

Lecture 22

Expectation Maximization (EM)

ECEN 5283 Computer Vision

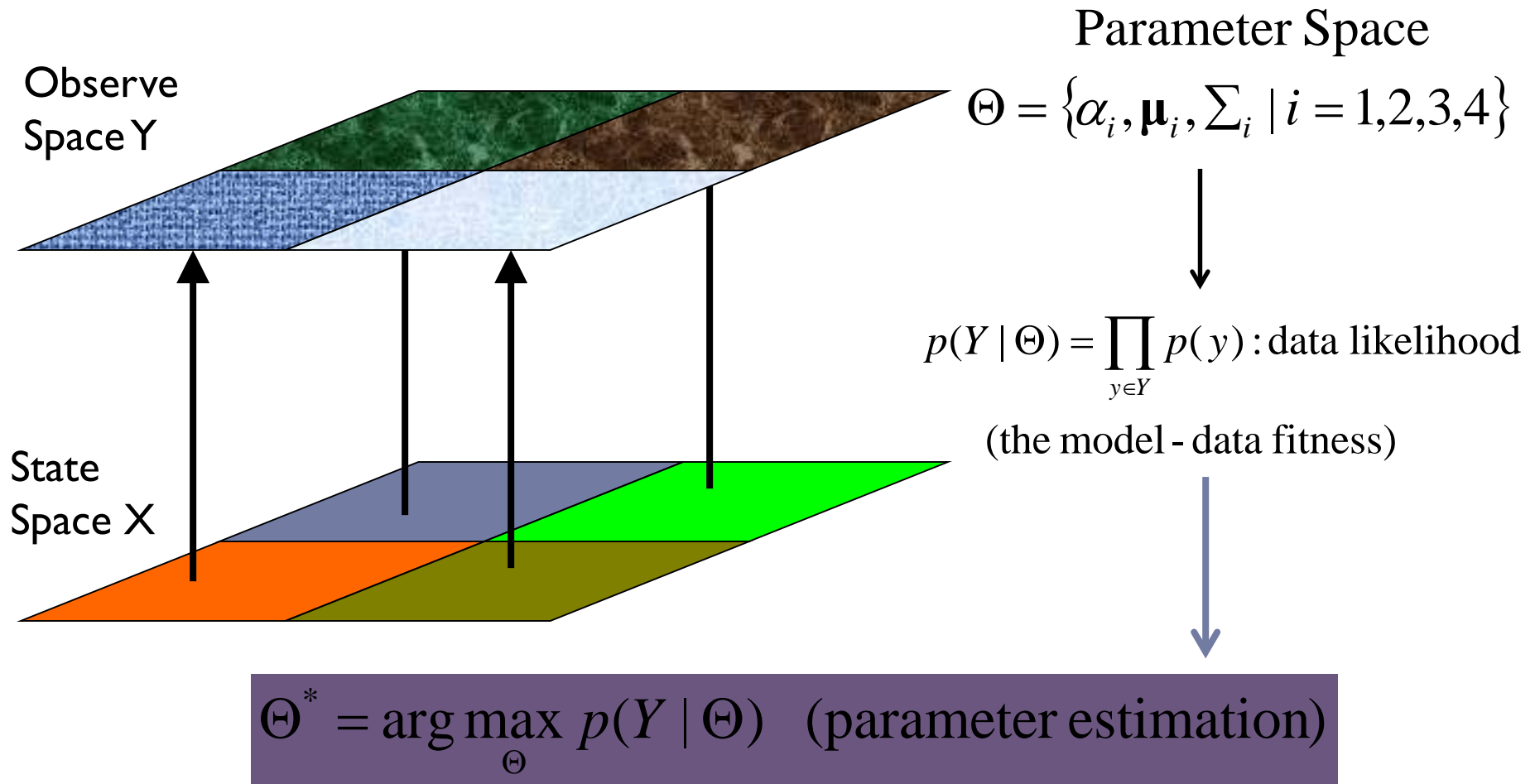
Dr. Guoliang Fan
School of Electrical and Computer Engineering
Oklahoma State University



Goals

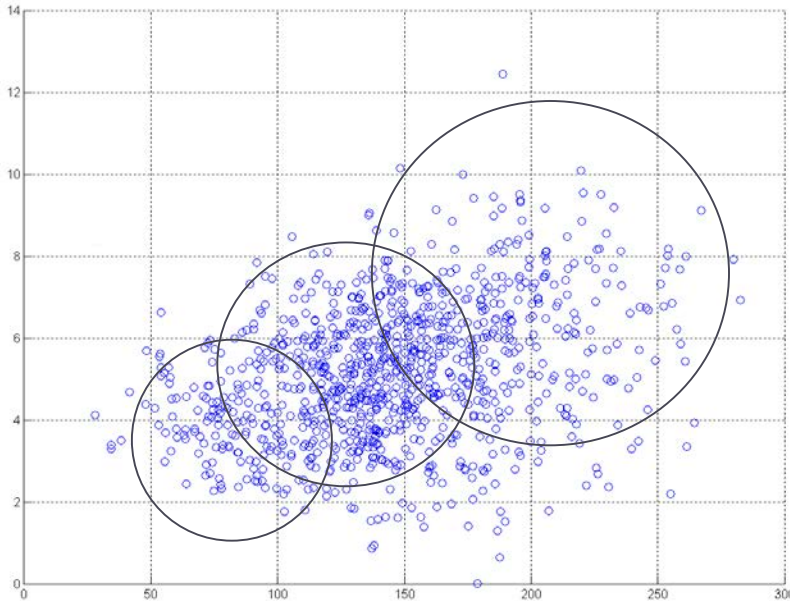
- ▶ To review the missing data problem and its two major issues.
- ▶ To introduce a soft-clustering algorithm, i.e., Expectation Maximization (EM) algorithm.

