



EE 6885 Statistical Pattern Recognition

Fall 2005
Prof. Shih-Fu Chang
<http://www.ee.columbia.edu/~sfchang>

Lecture 6 (9/28/05)

EE6887-Chang

6-1

■ Reading

- EM for Missing Features
 - Textbook, DHS 3.9
- Bayesian Parameter Estimation and Sufficient Statistics
 - Textbook, DHS 3.4-3.6
- Application: Bayesian Face Detection
 - Reference paper (see course web site)

EE6887-Chang

6-2

Parameter Estimation: ML and EM

- Given Data D

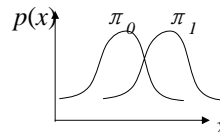
$$\text{Find } \hat{\theta} = \arg \max_{\theta} p(D | \theta)$$

- Gaussian

$$\Rightarrow \hat{\mu} = (1/n) \sum_k \bar{x}_k \quad \hat{\Sigma} = (1/n) \sum_k (\bar{x}_k - \mu)(\bar{x}_k - \mu)^t$$

- Mixture of Gaussian

$$l = \sum_{n=1}^N \log(\pi_0 N(x_n | \mu_0, \Sigma_0) + \pi_1 N(x_n | \mu_1, \Sigma_1))$$



- Derive Auxiliary Function

$$Q(\theta | \theta_t) = \sum_{n=1}^N \sum_z \underbrace{p(z | x_n, \theta_t)} \underbrace{\log p(x_n, z | \theta)}$$

$$\theta_{t+1} = \arg \max_{\theta} Q(\theta | \theta_t)$$

EE6887-Chang

6-3

EM for Missing Feature

$D = \{D_g, D_b\}$ D_g : good feature, D_b : missing feature

$$Q(\theta; \theta^i) = E_{D_b} [\ln p(D_g, D_b; \theta) | D_g; \theta^i] \quad \text{Marginalize over the missing feature}$$

- Example, 2-D Gaussian

$$D = \{x_1, x_2, x_3, x_4\} = \left\{ \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} * \\ 4 \end{pmatrix} \right\} \quad \theta = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \sigma_1^2 \\ \sigma_2^2 \end{bmatrix} \quad \theta^0 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

$$Q(\theta; \theta^0) = E_{x_{41}} [\ln p(x_g, x_{41}; \theta) | x_g; \theta^0]$$

$$= \int_{-\infty}^{\infty} \left[\sum_{k=1}^3 [\ln p(x_k | \theta)] + \ln p(x_4 | \theta) \right] p(x_{41} | \theta^0; x_{42} = 4) dx_{41}$$

$$= \sum_{k=1}^3 [\ln p(x_k | \theta)] + \int_{-\infty}^{\infty} \ln p\left(\begin{matrix} x_{41} \\ 4 \end{matrix} \middle| \theta\right) \frac{p\left(\begin{matrix} x_{41} \\ 4 \end{matrix} \middle| \theta^0\right)}{\int_{-\infty}^{\infty} p\left(\begin{matrix} x_{41} \\ 4 \end{matrix} \middle| \theta^0\right) dx_{41}} dx_{41} \quad \theta^i = \begin{bmatrix} 0.75 \\ 2.0 \\ 0.938 \\ 2.0 \end{bmatrix} \quad \theta = \begin{bmatrix} 1.0 \\ 2.0 \\ 0.5 \\ 2.0 \end{bmatrix}$$

- E step and M step?

See figure in Sec. 3.9

EE6887-Chang

6-4

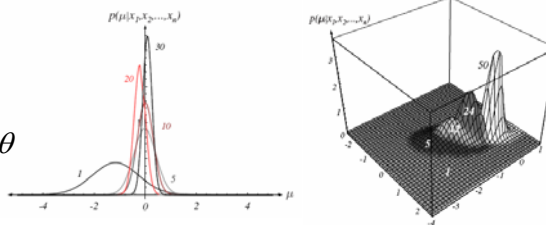
Bayesian Parameter Estimation

- Instead of fixed unknown, we assume θ as random variable with distribution $p(\theta)$
- Given samples D , find the maximal posterior

$$p(\theta | D) = \frac{p(D | \theta) p(\theta)}{\int p(D | \theta) p(\theta) d\theta} \quad p(D | \theta) = \prod_{k=1}^n p(x_k | \theta)$$

