



EE 6885 Statistical Pattern Recognition

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

Lecture 5 (9/21/05)

EE6887-Chang

4-1

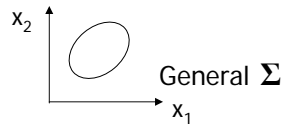
- Reading
 - Model Parameter Estimation
 - ML Estimation, Chap. 3.2
 - Mixture of Gaussian and EM
 - Reference Book, HTF Chap. 8.5
 - Textbook, DHS 3.9
- Homework #2 due 2005-09-28, Wed
- No class/office hours next Monday, 2005-09-26

EE6887-Chang

5-2

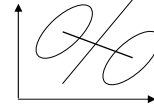
Multi-variate Gaussian

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} e^{\left(\frac{-1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)\right)}$$



Bayesian Classifiers

Decision Boundaries for Gaussians



$$w^T(x - x_0) = 0 \quad w = \Sigma^{-1}(\mu_i - \mu_j) \quad x_0 = \frac{1}{2}(\mu_i + \mu_j) - \frac{\ln[P(\omega_i)/P(\omega_j)]}{(\mu_i - \mu_j)^T \Sigma^{-1}(\mu_i - \mu_j)} (\mu_i - \mu_j)$$

Missing Features by Marginalization

- $x = [x_g, x_b]$, x_g : good features, x_b : bad features

$$\text{compute } P(w_i | x_g) = \frac{p(w_i, x_g)}{p(x_g)} = \frac{\int p(w_i, x_g, x_b) dx_b}{p(x_g)} = \frac{\int p(x_g, x_b | w_i) p(w_i) dx_b}{\int p(x_g, x_b) dx_b}$$

Joint prob. of (ω_i, x_g, x_b) marginalized over x_b

EE6887-Chang

5-3

Parameter Estimation

Parametric form of distribution

e.g., $p(x | w_i) \sim N(\mu_i, \Sigma_i) \quad p(x | w_i) = p(x | w_i, \theta_i)$

How to estimate θ_i ?

→ learn from data samples $D = \{x_1, x_2, \dots, x_n\}$

Likelihood $l(\theta) = p(D | \theta) = \prod_{k=1}^n p(x_k | \theta)$ assume x_1, \dots, x_n independent

$$\begin{aligned} \text{Find } \hat{\theta} &= \arg \max_{\theta} p(D | \theta) = \arg \max_{\theta} \prod_{k=1}^n p(x_k | \theta) \\ &= \arg \max_{\theta} \sum_1^n \ln p(x_k | \theta) \end{aligned}$$

Use gradient operation

$$\nabla_{\theta} = \begin{bmatrix} \partial / \partial \theta_1 \\ \dots \\ \partial / \partial \theta_p \end{bmatrix} \quad \nabla_{\theta} l(\theta) = 0$$

EE6887-Chang

5-4

Case I: Gaussian: μ unknown

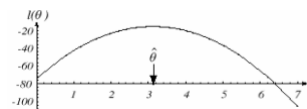
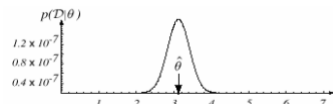
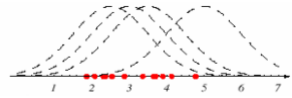
$$\ln P(x_k | \mu) = -\frac{1}{2} \ln[(2\pi)^d |\Sigma|] - \frac{1}{2} (x_k - \mu)^t \Sigma^{-1} (x_k - \mu)$$

$$\text{and } \nabla_{\theta\mu} \ln P(x_k | \mu) = \Sigma^{-1} (x_k - \mu)$$

$$\nabla_{\theta\mu} \sum_k \ln P(x_k | \mu) = 0$$

$$\sum_k \Sigma^{-1} (x_k - \mu) = 0$$

$$\hat{\mu} = \frac{\sum_k x_k}{n}$$



ML estimator of the Gaussian mean is the sample mean

EE6887-Chang

5-5

Case II: Gaussian Case: *unknown* μ and σ (1D)

$$\theta = (\mu, \sigma^2) \quad l = \ln P(x_k | \theta) = -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (x_k - \mu)^2$$

$$\nabla_{\theta} l = \begin{pmatrix} \frac{\partial}{\partial \mu} (\ln P(x_k | \theta)) \\ \frac{\partial}{\partial \sigma^2} (\ln P(x_k | \theta)) \end{pmatrix} = \begin{bmatrix} \frac{1}{\sigma^2} (x_k - \mu) \\ -\frac{1}{2\sigma^2} - \frac{(x_k - \mu)^2}{2(\sigma^2)^2} \end{bmatrix}$$

