

语音及语言信息处理国家工程实验室

Pattern Classification (V)

杜俊

jundu@ustc.edu.cn



中国科学技术大学
安徽科大讯飞信息科技
股份有限公司



Outline

- Bayesian Decision Theory
 - How to make the optimal decision?
 - Maximum *a posteriori* (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

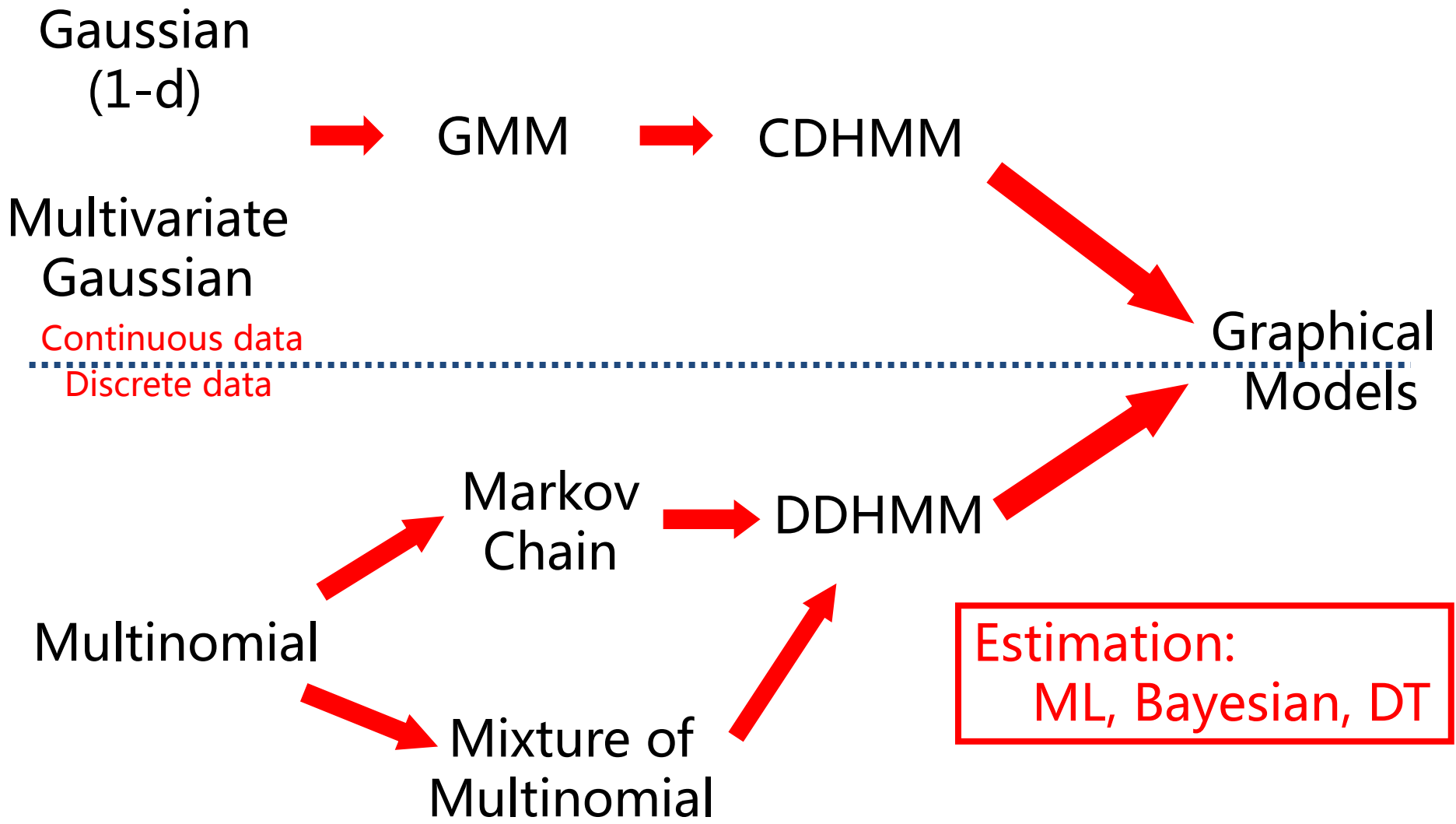


Plug-in MAP Decision Rule

$$\begin{aligned} C_p &= \arg \max_{C_i} p(C_i | X) = \arg \max_{C_i} P(C_i) \cdot p(X | C_i) \\ &\approx \arg \max_{C_i} \bar{P}_{\Gamma_i}(C_i) \cdot \bar{p}_{\Lambda_i}(X | C_i) \end{aligned}$$

