

语音及语言信息处理国家工程实验室

# **Pattern Classification (V)**











# Outline



- Bayesian Decision Theory
  - How to make the optimal decision?
  - Maximum *a posterior* (MAP) decision rule
- Generative Models
  - Joint distribution of observation and label sequences
  - Model estimation: MLE, Bayesian learning, discriminative training
- Discriminative Models
  - Model the posterior probability directly (discriminant function)
  - Logistic regression, support vector machine, neural network







$$C_{P} = \underset{C_{i}}{\operatorname{arg\,max}} p(C_{i} \mid X) = \underset{C_{i}}{\operatorname{arg\,max}} P(C_{i}) \cdot p(X \mid C_{i})$$
  
$$\approx \underset{C_{i}}{\operatorname{arg\,max}} \overline{P}_{\Gamma_{i}}(C_{i}) \cdot \overline{p}_{\Lambda_{i}}(X \mid C_{i})$$







# Example: Gaussian Mixture Model (I)

- Gaussian distribution (univariate/multivariate) is a single mode distribution.
- In many cases, the true distribution of data is complicate and has multiple modes in nature.



- For this kind of applications, better to use a more flexible model
  - Gaussian mixture model (GMM)
  - A GMM can be tuned to approximate any arbitrary distribution



# $p(x) = \bigotimes_{k=1}^{K} W_k \times N(x \mid \mathcal{M}_k, \mathcal{S}_k^2)$

Univariate density

Gaussian mixture model (GMM)

$$p(\mathbf{x}) = \sum_{k=1}^{K} \omega_k \cdot N(\mathbf{x} \,|\, \mathbf{\mu}_k, \boldsymbol{\Sigma}_k)$$

- GMM is a mixture of single Gaussian distribution (each one is called mixture component) which have different means and variances.
- $\omega_k$  is called mixture weight, prior probability of each mixture component.

$$\sum_{k=1}^{n} \omega_k = 1$$

• GMM is widely used for speaker recognition, audio classification, audio segmentation, etc.





# Example: Gaussian Mixture Model (III)



- However, estimation of a GMM is not trivial.
- Consider a simple case:
  - We have a set of training data  $D = \{x_1, x_2, \dots, x_n\}$
  - Use a 2-mixture GMM to model it:

$$p(x) = \frac{0.3}{\sqrt{2\rho S_1^2}} e^{-\frac{(x-m_1)^2}{2S_1^2}} + \frac{0.7}{\sqrt{2\rho S_2^2}} e^{-\frac{(x-m_2)^2}{2S_2^2}}$$

- We try to get ML estimate of  $\mu_1$ ,  $\sigma_1$ ,  $\mu_2$ ,  $\sigma_2$  from training data.
- Simple maximization based on differential calculus does not work.
  - For each  $x_i$ , we do not know which mixture it comes from. The number of item in likelihood function  $p(D|\mu_1, \sigma_1, \mu_2, \sigma_2)$  increases exponentially as we observe more and more data.
  - No simple solution.
- Need alternative method Expectation Maximization (EM) algorithm

#### The Expectation-Maximization Algorithm (I)

- EM algorithm is an iterative method of obtaining maximum likelihood estimate of model parameters.
- EM suits best to the so-called *missing data* problem:
  - Only observe a subset of features, called *observed*, X.
  - Other features are *missing* or unobserved, denoted as Y.
  - The complete data Z={X,Y}.
  - If given the complete data Z, it is usually easy to obtain ML estimation of model parameters.
  - How to do ML estimation based on observed X only??



### The Expectation-Maximization Algorithm (II)

- Initialization: find an initial values for unknown parameters
- EM algorithm consists of two steps:
  - Expectation (E-step): the expectation is calculated with respect of the missing data Y, using the current estimate of the unknown parameters and conditioned upon the observed X.
  - Maximization (M-step): provides a new estimate of unknown parameters (better than the initial ones) in terms of maximizing the above expectation → increasing likelihood function of observed.
- Iterate until convergence



#### EM Algorithm: E-step



• E-step: form an auxiliary function

$$Q(\theta; \theta^{(i)}) = E_{Y} \Big[ \ln p(X, Y | \theta) | X, \theta^{(i)} \Big]$$

