in the same clique

Gradient Ascent Summary

- Need to run inference over each data instance at every iteration
- Pros ensure that parameters define legal
 Flexible, can be extended to non table CPDs
- Cons
 - Constrained optimization: need tCPDs
 - For reasonable convergence, need to combine with advanced methods (conjugate gradient, line search)

Expectation Maximization (EM)

- Special-purpose algorithm designed for optimizing likelihood functions
- Intuition
 - Parameter estimation is easy given complete data
 - Computing probability of missing data is "easy" (=inference) given parameters

Example

MLE estimate for $\theta_{d^1|c^0}$ if all data were fully observed:

$$\hat{\theta}_{d^{1}|c^{0}} = \frac{M[d^{1}, c^{0}]}{M[c^{0}]} = \frac{\sum_{m=1}^{M} \mathbb{I}\{\xi[m]\langle D, C\rangle = \langle d^{1}, c^{0}\rangle\}}{\sum_{m=1}^{M} \mathbb{I}\{\xi[m]\langle C\rangle = c^{0}\}}$$

Now assume we have a sample $\boldsymbol{o} = \langle a^1, ?, ?, d^0 \rangle$

Four possible assignments of *b*, *c*

-If we knew the true assignment we could compute the MLE parameters

-If we knew the parameters we could compute the probability of each assingment



$$\begin{array}{ccc} \boldsymbol{\theta}_{a^{1}} & \boldsymbol{\theta}_{b^{1}} \\ \boldsymbol{\theta}_{d^{1}|c^{0}} & \boldsymbol{\theta}_{d^{1}|c^{1}} \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{0}} & \boldsymbol{\theta}_{c^{1}|a^{1},b^{0}} \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{1}} & \boldsymbol{\theta}_{c^{1}|a^{1},b^{1}} \end{array}$$

Example

Assume we are given estimates for the parameters

 $\boldsymbol{o} = \langle a^1, ?, ?, d^0 \rangle$

 $Q(B,C) = P(B,C \mid a^1, d^0, \boldsymbol{\theta})$



$$\begin{aligned} \boldsymbol{\theta}_{a^{1}} &= 0.3 \quad \boldsymbol{\theta}_{b^{1}} &= 0.9 \\ \boldsymbol{\theta}_{d^{1}|c^{0}} &= 0.1 \quad \boldsymbol{\theta}_{d^{1}|c^{1}} &= 0.8 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{0}} &= 0.83 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{0}} &= 0.6 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{1}} &= 0.09 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{1}} &= 0.2 \end{aligned}$$

Example

Assume we are given estimates for the parameters

$$\boldsymbol{o} = \langle a^1, ?, ?, d^0 \rangle$$

$$Q(B,C) = P(B,C \mid a^1, d^0, \boldsymbol{\theta})$$

$$Q(\langle b^{1}, c^{1} \rangle) = 0.3 \cdot 0.9 \cdot 0.2 \cdot 0.2/0.2196 = 0.0492$$

$$Q(\langle b^{1}, c^{0} \rangle) = 0.3 \cdot 0.9 \cdot 0.8 \cdot 0.9/0.2196 = 0.8852$$

$$Q(\langle b^{0}, c^{1} \rangle) = 0.3 \cdot 0.1 \cdot 0.6 \cdot 0.2/0.2196 = 0.0164$$

$$Q(\langle b^{0}, c^{0} \rangle) = 0.3 \cdot 0.1 \cdot 0.4 \cdot 0.9/0.2196 = 0.0492$$



$$\begin{aligned} \boldsymbol{\theta}_{a^1} &= 0.3 \quad \boldsymbol{\theta}_{b^1} &= 0.9 \\ \boldsymbol{\theta}_{d^1|c^0} &= 0.1 \quad \boldsymbol{\theta}_{d^1|c^1} &= 0.8 \\ \boldsymbol{\theta}_{c^1|a^0,b^0} &= 0.83 \quad \boldsymbol{\theta}_{c^1|a^1,b^0} &= 0.6 \\ \boldsymbol{\theta}_{c^1|a^0,b^1} &= 0.09 \quad \boldsymbol{\theta}_{c^1|a^1,b^1} &= 0.2 \end{aligned}$$