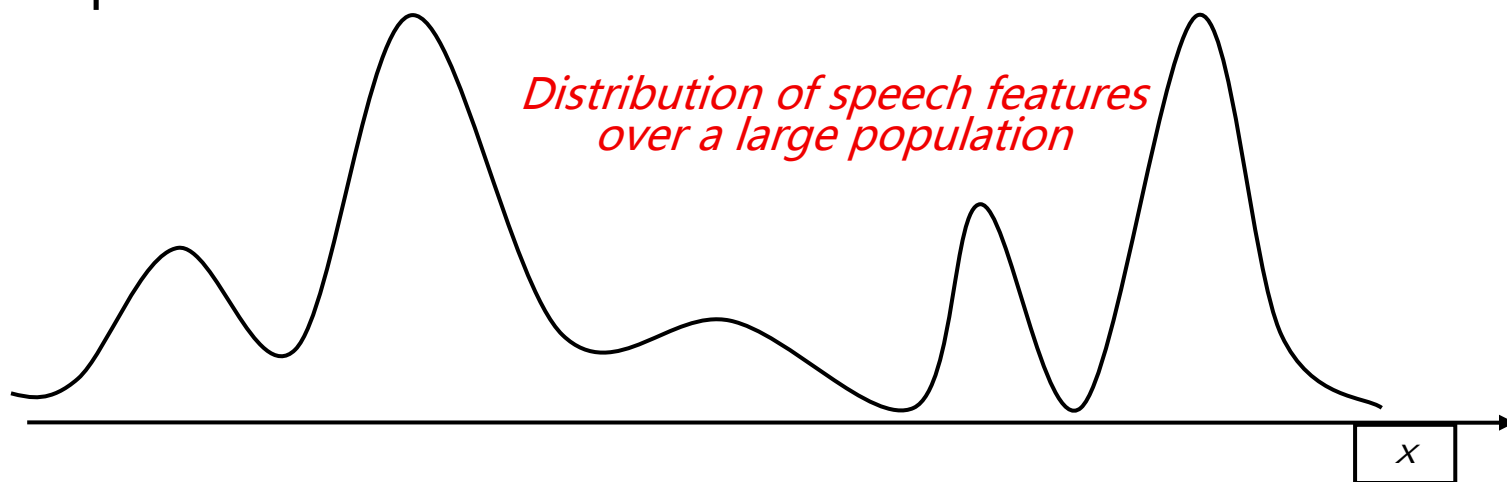


Statistical Models: Roadmap



Example: Gaussian Mixture Model (I)

- Gaussian distribution (univariate/multivariate) is a single mode distribution.
- In many cases, the true distribution of data is complicated 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



Example: Gaussian Mixture Model (II)

- Gaussian mixture model (GMM)

- Univariate density

$$p(x) = \sum_{k=1}^K \omega_k \cdot N(x | \mu_k, \sigma_k^2)$$

- Multivariate density

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

- GMM is a mixture of single Gaussian distribution (each one is called mixture component) which have different means and variances.

- ω_k is called mixture weight, prior probability of each mixture component.

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

- 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:

$$p(x) = \frac{0.3}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} + \frac{0.7}{\sqrt{2\pi\sigma_2^2}} e^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}}$$

- 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 \rightarrow increasing likelihood function of observed.
- Iterate until convergence



EM Algorithm: E-step

- E-step: form an auxiliary function

$$Q(\theta; \theta^{(i)}) = E_Y \left[\ln p(X, Y | \theta) \mid X, \theta^{(i)} \right]$$

- 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.
- $Q(\theta; \theta^{(i)})$ is a function of θ with $\theta^{(i)}$ assumed to be fixed.
- If missing data Y is continuous:

$$Q(\theta; \theta^{(i)}) = \int_{\Lambda_Y} \ln p(X, Y | \theta) \cdot p(Y | X, \theta^{(i)}) dY$$

- If missing data Y is discrete:

$$Q(\theta; \theta^{(i)}) = \sum_Y \ln p(X, Y | \theta) \cdot p(Y | X, \theta^{(i)})$$



EM Algorithm: M-step

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

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

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

$$p(X | \theta^{(i+1)}) \geq p(X | \theta^{(i)})$$

- Replace $\theta^{(i+1)}$ with $\theta^{(i)}$ and iterate until convergence.