- The expectation of log-likelihood function of complete data is calculated based on the current estimate of unknown parameter, and conditioned on the observed data.
- $Q(q;q^{(i)})$  is a function of  $\theta$  with  $q^{(i)}$  assumed to be fixed.
- If missing data Y is continuous:

$$Q(q;q^{(i)}) = \bigcup_{L_Y} \ln p(X,Y|q) \times p(Y|X,q^{(i)}) \,\mathrm{d}Y$$

– If missing data Y is discrete:

$$Q(q;q^{(i)}) = \mathop{\text{a}}_{Y} \ln p(X,Y|q) \times p(Y|X,q^{(i)})$$



#### EM Algorithm: M-step

• M-step: choose a new estimate  $q^{(i+1)}$  which maximizes  $Q(q;q^{(i)})$ 

$$Q^{(i+1)} = \arg\max_{q} Q(q; q^{(i)})$$

-  $q^{(i+1)}$  is a better estimate in terms of increasing likelihood value  $p(X | \theta)$  than  $q^{(i)}$ 

$$p(X \mid \theta^{(i+1)}) \ge p(X \mid \theta^{(i)})$$

• Replace  $q^{(i+1)}$  with  $q^{(i)}$  and iterate until convergence.









- EM algorithm guarantees that the log-likelihood of the observed data p(X|θ) will increase monotonically.
- EM algorithm may converge to a local maximum or global maximum. And convergence rate is reasonably good.
- Applications of the EM algorithm:
  - ML estimation of some complicated models,
    - e.g., GMM, HMM, ... (in general mixture models of e-family)
  - ET (emission tomography) image reconstruction
  - Active Noise Cancellation (ANC)
  - Spread-spectrum multi-user communication



#### An Application of EM Algorithm: ML Estimation of Multivariate GMM(I)



• We decide to model the data by using multivariate GMM:

$$p(X) = \mathop{\text{a}}\limits_{k=1}^{K} \mathcal{W}_{k} \times N(X \mid \mathcal{M}_{k}, \overset{\text{a}}{a}_{k}) \qquad (\text{with } \mathop{\text{a}}\limits_{k=1}^{K} \mathcal{W}_{k} = 1)$$

- Problem: use data set D to estimate GMM model parameters, including ωk, μk, Σk (k=1,2,...,K).
- If we know the label of mixture component label It from which each data Xt come from, the estimation is easy.
- Since the mixture component is not available in training set, we treat it as missing data:
  - Observed data:  $D = \{X_1, X_2, \dots, X_T\}$ .
  - Missing data: L={I1, I2, ..., IT}.
  - Complete data: {D,L} = {X1, I1, X2, I2, ..., XT, IT}

#### An Application of EM algorithm: ML estimation of Multivariate GMM(II)



• E-step:

$$Q(\{\omega_{k},\mu_{k},\sum_{k}\}|\{\omega_{k}^{(i)},\mu_{k}^{(i)},\sum_{k}^{(i)}\}) = \sum_{L} \ln p(D,L|\{\omega_{k},\mu_{k},\sum_{k}\}) \cdot p(L|D,\{\omega_{k}^{(i)},\mu_{k}^{(i)},\sum_{k}^{(i)}\})$$

$$= C + \sum_{L} \left[\sum_{t=1}^{T} [\ln \omega_{l_{t}} - \frac{1}{2}\ln|\sum_{l_{t}}| - \frac{1}{2} \cdot (X_{t} - \mu_{l_{t}})^{\mathrm{T}} \sum_{l_{t}}^{-1} (X_{t} - \mu_{l_{t}})] \cdot \prod_{t=1}^{T} p(l_{t}|X_{t},\{\omega_{k}^{(i)},\mu_{k}^{(i)},\sum_{k}^{(i)}\})\right]$$

$$= C + \sum_{k=1}^{K} \sum_{t=1}^{T} [\ln \omega_{k} - \frac{1}{2}\ln|\sum_{k}| - \frac{1}{2} \cdot (X_{t} - \mu_{k})^{\mathrm{T}} \sum_{k}^{-1} (X_{t} - \mu_{k})] \cdot p(l_{t} = k \mid X_{t}, \omega_{k}^{(i)}, \mu_{k}^{(i)}, \sum_{k}^{(i)})$$