Assume we are given estimates for the parameters

 $o' = \langle ?, b^1, ?, d^1 \rangle$

 $Q'(A,C) = P(A,D \mid b^1,d^1,\boldsymbol{\theta})$



$$\begin{aligned} \boldsymbol{\theta}_{a^{1}} &= 0.3 \quad \boldsymbol{\theta}_{b^{1}} &= 0.9 \\ \boldsymbol{\theta}_{d^{1}|c^{0}} &= 0.1 \quad \boldsymbol{\theta}_{d^{1}|c^{1}} &= 0.8 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{0}} &= 0.83 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{0}} &= 0.6 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{1}} &= 0.09 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{1}} &= 0.2 \end{aligned}$$

Assume we are given estimates for the parameters

 $o' = \langle ?, b^1, ?, d^1 \rangle$

 $Q'(A,C) = P(A,D \mid b^1,d^1,\boldsymbol{\theta})$

 $\begin{array}{l} Q'(\langle a^1, c^1 \rangle) = 0.3 \cdot 0.9 \cdot 0.2 \cdot 0.8 / 0.1675 = 0.2579 \\ Q'(\langle a^1, c^0 \rangle) = 0.3 \cdot 0.9 \cdot 0.8 \cdot 0.1 / 0.1675 = 0.1290 \\ Q'(\langle a^0, c^1 \rangle) = 0.7 \cdot 0.9 \cdot 0.09 \cdot 0.8 / 0.1675 = 0.2708 \\ Q'(\langle a^0, c^0 \rangle) = 0.7 \cdot 0.9 \cdot 0.91 \cdot 0.1 / 0.1675 = 0.3423 \end{array}$

This is like having four data instances with weights.



$$\begin{aligned} \boldsymbol{\theta}_{a^{1}} &= 0.3 \quad \boldsymbol{\theta}_{b^{1}} &= 0.9 \\ \boldsymbol{\theta}_{d^{1}|c^{0}} &= 0.1 \quad \boldsymbol{\theta}_{d^{1}|c^{1}} &= 0.8 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{0}} &= 0.83 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{0}} &= 0.6 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{1}} &= 0.09 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{1}} &= 0.2 \end{aligned}$$

 $\begin{array}{l} o = \langle a^1, ?\, , ?\, , d^0\rangle \\ o' = \langle ?\, , b^1, ?\, , d^1\rangle \end{array}$

This is like having the data

$\langle a^1, b^1, c^1, d^0 \rangle$	0.0492
$\langle a^1, b^1, c^0, d^0 \rangle$	0.8852
$\langle a^1, b^0, c^1, d^0 \rangle$	0.0164
$\langle a^1, b^0, c^0, d^0 \rangle$	0.0492

$\langle a^1, b^1, c^1, d^1 \rangle$	0.2579
$\langle a^1, b^1, {\rm c}^0, d^1\rangle$	0.1290
$\langle a^0, b^1, c^1, d^1 \rangle$	0.2708
$\langle a^0, b^1, c^0, d^1 \rangle$	0.3423



$$\begin{aligned} \boldsymbol{\theta}_{a^{1}} &= 0.3 \quad \boldsymbol{\theta}_{b^{1}} &= 0.9 \\ \boldsymbol{\theta}_{d^{1}|c^{0}} &= 0.1 \quad \boldsymbol{\theta}_{d^{1}|c^{1}} &= 0.8 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{0}} &= 0.83 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{0}} &= 0.6 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{1}} &= 0.09 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{1}} &= 0.2 \end{aligned}$$

 $\begin{array}{l} o = \langle a^1, ?\, , ?\, , d^0\rangle \\ o' = \langle ?\, , b^1, ?\, , d^1\rangle \end{array}$

This is like having the data

$\langle a^1, b^1, c^1, d^0 \rangle$	0.0492
$\langle a^1, b^1, {\rm c}^0, d^0\rangle$	0.8852
$\langle a^1, b^0, c^1, d^0 \rangle$	0.0164
$\langle a^1, b^0, c^0, d^0 \rangle$	0.0492

$\langle a^1, b^1, c^1, d^1 \rangle$	0.2579
$\langle a^1, b^1, {\rm c}^0, d^1\rangle$	0.1290
$\langle a^0, b^1, c^1, d^1 \rangle$	0.2708
$\langle a^0, b^1, c^0, d^1 \rangle$	0.3423