$$\sum_k \nabla_{\theta} \ln P(x_k | \theta) = 0 \quad \Rightarrow \quad \hat{\mu} = \frac{\sum_k x_k}{n} \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})^2$$

$$\text{Exp}(\hat{\sigma}^2) = \frac{n-1}{n} \sigma^2 \neq \sigma^2 \quad \text{if } n \rightarrow \infty \text{ then } \text{Exp}(\hat{\sigma}^2) \rightarrow \sigma^2$$

Asymptotically unbiased

- Multi-Dimensional $\theta = (\mu, \Sigma)$

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

ML estimator: mean \rightarrow sample mean, variance \rightarrow biased sample variance

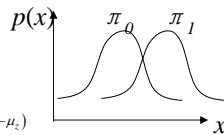
EE6887-Chang

5-6

Mixture Of Gaussians

- Real distributions seldom follow a single Gaussian
→ mixture of Gaussians

$$\begin{aligned}
 p(x) &= \sum_z p(x, z) = \sum_z p(z) p(x|z) \\
 &= \sum_z \pi_z N(x|\mu_z, \Sigma_z) = \sum_{z=1}^Z \pi_z \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma_z|}} e^{-\frac{1}{2}(x-\mu_z)^T \Sigma_z^{-1} (x-\mu_z)}
 \end{aligned}$$



- Given data x_1, \dots, x_N , define log-likelihood:

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

- Posterior probability of x being generated by a specific component

$$\underset{\substack{\text{(responsibility of} \\ \text{component } i)}}}{\text{posterior}} = \tau^i = p(z_i = 1|x, \theta), \quad \theta = \{\mu_0, \Sigma_0, \mu_1, \Sigma_1\}$$

EE6887-Chang

5-7

Derivation of the E-M solution

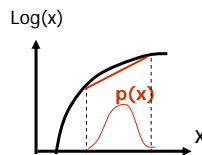
$$\boxed{\text{log likelihood}} \quad l(\theta) = \sum_{n=1}^N \log \sum_z p(x_n, z|\theta)$$

- Maximization of $l(\theta)$ directly is hard due to log_of_sum
- Instead, look at

$$\Delta l(\theta) = l(\theta) - l(\theta_i), \quad \theta_i : \text{current estimation of } \theta$$

- Jensen's Inequality

$$\begin{aligned}
 \text{If } f \text{ is concave, } f(\mathbb{E}\{x\}) &\geq \mathbb{E}\{f\{x\}\} \\
 f(\mathbb{E}\{g(x)\}) &\geq \mathbb{E}\{f\{g(x)\}\} \\
 \text{e.g., } f(x) &= \log(x)
 \end{aligned}$$



$$\log\left(\sum_i p_i x_i\right) \geq \sum_i p_i \log(x_i), \text{ where } \sum_i p_i = 1$$

$$\text{If } f \text{ is convex, } f(\mathbb{E}\{x\}) \leq \mathbb{E}\{f\{x\}\}$$

EE6887-Chang

5-8

Auxiliary Function in E-M

$$\begin{aligned}
 \Delta l(\theta) &= l(\theta) - l(\theta_t) = \sum_{n=1}^N \log p(x_n | \theta) - \sum_{n=1}^N \log p(x_n | \theta_t) \\
 &= \sum_{n=1}^N \log \frac{p(x_n | \theta)}{p(x_n | \theta_t)} = \sum_{n=1}^N \log \sum_z \frac{p(x_n, z | \theta)}{p(x_n | \theta_t)} \quad \text{marginalization} \\
 &= \sum_{n=1}^N \log \sum_z \frac{p(x_n, z | \theta)}{p(x_n | \theta_t)} \frac{p(x_n, z | \theta_t)}{p(x_n, z | \theta_t)} \\
 &= \sum_{n=1}^N \log \sum_z p(z | x_n, \theta_t) \frac{p(x_n, z | \theta)}{p(x_n, z | \theta_t)} \\
 &\geq \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log \frac{p(x_n, z | \theta)}{p(x_n, z | \theta_t)} \quad \text{Jensen's inequality} \\
 &= Q(\theta | \theta_t)
 \end{aligned}$$