Issue (1) Parameter Estimation



To find the parameter Θ that can best explain the current observation Y .

Parameter Estimation Example



Labels and Data

$$(X, Y) = ((x_l, y_l) \mid l = 1, \dots, N)$$

x_l : label, y_l : observation

Parameter Space

$$\Theta = \{\alpha_i, \boldsymbol{\mu}_i, \Sigma_i \mid i = 1, 2, 3\}$$



$$\Theta^* = \arg \max_{\Theta} p(Y \mid \Theta)$$

(maximum likelihood estimation)

$\alpha_i = P(x = i)$
(prior probability)

$p(y \mid x = i) = N(y \mid \boldsymbol{\mu}_i, \Sigma_i)$
(likelihood function)

Missing Data Problem:

Issue (2) Data Classification



- ▶ After parameter estimation, we need to decide the class label of each data sample according to estimated parameters.

$$\Theta = \{\alpha_i, \mu_i, \Sigma_i \mid i = 1, \dots, k\}$$
$$Y = (y_1, y_2, \dots, y_N) \xrightarrow{\hspace{2cm}} X = (x_1, x_2, \dots, x_N)$$

- ▶ We need to compute **the posterior probability of data sample y belonging to class x** . (This is the estimate of the missing data)

$$\alpha_i = \Pr(x = i) = p(x)$$

(prior probability)

$$p(y \mid x = i) = N(y \mid \mu_i, \Sigma_i)$$

(likelihood function)

$$p(x \mid y) = \frac{p(x, y)}{p(y)} = \frac{p(y \mid x)p(x)}{\sum_{i=1}^k p(y \mid x = i)p(x = i)} \quad \text{(posterior probability)}$$

$$x^* = \arg_{x \in X} \max p(x \mid y) \quad \text{(maximum *a posteriori* or MAP)}$$

Data Classification Example

- ▶ If we have **six samples** and **three classes**, the missing data indicates the class label for each pixel. Hopefully, the estimated missing data will be close it.

$$\mathbf{I}^0 = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}_{6 \times 3}$$

The true missing data

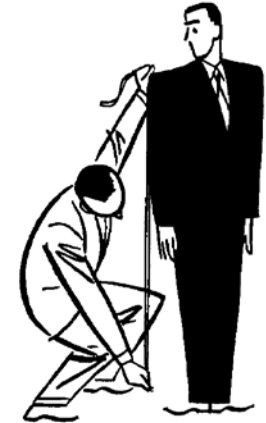
$$\mathbf{I} = \begin{pmatrix} 0.1 & 0.8 & 0.1 \\ 0.3 & 0.6 & 0.1 \\ 0.4 & 0.3 & 0.3 \\ 0.7 & 0.2 & 0.1 \\ 0.3 & 0.2 & 0.5 \\ 0.05 & 0.05 & 0.9 \end{pmatrix}_{6 \times 3}$$

The estimated missing data

$$x_l = \arg_{m \in \{1, \dots, g\}} \max \mathbf{I}(l, m)$$

EM Formulation: Objective Function

- ▶ *How does a tailor make a cloth?*
 - ▶ To make a cloth that fits best to the body
- ▶ *How to estimate the parameters Θ ?*
 - ▶ Estimating Θ that best fits the data:



$$Y = (y_1, y_2, \dots, y_N)$$

- ▶ *How to evaluate the fitness between Θ and Y ?*
 - ▶ The fitness between Θ (model) and Y (data) is reflected by the likelihood of Y given Θ . Therefore, parameter estimation is:

$$\Theta^* = \arg \max_{\Theta} p(Y | \Theta) = \arg \max_{\Theta} \prod_{i=1}^N p(y_i | \Theta) = \arg \max_{\Theta} \sum_{i=1}^N \log(p(y_i | \Theta))$$

$$p(Y | \Theta) = \prod_{i=1}^N p(y_i | \Theta) \text{ (under the independent assumption)}$$