EM Algorithm

- EM algorithm guarantees that the log-likelihood of the observed data $p(X|\theta)$ 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)

- Assume we observe a data set $D = \{X_1, X_2, \dots, X_T\}$ (a set of vectors)
- We decide to model the data by using multivariate GMM:

$$p(X) = \sum_{k=1}^K \omega_k \cdot N(X | \mu_k, \Sigma_k) \quad \left(\text{with } \sum_{k=1}^K \omega_k = 1\right)$$

- Problem: use data set D to estimate GMM model parameters, including $\omega_k, \mu_k, \Sigma_k$ ($k=1, 2, \dots, K$).
- If we know the label of mixture component label l_t from which each data X_t 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: $L = \{l_1, l_2, \dots, l_T\}$.
 - Complete data: $\{D, L\} = \{X_1, l_1, X_2, l_2, \dots, X_T, l_T\}$



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

- E-step:

$$\begin{aligned} Q(\{\omega_k, \mu_k, \Sigma_k\} | \{\omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)}\}) &= \sum_L \ln p(D, L | \{\omega_k, \mu_k, \Sigma_k\}) \cdot p(L | D, \{\omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)}\}) \\ &= C + \sum_L \left[\sum_{t=1}^T [\ln \omega_{l_t} - \frac{1}{2} \ln |\Sigma_{l_t}| - \frac{1}{2} \cdot (X_t - \mu_{l_t})^T \Sigma_{l_t}^{-1} (X_t - \mu_{l_t})] \cdot \prod_{t=1}^T p(l_t | X_t, \{\omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)}\}) \right] \\ &= C + \sum_{k=1}^K \sum_{t=1}^T [\ln \omega_k - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} \cdot (X_t - \mu_k)^T \Sigma_k^{-1} (X_t - \mu_k)] \cdot p(l_t = k | X_t, \omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)}) \end{aligned}$$

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

- M-step:

$$\frac{\partial Q}{\partial \mu_k} = 0 \Rightarrow \mu_k^{(i+1)} = \frac{\sum_{t=1}^T X_t \cdot p(l_t = k | X_t, \omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)})}{\sum_{t=1}^T p(l_t = k | X_t, \omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)})}$$

$$\frac{\partial Q}{\partial \Sigma_k} = 0 \Rightarrow \Sigma_k^{(i+1)} = \frac{\sum_{t=1}^T (X_t - \mu_k^{(i+1)})^T \cdot (X_t - \mu_k^{(i+1)}) \cdot p(l_t = k | X_t, \omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)})}{\sum_{t=1}^T p(l_t = k | X_t, \omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)})}$$

$$\frac{\partial}{\partial \omega_k} [Q - \lambda (\sum_{k=1}^K \omega_k - 1)] = 0 \Rightarrow \omega_k = \frac{\sum_{t=1}^T p(l_t = k | X_t, \omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)})}{\sum_{k=1}^K \sum_{t=1}^T p(l_t = k | X_t, \omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)})} = \frac{\sum_{t=1}^T p(l_t = k | X_t, \omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)})}{T}$$

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

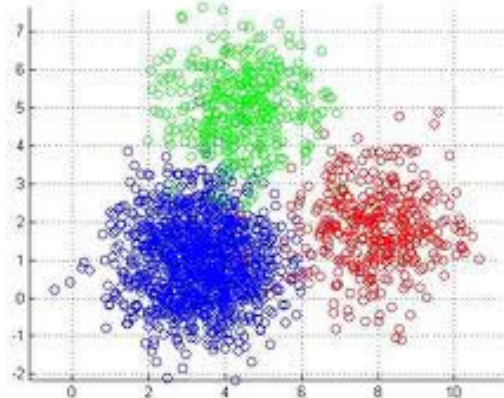
- where

$$p(l_t = k | X_t, \{\omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)}\}) = \frac{\omega_k^{(i)} \cdot N(X_t | \mu_k^{(i)}, \Sigma_k^{(i)})}{\sum_{k=1}^K \omega_k^{(i)} \cdot N(X_t | \mu_k^{(i)}, \Sigma_k^{(i)})}$$

- Iterative ML estimation of GMM
 - Initiation: choose $\{\omega_k^{(0)}, \mu_k^{(0)}, \Sigma_k^{(0)}\}$. 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
$$\{\omega_k^{(i)}, \mu_k^{(i)}, \Sigma_k^{(i)}\} \Rightarrow \{\omega_k^{(i+1)}, \mu_k^{(i+1)}, \Sigma_k^{(i+1)}\}$$
 - $i++$, 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.



Next Model: Hidden Markov Model

- HMM has been applied in many areas to model different types of data, such as speech, language, DNA, protein, etc.
- HMM is an extension of Markov chain model.
- Also needs EM algorithm in ML estimation.
- We will talk about HMM in detail as a future topic.



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.