■ Note there is no log_of_sum.
 So taking derivative is easier

EE6887-Chang

5-9

E-M improves likelihood

- Auxiliary function derived based on Jensen's Inequality,

$$Q(\theta | \theta_t) = \sum_{n=1}^N \sum_z \underbrace{p(z | x_n, \theta_t)}_{\text{expectation over } z \text{ with current } \theta_t} \log \underbrace{p(x_n, z | \theta)}_{\text{joint likelihood of observed \& hidden}} + \text{const}$$

- Now estimate θ_{t+1} by maximizing Q

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

- So in the expectation step, compute τ_n^z , the 'responsibility' of component z for sample x_n
- In the maximization step, take derivative of Q to θ , and find the new estimate for θ (Note only sum_of_log is involved)

EE6887-Chang

5-10

EM Always Improves Likelihood

- Why does EM always improve $l(\theta)$?

$$\Delta l(\theta_{t+1}) = l(\theta_{t+1}) - l(\theta_t) \geq Q(\theta_{t+1} | \theta_t)$$

$$Q(\theta_{t+1} | \theta_t) = \max_{\theta} Q(\theta | \theta_t) \geq Q(\theta_t | \theta_t) = 0 \quad \therefore \Delta l(\theta_{t+1}) \geq 0$$

$$Q(\theta | \theta_t) = \sum_{n=1}^N \sum_z \underbrace{p(z | x_n, \theta_t)}_{\text{expectation over } z \text{ with current } \theta_t} \underbrace{\log p(x_n, z | \theta)}_{\text{joint likelihood of observed \& hidden}} + \text{const}$$

- General steps of EM:
 - Define likelihood model with parameters θ
 - Identify hidden variables z
 - Derive the auxiliary function and the E and M equations
 - In each iteration, estimate the posteriors of hidden variables
 - Re-estimate the model parameters. Repeat until stop

EE6887-Chang

5-11

Expectation-Maximization (E-M) Solution of GMM

- EM for estimating θ and τ_i .
 - Follow 'divide and conquer' principle. In iteration step t:

$$\text{Expectation: } \tau_n^{i(t)} = \frac{\pi_i^{(t)} N(x_n | \mu_i^{(t)}, \Sigma_i^{(t)})}{\sum_j \pi_j^{(t)} N(x_n | \mu_j^{(t)}, \Sigma_j^{(t)})} \quad \text{Weight from component } i$$

$$\text{Maximization: } \mu_i^{(t+1)} = \frac{\sum_n \tau_n^{i(t)} x_n}{\sum_n \tau_n^{i(t)}} \quad \Sigma_i^{(t+1)} = \frac{\sum_n \tau_n^{i(t)} (x_n - \mu_i^{(t)}) (x_n - \mu_i^{(t)})^T}{\sum_n \tau_n^{i(t)}}$$

Divide data to each group,
Compute mean and variance
from each group

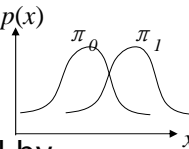
$$\pi_i^{(t+1)} = \frac{\sum_n \tau_n^{i(t)}}{N}$$

EE6887-Chang

5-12

GMM for Clustering

- Given the estimated GMM model, compute posteriors $\tau^i = p(z_i = 1 | x, \theta)$, $\theta = \{\mu_0, \Sigma_0, \mu_1, \Sigma_1\}$
- Estimate the probability that x is generated by cluster i



$$\text{Expectation: } \tau_n^{i(t)} = \frac{\pi_i^{(t)} N(x_n | \mu_i^{(t)}, \Sigma_i^{(t)})}{\sum_j \pi_j^{(t)} N(x_n | \mu_j^{(t)}, \Sigma_j^{(t)})}$$

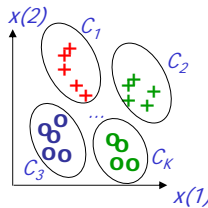
- Each sample is assigned to every cluster with a 'soft' decision.

EE6887-Chang

5-13

Comparison: K-Mean Clustering

- Training data $\{x_i\} + \{label(i) ?\}$
- Unsupervised learning
- K-mean clustering
 - Fix K values
 - Initialize the representative of each cluster
 - Map samples to closest cluster (hard decision)
 - Re-compute the centers



x_1, x_2, \dots, x_N samples

for $i=1, 2, \dots, N$,

$x_i \rightarrow C_k$, if $Dist(x_i, C_k) < Dist(x_i, C_{k'}), k \neq k'$

end

EE6887-Chang

5-14