

语音及语言信息处理国家工程实验室

# **Pattern Classification (III)**







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 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







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





## **Generative VS. Discriminative Models**



- Generative Models
  - Joint distribution of observation and label sequences

 $p(X, \mathbf{C}_i) = P(\mathbf{C}_i) \cdot p(X | \mathbf{C}_i)$ 

- Discriminative Models
  - Model the posterior probability directly (discriminant function)
  - Model the boundaries of data sets

 $g_i(X) \longrightarrow p(C_i | X)$ 





# Useful Models (I)

- Choose a proper model based on the nature of observation data
- Some useful statistical models for a variety of data types:
  - Gaussian (Normal) distribution
    - ➔ Uni-modal continuous feature scalars
  - Multivariate Gaussian (Normal) distribution
    - ➔ Uni-modal continuous feature vectors
  - Gaussian Mixture models (GMM)
    - ➔ Multi-modal continuous feature scalars/vectors







# Useful Models (II)



- Markov chain model: discrete sequential data
  - N-gram model in language modeling
- Hidden Markov model (HMM): ideal for various kinds of sequential observation data; provides better modeling capability than simple Markov chain model :
  - Model speech signals for recognition (speech recognition)
  - Model sign/gesture for recognition (sign language recognition)
  - Model biological data (DNA & protein sequence): profile HMM

- ...



### Useful Models (III)



- Markov random field: multi-dimensional spatial data
  - Model image data: e.g., used for OCR, etc
  - HMM is a special case of Markov random field
- Graphical models (a.k.a., Bayesian networks, Belief networks)
  - Widely used in machine learning, data mining
  - High-dimensional data (discrete or continuous)
  - Model a very complex stochastic process
  - Automatically learn dependency from data
  - HMM is also a special case of graphical model
- Discriminant functions
  - Linear regression
  - Logistic regression
  - Support vector machine
  - Neural network



# **Discriminant Functions (I)**



- Instead of designing a classifier based on probability distribution, we can build an ad-hoc classifier based on some discriminant functions to model class boundary info directly.
  - A set of discriminant functions  $g_i(\mathbf{x}; \theta_i)$  for each class Ci
  - $g_i(\mathbf{x}; \theta_i)$  has a pre-defined function form with unknown parameters  $\theta_i$  which should be estimated from training data
  - For an unknown pattern with feature vector **x**, the decision is

$$C_{\mathbf{x}} = \arg\max_{i} g_{i}(\mathbf{x}; \theta_{i})$$



# **Discriminant Functions (II)**



- Examples
  - Linear discriminant function:

$$g(\mathbf{x}) = \mathbf{w}^{\mathrm{T}}\mathbf{x} + w_0 = \langle \mathbf{w}, \mathbf{x} \rangle + w_0$$

– Quadratic discrimiant function: (2<sup>nd</sup> order)

$$g(\mathbf{x}) = \mathbf{x}^{\mathrm{T}} \mathbf{W} \mathbf{x} + \mathbf{w}^{\mathrm{T}} \mathbf{x} + w_0 = \langle \mathbf{x} \otimes \mathbf{x}, \operatorname{vec}(\mathbf{W}) \rangle + \langle \mathbf{w}, \mathbf{x} \rangle + w_0$$

- Polynomial discriminant function: (N-th order)
- Neural network: (arbitrary nonlinear function)
- Optimal MAP classifier is a special case when choosing discriminant functions as the class posterior probabilities



## **Discriminant Functions (III)**



- Unknown parameters of discriminant functions are estimated based on optimizing an objective function via some gradient descent methods:
  - Linear regression: Achieving a good mapping
  - Logistic regression: Minimizing empirical classification errors
  - Support vector machine (SVM): Maximizing separation margin
  - Neural network: MMSE or cross-entropy minimization



#### **Linear Regression**



• Find a good mapping from **x** to y









# Logistic Regression



NEL-SL

• Counting errors in training samples

$$(\mathbf{x}_{i}, y_{i}) \Rightarrow \begin{cases} g_{i} = -y_{i} \mathbf{w}^{\mathrm{T}} \mathbf{x}_{i} < 0 & \text{Correct Classification} \\ g_{i} = -y_{i} \mathbf{w}^{\mathrm{T}} \mathbf{x}_{i} > 0 & \text{Wrong Classification} \end{cases}$$
$$\mathbf{w}^{*} = \arg\min \sum H(g_{i}) = \arg\min \sum H(-y_{i} \mathbf{w}^{\mathrm{T}} \mathbf{x}_{i}) \qquad 0 \end{cases}$$

$$\mathbf{w}^{\prime} = \underset{\mathbf{w}}{\operatorname{arg\,min}} \sum_{i} H(g_{i}) = \underset{\mathbf{w}}{\operatorname{arg\,min}} \sum_{i} H(-y_{i} \mathbf{w}^{\prime} \mathbf{x}_{i})$$

$$\overset{0.8}{\underset{0.6}{\overset{0.6}{\phantom{0}}}}$$

$$\mathbf{w}^{*} = \arg\min_{\mathbf{w}} \sum_{i} l(g_{i}) = \arg\min_{\mathbf{w}} \sum_{i} l(-y_{i} \mathbf{w}^{T} \mathbf{x}_{i})$$



logistic sigmoid function







adding per high app

## Support Vector Machine (II)

- The decision boundary H should be as far away from the data of both classes as possible
- We should maximize the margin:  $m = \frac{1}{11}$



# Support Vector Machine (III)

The decision boundary can be found by solving the following constrained optimization problem:

Minimize 
$$rac{1}{2}||\mathbf{w}||^2 = \mathbf{w}^{ ext{T}}\mathbf{w}$$
 subject to  $y_i(\mathbf{w}^T\mathbf{x}_i+b) \geq 1$   $orall i$ 

• Convert to its dual problem (using KKT condition):

max. 
$$W(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1,j=1}^{n} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$
 subject to  $\alpha_i \ge 0, \sum_{i=1}^{n} \alpha_i y_i = 0$ 

• A standard quadratic programming (QP) problem

$$\mathbf{w}^* = \sum_i \alpha_i y_i \mathbf{x}_i$$

• Support vectors:  $\alpha_i \neq 0$ 



# Support Vector Machine (IV)

• We allow "error"  $\xi_i$  in classification  $\rightarrow$  soft-margin SVM





#### Support Vector Machine (V)

• Soft-margin SVM can be formulated as:

$$\mathbf{w}^* = \min_{\mathbf{w}, \mathbf{X}_i} \quad \left[ \frac{1}{2} \| \mathbf{w} \|^2 + C \cdot \sum_i \mathbf{X}_i \right]$$

subject to

$$y_i(x_iw^T + b) > 1 - X_i \quad X_i > 0 \quad ("i)$$

- C is the regularization coefficient
- $\sum_{i} \xi_{i}$  is the upper bound of training classification errors





## Support Vector Machine (VI)



- For nonlinear separation boundary
  - use a Kernel function





## **Neural Network**

- Feed-forward multilayer perceptron (MLP)
- Error back-propagation (BP)