We can now compute expected counts,

$$\bar{M}_{\boldsymbol{\theta}}[\boldsymbol{y}] = \sum_{m=1}^{M} \sum_{\boldsymbol{h}[m] \in \operatorname{Val}(\boldsymbol{H}[m])} Q(\boldsymbol{h}[m]) \mathbb{I}\{\boldsymbol{\xi}[m] \langle \boldsymbol{Y} \rangle = \boldsymbol{y}\}$$



$$\begin{aligned} \boldsymbol{\theta}_{a^{1}} &= 0.3 \quad \boldsymbol{\theta}_{b^{1}} &= 0.9 \\ \boldsymbol{\theta}_{d^{1}|c^{0}} &= 0.1 \quad \boldsymbol{\theta}_{d^{1}|c^{1}} &= 0.8 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{0}} &= 0.83 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{0}} &= 0.6 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{1}} &= 0.09 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{1}} &= 0.2 \end{aligned}$$

 $\begin{array}{l} o = \langle a^1, ?\, , ?\, , d^0\rangle \\ o' = \langle ?\, , b^1, ?\, , d^1\rangle \end{array}$

This is like having the data

$\langle a^1, b^1, c^1, d^0 \rangle$	0.0492
$\langle a^1, b^1, \mathbf{c}^0, d^0 \rangle$	0.8852
$\langle a^1, b^0, \mathbf{c}^1, d^0 \rangle$	0.0164
$\langle a^1, b^0, \mathbf{c}^0, d^0 \rangle$	0.0492

$\langle a^1, b^1, \mathbf{c}^1, d^1 \rangle$	0.2579
$\langle a^1, b^1, \mathbf{c}^0, d^1 \rangle$	0.1290
$\langle a^0, b^1, c^1, d^1 \rangle$	0.2708
$\langle a^0, b^1, c^0, d^1 \rangle$	0.3423

Parameters:

$$\begin{aligned} \boldsymbol{\theta}_{a^{1}} &= 0.3 \quad \boldsymbol{\theta}_{b^{1}} &= 0.9 \\ \boldsymbol{\theta}_{d^{1}|c^{0}} &= 0.1 \quad \boldsymbol{\theta}_{d^{1}|c^{1}} &= 0.8 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{0}} &= 0.83 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{0}} &= 0.6 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{1}} &= 0.09 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{1}} &= 0.2 \end{aligned}$$



We can now compute expected counts,

$$\bar{M}_{\boldsymbol{\theta}}[\boldsymbol{y}] = \sum_{m=1}^{M} \sum_{\boldsymbol{h}[m] \in \operatorname{Val}(\boldsymbol{H}[m])} Q(\boldsymbol{h}[m]) \mathbb{I}\{\boldsymbol{\xi}[m] \langle \boldsymbol{Y} \rangle = \boldsymbol{y}\}$$

 $\begin{array}{l} o = \langle a^1, ?\, , ?\, , d^0 \rangle \\ o' = \langle ?\, , b^1, ?\, , d^1 \rangle \end{array}$

We can now compute expected counts, $\bar{M}_{\theta}[\mathbf{y}] = \sum_{m=1}^{M} \sum_{\mathbf{h}[m] \in \operatorname{Val}(\mathbf{H}[m])} Q(\mathbf{h}[m]) \mathbb{I}\{\xi[m] \langle \mathbf{Y} \rangle = \mathbf{y}\}$

$$\begin{split} \overline{M}_{\theta}[d^{1},c^{0}] &= Q'(\langle d^{1},c^{0}\rangle) + Q'(\langle d^{1},c^{0}\rangle) = 0.1290 + 0.3423 = 0.4713 & \text{wrong} \\ formula \text{ in } \\ \overline{M}_{\theta}[c^{0}] &= Q(\langle b^{1},c^{0}\rangle) + Q(\langle b^{0},c^{0}\rangle) + Q'(\langle a^{1},c^{0}\rangle) + Q'(\langle a^{0},c^{0}\rangle) & \text{ the book} \\ &= 0.8852 + 0.0492 + 0.1290 + 0.3423 = 1.4057 \end{split}$$

Once you have the expected counts, you can do MLE estimation and update the parameters

$$\begin{aligned} \boldsymbol{\theta}_{a^{1}} &= 0.3 \quad \boldsymbol{\theta}_{b^{1}} &= 0.9 \\ \boldsymbol{\theta}_{d^{1}|c^{0}} &= 0.1 \quad \boldsymbol{\theta}_{d^{1}|c^{1}} &= 0.8 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{0}} &= 0.83 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{0}} &= 0.6 \\ \boldsymbol{\theta}_{c^{1}|a^{0},b^{1}} &= 0.09 \quad \boldsymbol{\theta}_{c^{1}|a^{1},b^{1}} &= 0.2 \end{aligned}$$

$$\bar{\theta}_{d^1|c^0} = \frac{\bar{M}_{\theta}[d^1, c^0]}{\bar{M}_{\theta}[c^0]}$$