EM Formulation: Data log-likelihood



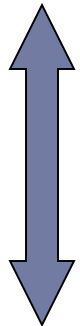
$$\begin{aligned}\log p(Y | \Theta) &= \log \prod_{j=1}^N p(y_j | \Theta) = \sum_{j=1}^N \log p(y_j | \Theta) \\&= \sum_{j=1}^N \log \left(\sum_{i=1}^k p(y_j, x_j = i | \Theta) \right) \\&= \sum_{j=1}^N \log \left(\sum_{i=1}^k p(y_j | x_j = i, \Theta) p(x_j = i | \Theta) \right) \\&= \sum_{j=1}^N \log \left(\sum_{i=1}^k \underbrace{p(y_j | x_j = i, \Theta)}_{\text{Likelihood function}} \underbrace{\alpha_i}_{\text{prior}} \right)\end{aligned}$$

EM Formulation: Likelihood Function



- ▶ What is a likelihood function?
 - ▶ The likelihood function indicates how likely a particular distribution is to produce an observed sample. It is like a ruler for the tailor.
- ▶ What it can do for EM?
 - ▶ It allows us to estimate unknown parameters based on known outcomes.

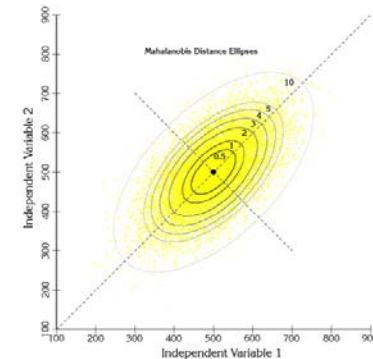
$$p(y | \theta_m) = \frac{1}{(2\pi)^{d/2} \det(\Sigma_m)^{1/2}} \exp \left\{ -\frac{1}{2} (y - \mu_m)^T \Sigma_m^{-1} (y - \mu_m) \right\}$$



The dimension of y

$$p(y | \theta_m) = p(y | x = m) \text{ (likelihood function)}$$

Mahalanobis Distance





EM Formulation: Two-Step Iteration

- ▶ We use a two-step iteration to solve the missing data problem
 - ▶ Initialize the parameters (it is like to initialize the centers in k-means)

$$\Theta = \{\alpha_i, \mu_i, \Sigma_i \mid i = 1, \dots, k\}$$

- ▶ Step 1: Estimate the missing data (x) in terms of the posterior probability of each data sample (y) (it is like to classify each data point in k-means.)

$$p(x = i \mid y) \quad \{i = 1, \dots, k\}$$

(estimate the probability of each sample belonging to different classes)

- ▶ Step 2: From the estimated missing data, to obtain the maximum likelihood estimate of the parameters (it is like to update the centers in k-means)

$$\Theta^* = \arg \max_{\Theta} \log p(Y \mid \Theta)$$

(update the parameters to better fit the data and the model.)

EM Algorithm: E-step

- **Initialization:** set $s=0$ and

$$\Theta^0 = (\alpha_1^{(0)}, \alpha_2^{(0)}, \dots, \alpha_g^{(0)}, \theta_1^{(0)}, \theta_2^{(0)}, \dots, \theta_g^{(0)}).$$

- **Expectation (E-step):**

$$I(l, m) = \frac{\alpha_m^{(s)} p(y_l | \theta_m^{(s)})}{\sum_{i=1}^k \alpha_i^{(s)} p(y_l | \theta_i^{(s)})} = p(x_l = m | y_l, \Theta^{(s)})$$

Posterior probability

$$p(x | y) = \frac{p(x, y)}{p(y)} = \frac{p(y | x) p(x)}{\sum_x p(y | x) p(x)}$$

$$p(y_l | \theta_m^{(s)}) = \frac{\exp\left\{-\frac{1}{2}(y_l - \mu_m)^T \Sigma_m^{-1}(y_l - \mu_m)\right\}}{(2\pi)^{d/2} \det(\Sigma_l)^{1/2}}$$

Likelihood function

EM Algorithm: M-step

► Maximization (M-step): $\Theta^* = \arg \max_{\Theta} p(Y | \Theta)$

$$\frac{\partial \log p(\mathbf{Y} | \Theta^{(s+1)})}{\partial \alpha_m} = 0 \longrightarrow \alpha_m^{(s+1)} = \frac{1}{N} \sum_{l=1}^N p(x_l = m | y_l, \Theta^{(s)})$$

$$\frac{\partial \log p(\mathbf{Y} | \Theta^{(s+1)})}{\partial \mu_m} = 0 \longrightarrow \mu_m^{(s+1)} = \frac{\sum_{l=1}^N y_l p(x_l = m | y_l, \Theta^{(s)})}{\sum_{l=1}^N p(x_l = m | y_l, \Theta^{(s)})}$$

$$\frac{\partial \log p(\mathbf{Y} | \Theta^{(s+1)})}{\partial \Sigma_m} = 0 \longrightarrow \Sigma_m^{(s+1)} = \frac{\sum_{l=1}^N p(x_l = m | y_l, \Theta^{(s)}) \{ (y_l - \mu_m^{(s)}) (y_l - \mu_m^{(s)})^T \}}{\sum_{l=1}^N p(x_l = m | y_l, \Theta^{(s)})}$$