#### An Application of EM Algorithm: ML Estimation of Multivariate GMM(III)



• M-step:

$$\begin{aligned} \frac{\partial Q}{\partial \mu_{k}} &= 0 \Rightarrow \mu_{k}^{(i+1)} = \frac{\sum_{t=1}^{T} X_{t} \cdot p(l_{t} = k \mid X_{t}, \omega_{k}^{(i)}, \mu_{k}^{(i)}, \Sigma_{k}^{(i)})}{\sum_{t=1}^{T} p(l_{t} = k \mid X_{t}, \omega_{k}^{(i)}, \mu_{k}^{(i)}, \Sigma_{k}^{(i)})} \\ \frac{\partial Q}{\partial \Sigma_{k}} &= 0 \Rightarrow \Sigma_{k}^{(i+1)} = \frac{\sum_{t=1}^{T} (X_{t} - \mu_{k}^{(i+1)})^{T} \cdot (X_{t} - \mu_{k}^{(i+1)}) \cdot p(l_{t} = k \mid X_{t}, \omega_{k}^{(i)}, \mu_{k}^{(i)}, \Sigma_{k}^{(i)})}{\sum_{t=1}^{T} p(l_{t} = k \mid X_{t}, \omega_{k}^{(i)}, \mu_{k}^{(i)}, \Sigma_{k}^{(i)})} \\ \frac{\partial}{\partial \omega_{k}} [Q - \lambda(\sum_{k=1}^{K} \omega_{k} - 1)] = 0 \Rightarrow \omega_{k} = \frac{\sum_{t=1}^{T} p(l_{t} = k \mid X_{t}, \omega_{k}^{(i)}, \mu_{k}^{(i)}, \Sigma_{k}^{(i)})}{\sum_{k=1}^{K} \sum_{t=1}^{T} p(l_{t} = k \mid X_{t}, \omega_{k}^{(i)}, \mu_{k}^{(i)}, \Sigma_{k}^{(i)})} = \frac{\sum_{t=1}^{T} p(l_{t} = k \mid X_{t}, \omega_{k}^{(i)}, \mu_{k}^{(i)}, \Sigma_{k}^{(i)})}{T} \end{aligned}$$



#### An Application of EM Algorithm: ML Estimation of Multivariate GMM(IV)

• where

- Iterative ML estimation of GMM
  - Initiation: choose  $\{W_k^{(0)}, m_k^{(0)}, a_k^{(0)}\}$ . Usually use vector clustering algorithm (such as K-means) to cluster all data into K clusters. Each cluster is used to train for one Gaussian mix.
  - i=0;
  - Use EM algorithm to refine model estimation

 $\{W_k^{(i)}, M_k^{(i)}, \mathring{a}_k^{(i)}\} \triangleright \{W_k^{(i+1)}, M_k^{(i+1)}, \mathring{a}_k^{(i+1)}\}$ 

i++, go back until convergence.



#### **GMM Initialization: K-Means Clustering**



- K-Means Clustering: a.k.a. unsupervised learning
- Cluster a data set into many homogeneous groups
- K-Means algorithm:
  - step 1: assign all data into one group; calculate centroid.
  - step 2: choose a group and split.
  - step 3: re-assign all data to groups.
  - step 4: calculate centroids for all groups.
  - step 5: go back to step 3 until convergence.
  - step 6: stop until K classes
- Basics for clustering:
  - distance measure
  - centroid calculation
  - choose a group and split





#### Applications of GMM

- GMM is widely used to model speech or audio signals. In many cases, we use diagonal covariance matrices in GMM for simplicity.
- Speaker recognition:
  - Collect some speech signals from all known speakers.
  - Train a GMM for each known speaker by using his/her voice.
  - For an unknown speaker, prompt him/her to say sth.
  - Classify it based on all trained GMM' s and determine the speaker' s identity or reject.
- Audio classification:
  - Classify a continuous audio/video stream (from radio or TV) into some homogeneous segments: anchor' s speech, in-field interview, telephone interview, music, commercial ads, sports, etc.
  - For each category, train a GMM based on training data.
  - Use all trained GMMs to scan an unknown audio stream to segment it.



#### Project: Building a 2-Class Classifier

- Given some data from two classes
- Build a classifier with multivariate Gaussian models
  - ML estimation
  - Test with the plug-in MAP decision rule
- Improve it with GMM models
  - Initialize GMM with the K-means clustering
  - Estimate GMM with the EM algorithm
  - Investigate GMM with the mixture number = 2, 4, 8.
- Improve the Gaussian classifier with discriminative training (minimum classification error estimation)
- Preferably programming with C/C++
- Report all of your experiments and your best classifier.