EM Overview

- Pick a starting point for parameters
- Iterate:
 - E-step (Expectation): "Complete" the data using current parameters
 - M-step (Maximization): Estimate parameters relative to data completion
- Guaranteed to improve $L(\theta : D)$ at each iteration

Example: Clustering



Clustering methods

- hard clustering: clusters do not overlap
- soft clustering: clusters may overlap

Mixture models

- probabilistically-grounded way of doing soft clustering
- each cluster: a generative model (Gaussian or multinomial)
- parameters (e.g. mean/covariance are unknown)

Example: Gaussian Mixture models

 $\circ \circ$

- Observations $x_1 \dots x_n$
 - K = 2 Gaussians with unknown μ, σ^2
 - estimation trivial if we know the source of each observation

Example: Gaussian Mixture models

- Observations $x_1 \dots x_n$
 - K = 2 Gaussians with unknown μ, σ^2
 - estimation trivial if we know the source of each observation

$$\mu_b = \frac{x_1 + x_2 + \dots + x_{n_b}}{n_b}$$
$$\sigma_b^2 = \frac{(x_1 - \mu_1)^2 + \dots + (x_n - \mu_n)^2}{n_b}$$

Example: Gaussian Mixture models

- Observations x_1, \dots, x_n
 - K = 2 Gaussians with unknown μ, σ^2
 - estimation trivial if we know the source of each observation

- What if we don't know the source?
- If we knew parameters of the Gaussians (μ, σ^2) , we can guess whether point is more likely to be *a* or *b*.

EM for GMMs

- Chicken and egg problem
 - need (μ_a , σ_a^2) and (μ_b , σ_b^2) to guess source of points
 - need to know source to estimate (μ_a, σ_a^2) and (μ_b, σ_b^2)
- EM algorithm
 - start with two randomly placed Gaussians $(\mu_{a'}\sigma_a^2), (\mu_b, \sigma_b^2)$
 - for each point: $P(b \mid x_i)$ = does it look like it came from b ?
 - adjust $(\mu_{a'}, \sigma_a^2)$ and (μ_b, σ_b^2) to fit points assigned to them



$$P(x_{i} \mid b) = \frac{1}{\sqrt{2\pi\sigma_{b}^{2}}} \exp\left(-\frac{(x_{i} - \mu_{s})^{2}}{2\sigma_{b}^{2}}\right)$$
$$b_{i} = P(b \mid x_{i}) = \frac{P(x_{i} \mid b)P(b)}{P(x_{i} \mid b)P(b) + P(x_{i} \mid a)P(a)}$$

$$a_i = P(a \mid x_i) = 1 - b_i$$



$$P(x_{i} \mid b) = \frac{1}{\sqrt{2\pi\sigma_{b}^{2}}} \exp\left(-\frac{(x_{i} - \mu_{s})^{2}}{2\sigma_{b}^{2}}\right)$$
$$b_{i} = P(b \mid x_{i}) = \frac{P(x_{i} \mid b)P(b)}{P(x_{i} \mid b)P(b) + P(x_{i} \mid a)P(a)}$$

$$a_i = P(a \mid x_i) = 1 - b_i$$



$$P(x_{i} | b) = \frac{1}{\sqrt{2\pi\sigma_{b}^{2}}} \exp\left(-\frac{(x_{i} - \mu_{s})^{2}}{2\sigma_{b}^{2}}\right)$$
$$b_{i} = P(b | x_{i}) = \frac{P(x_{i} | b)P(b)}{P(x_{i} | b)P(b) + P(x_{i} | a)P(a)}$$

$$a_{i} = P(a \mid x_{i}) = 1 - b_{i}$$

$$\mu_{b} = \frac{b_{1}x_{1} + b_{2}x_{2} + \dots + b_{n}x_{n}}{b_{1} + b_{2} + \dots + b_{n}}$$

$$\sigma_{b}^{2} = \frac{b_{1}(x_{1} - \mu_{1})^{2} + \dots + b_{n}(x_{n} - \mu_{n})^{2}}{b_{1} + b_{2} + \dots + b_{n}}$$

$$\mu_{a} = \frac{a_{1}x_{1} + a_{2}x_{2} + \dots + a_{n}x_{n_{n}}}{a_{1} + a_{2} + \dots + a_{n}}$$