- Posterior gives the probability distribution of θ . When the sample size increases, the distribution becomes peaky.
- Difference from ML estimation?

$$p(x | D) = \int p(x | \theta) p(\theta | D) d\theta$$



EE6887-Chang

6-5

Example: Gaussian

- Univariate case: μ is the only unknown

$$p(x | \mu) \sim N(\mu, \sigma^2) \quad p(\mu) \sim N(\mu_0, \sigma_0^2)$$

$$\text{posterior } p(\mu | D) \propto p(D | \mu) p(\mu)$$

$$= \prod_{k=1}^n p(x_k | \mu) p(\mu)$$

$$= \prod_{k=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x_k - \mu}{\sigma}\right)^2\right] \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right]$$

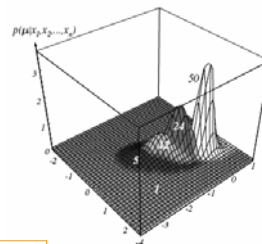
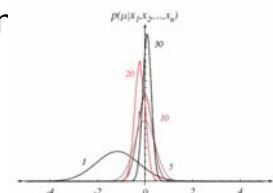
$$= \frac{1}{\sqrt{2\pi}\sigma_n} \exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_n}{\sigma_n}\right)^2\right] \sim N(\mu_n, \sigma_n^2)$$

Reproducing density

$$\mu_n = \left(\frac{n\sigma_0^2}{n_0\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \cdot \mu_0 \quad \text{where } \hat{\mu}_n = \frac{1}{n} \sum_{k=1}^n x_k$$

$$\text{and } \sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}$$

what if $n \uparrow, \sigma_0 \uparrow, \sigma \uparrow$?



EE6887-Chang

6-6

Example: Gaussian Case (Cont.)

- Class conditional density

$$p(x|D) = \int p(x|\mu)p(\mu|D)d\theta$$

$$p(x|D) \sim N(\mu_n, \sigma^2 + \sigma_n^2) \quad \text{vs.} \quad p(x|\mu) \sim N(\mu, \sigma^2)$$

- Mean is replaced, variance is increased.
- Bayesian Classifier: compute posterior $p(x|D_i, \omega_i)p(\omega_i)$

- Multi-variate case

$$\text{posterior } p(\mu|D) \sim N(\mu_n, \Sigma_n)$$

$$\mu_n = \Sigma_0(\Sigma_0 + \frac{1}{n}\Sigma)^{-1} \hat{\mu}_n + \frac{1}{n}\Sigma(\Sigma_0 + \frac{1}{n}\Sigma)^{-1}\mu_0 \quad \text{where} \quad \hat{\mu}_n = \frac{1}{n}\sum_{k=1}^n x_k$$

$$\Sigma_n = \Sigma_0(\Sigma_0 + \frac{1}{n}\Sigma)^{-1} \frac{1}{n}\Sigma$$

- Class conditional density $p(x|D) \sim N(\mu_n, \Sigma + \Sigma_n)$

EE6887-Chang

6-7

Sufficient Statistics

- sufficient statistics \mathbf{s}** : function of samples that contains all of the information relevant to estimation of parameters θ
- Example: Gaussian $N(\theta, \Sigma)$ θ is unknown

$$\begin{aligned} P(D|\theta) &= \prod_{k=1}^n \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x}_k - \theta)' \Sigma^{-1}(\mathbf{x}_k - \theta)\right) \\ &= \frac{1}{(2\pi)^{nd/2} |\Sigma|^{n/2}} \exp\left(-\frac{1}{2} \sum_{k=1}^n (\theta' \Sigma^{-1} \theta - 2\theta' \Sigma^{-1} \mathbf{x}_k + \mathbf{x}_k' \Sigma^{-1} \mathbf{x}_k)\right) \\ &= \exp\left[\frac{-n}{2} \theta' \Sigma^{-1} \theta + \theta' \Sigma^{-1} \left(\sum_{k=1}^n \mathbf{x}_k\right)\right] \times \frac{1}{(2\pi)^{nd/2} |\Sigma|^{n/2}} \exp\left(-\frac{1}{2} \sum_{k=1}^n (\mathbf{x}_k' \Sigma^{-1} \mathbf{x}_k)\right) \end{aligned}$$

$$g(\hat{\mu}_n, \theta) = \exp\left[\frac{-n}{2} (\theta' \Sigma^{-1} \theta - 2\theta' \Sigma^{-1} \hat{\mu}_n)\right] \quad h(D) \text{ independent of } \theta$$

$$\text{where } \hat{\mu}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k \quad \text{Sufficient statistics}$$

EE6887-Chang

6-8

Sufficient Statistics for Exponential Family

- Gaussian, exponential, Rayleigh, Poisson, etc.

$$p(x | \theta) = \alpha(x) \exp[a(\theta) + b(\theta)' c(x)]$$

$$p(D | \theta) = \prod_{k=1}^n \alpha(x_k) \exp[na(\theta) + b(\theta)' \sum_{k=1}^n c(x_k)] = g(s, \theta) h(D)$$

$$\text{where } s = \frac{1}{n} \sum_{k=1}^n c(x_k)$$

$$g(s, \theta) = \exp[n(a(\theta) + b(\theta)' s)]$$

$$h(D) = \prod_{k=1}^n \alpha(x_k)$$

- See Table 3.1 for the S.S. for different distributions.

EE6887-Chang

6-9

MAP detectors for face images

- H. Schneiderman & T. Kanade, "Probabilistic modeling of local appearance and spatial relationships for object recognition," CVPR '98.

$$\frac{P(\text{region} | \text{object})}{P(\text{region} | \overline{\text{object}})} > \lambda = \frac{P(\overline{\text{object}})}{P(\text{object})}$$

- Capture face information from 64x64 regions. Enough details for characterizing faces. Decompose into 16x16 subregions to simplify models and capture spatial locations.

$$P(\text{region} | \text{object}) = P(\{\text{subregion}\} | \text{object})$$

$$P(\{\text{subregion}\} | \text{object}) \approx \prod_{k=1}^{n_{\text{subs}}} P(\text{subregion}_k | \text{object})$$



- Does not consider subregion dependency. Less penalty for face detection due to consistent appearance and location of salient features, e.g., noises and eyes.
- Also capture the distinctive local appearances in the likelihood ratio

$$\prod_{k=1}^{n_{\text{subs}}} \frac{P(\text{subregion}_k | \text{object})}{P(\text{subregion}_k | \overline{\text{object}})}$$

EE6887-Chang

6-10

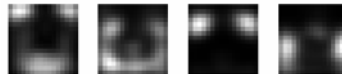
Bayesian face detector (Cont.)

- Simplify the subregion patterns by PCA projection, sparse coding, and quantization → 3.8M distinctive patterns of q_1

$$P(\text{pattern}, \text{pos} | \text{object}) \approx P(q_1, \text{pos} | \text{object})$$

- Estimate the location distributions from training data. $P(\text{pos} | q_1, \text{object})$

- Reduce the complexity of the above distributions by pattern clustering $P(\text{pos}' | q_1, \text{object}) \approx P(\text{pos}' | q_2, \text{object})$
 $q_2 = \text{VQ}(q_1)$



- Clustering helps discover salient locations

- Simple model for background and infrequent patterns

$$p(\text{pos} | q_1, \overline{\text{object}}) \approx \frac{1}{n_{\text{subs}}}$$

- Multi-resolution scaling to capture the face features at different scales



$$P(\text{region} | \text{object}) \approx \prod_{j=1}^{n_{\text{magn}}} P(\text{region}^j | \text{object})$$

EE6887-Chang

6-11

Bayesian face detector (Cont. 2)

- Final Bayesian Estimation

$$\prod_{j=1}^{n_{\text{magn}}} \prod_{i=1}^{n_{\text{subs}}} \frac{P(q_1^j | \text{object}) P(\text{pos}_i^j | q_2^j, \text{object})}{P(q_1^j | \overline{\text{object}})} > \lambda = \frac{P(\overline{\text{object}})}{P(\text{object})}$$

- Estimate class specific distribution using training frontal face and non-face images → simple frequency count

- For test images, search faces at multiple scales, e.g., 17 scales over 18x18 to 338x338.

Schneiderman & Kanade(125 images)		Rowley, Baluja, and Kanade [5] (130 images)	
Detection rate	False alarms	Detection rate	False alarms
93.0%	88	92.5%	862
90.5%	33	86.6%	79
77.0%	1	77.9%	2



- Key features:
 - joint local appearance-location statistical modeling
 - background modeling
 - discrete non-parameter distribution modeling

EE6887-Chang

6-12