$$\sigma_{a}^{2} = \frac{a_{1}(x_{1} - \mu_{1})^{2} + \dots + a_{n}(x_{n} - \mu_{n})^{2}}{a_{1} + a_{2} + \dots + a_{n}}$$



$$P(x_{i} \mid b) = \frac{1}{\sqrt{2\pi\sigma_{b}^{2}}} \exp\left(-\frac{(x_{i} - \mu_{s})^{2}}{2\sigma_{b}^{2}}\right)$$
$$b_{i} = P(b \mid x_{i}) = \frac{P(x_{i} \mid b)P(b)}{P(x_{i} \mid b)P(b) + P(x_{i} \mid a)P(a)}$$

$$a_{i} = P(a \mid x_{i}) = 1 - b_{i}$$

$$\mu_{b} = \frac{b_{1}x_{1} + b_{2}x_{2} + \dots + b_{n}x_{n}}{b_{1} + b_{2} + \dots + b_{n}}$$

$$\sigma_{b}^{2} = \frac{b_{1}(x_{1} - \mu_{1})^{2} + \dots + b_{n}(x_{n} - \mu_{n})^{2}}{b_{1} + b_{2} + \dots + b_{n}}$$

$$\mu_{a} = \frac{a_{1}x_{1} + a_{2}x_{2} + \dots + a_{n}x_{n_{n}}}{a_{1} + a_{2} + \dots + a_{n}}$$

$$\sigma_{a}^{2} = \frac{a_{1}(x_{1} - \mu_{1})^{2} + \dots + a_{n}(x_{n} - \mu_{n})^{2}}{a_{1} + a_{2} + \dots + a_{n}}$$



$$P(x_{i} \mid b) = \frac{1}{\sqrt{2\pi\sigma_{b}^{2}}} \exp\left(-\frac{(x_{i} - \mu_{s})^{2}}{2\sigma_{b}^{2}}\right)$$
$$b_{i} = P(b \mid x_{i}) = \frac{P(x_{i} \mid b)P(b)}{P(x_{i} \mid b)P(b) + P(x_{i} \mid a)P(a)}$$

$$a_{i} = P(a \mid x_{i}) = 1 - b_{i}$$

$$\mu_{b} = \frac{b_{1}x_{1} + b_{2}x_{2} + \dots + b_{n}x_{n}}{b_{1} + b_{2} + \dots + b_{n}}$$

$$\sigma_{b}^{2} = \frac{b_{1}(x_{1} - \mu_{1})^{2} + \dots + b_{n}(x_{n} - \mu_{n})^{2}}{b_{1} + b_{2} + \dots + b_{n}}$$

$$\mu_{a} = \frac{a_{1}x_{1} + a_{2}x_{2} + \dots + a_{n}x_{n_{n}}}{a_{1} + a_{2} + \dots + a_{n}}$$

$$\sigma_{a}^{2} = \frac{a_{1}(x_{1} - \mu_{1})^{2} + \dots + a_{n}(x_{n} - \mu_{n})^{2}}{a_{1} + a_{2} + \dots + a_{n}}$$

Expectation-Maximization

Iterate until convergence: On the t – th iteration let our estimates be

 $\lambda_t = \{\mu_1(t), \mu_2(t) \dots \mu_c(t)\} \qquad \begin{array}{l} \text{Just evaluate a} \\ \text{Gaussian at } x_k \end{array}$

E-step: Compute "expected" classes of all datapoints for each class

$$P(z_i \mid x_k, \lambda_t) = \frac{p(x_k \mid z_i, \lambda_t) P(z_i \mid \lambda_t)}{p(x_k \mid \lambda_t)} = \frac{p(x_k \mid z_i, \mu_i(t), \sigma^2 \mathbf{I}) p_i(t)}{\sum_{j=1}^c p(x_k \mid z_j, \mu_j(t), \sigma^2 \mathbf{I}) p_j(t)}$$

M-step: Estimate μ given our data's class membership distributions

$$\mu_i(t+1) = \frac{\sum_k P(z_i \mid x_k, \lambda_t) x_k}{\sum_k P(z_i \mid x_k, \lambda_t)}$$

Example: Gaussian Mixture Models



Example: Gaussian Mixture Models



EM Summary

- Need to run inference over each data instance at every iteration
- Pros
 - Easy to implement on top of MLE for complete data
 - Makes rapid progress, especially in early iterations
- Cons

- Convergence slows down at later iterations