

# Pattern Classification

All materials in these slides were taken from Pattern Classification (2nd ed) by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000 with the permission of the authors and the publisher

# Chapter 3: Maximum-Likelihood and Bayesian Parameter Estimation (part 2)

- Bayesian Estimation (BE)
- Bayesian Parameter Estimation: Gaussian Case
- Bayesian Parameter Estimation: General Estimation
- Problems of Dimensionality
- Computational Complexity
- Component Analysis and Discriminants
- Hidden Markov Models

## Bayesian Estimation (Bayesian learning to pattern classification problems)

• In MLE  $\theta$  was supposed fix

• In BE  $\theta$  is a random variable

- The computation of posterior probabilities
   P(ω<sub>i</sub> | x) lies at the heart of Bayesian
   classification
- Goal: compute  $P(\omega_i | x, D)$ Given the sample *D*, Bayes formula can be written  $p(\mathbf{x} | \omega_i, D) P(\omega_i | D)$

$$P(\omega_i \mid \mathbf{x}, D) = \frac{p(\mathbf{x} \mid \omega_i, D) P(\omega_i \mid D)}{\sum_{j=1}^{c} p(\mathbf{x} \mid \omega_j, D) P(\omega_j \mid D)}$$

Pattern Classification, Chapter 3

### • To demonstrate the preceding equation, use:

 $P(\mathbf{x}, D | \omega_i) = P(\mathbf{x} | D, \omega_i) P(D | \omega_i) \text{ (from def. of cond. prob.)}$   $P(\mathbf{x} | D) = \sum_j P(\mathbf{x}, \omega_j | D) \quad \text{(from law of total prob.)}$   $P(\omega_i) = P(\omega_i | D) \text{ (Training sample provides this!)}$ We assume that samples in different classes are independent Thus :

$$P(\omega_i \mid \mathbf{x}, D) = \frac{p(\mathbf{x} \mid \omega_i, D_i) P(\omega_i)}{\sum_{j=1}^{c} p(\mathbf{x} \mid \omega_j, D) P(\omega_j)}$$

Pattern Classification, Chapter 3

## Bayesian Parameter Estimation: Gaussian Case

**Goal:** Estimate  $\theta$  using the a-posteriori density  $P(\theta \mid D)$ 

• The univariate case:  $P(\mu \mid D)$  $\mu$  is the only unknown parameter

> P(x | μ) ~ N(μ,σ²) P(μ) ~ N(μ<sub>0</sub>,σ<sub>0</sub>²)

( $\mu_0$  and  $\sigma_0$  are known!)

$$P(\mu \mid \mathsf{D}) = \frac{P(\mathsf{D} \mid \mu)P(\mu)}{\int P(\mathsf{D} \mid \mu)P(\mu)d\mu}$$
$$= \alpha \prod_{k=1}^{k=n} P(x_k \mid \mu)P(\mu)$$

• But we know

$$p(x_k \mid \mu) \sim N(\mu, \sigma^2) \text{ and } p(\mu) \sim N(\mu_0, \sigma_0^2)$$

Plugging in their gaussian expressions and extracting out factors not depending on  $\mu$  yields:

$$p(\mu \mid D) = \alpha \exp\left(-\frac{1}{2}\left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\mu^2 - 2\left(\frac{1}{\sigma^2}\sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2}\right)\mu\right]\right]$$

(from eq. 29 page 93)

Pattern Classification, Chapter 3

**Bayes** Formula

### Observation: $p(\mu|D)$ is an exponential of a quadratic

$$p(\mu \mid D) = \alpha \exp\left(-\frac{1}{2}\left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\mu^2 - 2\left(\frac{1}{\sigma^2}\sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2}\right)\mu\right]\right)$$

It is again normal! It is called a reproducing density

$$P(\mu \mid \mathsf{D}) \sim N(\mu_n, \sigma_n^2)$$

$$p(\mu \mid D) = \alpha \exp\left(-\frac{1}{2}\left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\mu^2 - 2\left(\frac{1}{\sigma^2}\sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2}\right)\mu\right]\right)$$

 Identifying coefficients in the top equation with that of the generic Gaussian

$$p(\mu \mid D) = \frac{1}{\sqrt{2\pi\sigma_n}} \exp\left(-\frac{1}{2}\left[\frac{\mu - \mu_n}{\sigma_n}\right]\right)^2$$

Yields expressions for  $\mu_n$  and  $\sigma_n{}^2$ 

$$\frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \text{ and } \frac{\mu_n}{\sigma_n^2} = \frac{n}{\sigma^2}\hat{\mu}_n + \frac{\mu_0}{\sigma_0^2}$$

Pattern Classification, Chapter 3

7

### Solving for $\mu_n$ and $\sigma_n^2$ yields:

$$\mu_n = \left(\frac{n\sigma_0^2}{n_0\sigma_0^2 + \sigma^2}\right)\hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2}\mu_0$$
  
and 
$$\sigma_n^2 = \frac{\sigma_0^2\sigma^2}{n\sigma_0^2 + \sigma^2}$$

From these equations we see as *n* increases:

- the variance decreases monotonically
- the estimate of  $p(\mu|D)$  becomes more peaked



FIGURE 3.2. Bayesian learning of the mean of normal distributions in one and two dimensions. The posterior distribution estimates are labeled by the number of training samples used in the estimation. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

9

### The univariate case P(x | D)

- $P(\mu \mid D)$  computed (in preceding discussion)
- P(x | D) remains to be computed!

 $P(x | \mathbf{D}) = \int P(x | \mu) P(\mu | \mathbf{D}) d\mu$  is Gaussian

It provides: 
$$P(x | \mathbf{D}) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$$

We know  $\sigma^2$  and how to compute  $\mu_n$  and  $\sigma_n^2$ (Desired class-conditional density P(x | D<sub>j</sub>,  $\omega_j$ )) Therefore: using P(x | D<sub>j</sub>,  $\omega_j$ ) together with P( $\omega_j$ ) And using Bayes formula, we obtain the Bayesian classification rule:

$$M_{\substack{\omega_{j} \\ \omega_{j}}} \left[ P(\omega_{j} \mid x, \mathsf{D}] \equiv M_{\substack{\omega_{j} \\ \omega_{j}}} \left[ P(x \mid \omega_{j}, \mathsf{D}_{j}) P(\omega_{j}) \right] \right]$$

# 3.5 Bayesian Parameter Estimation: General Theory

- P(x | D) computation can be applied to any situation in which the unknown density can be parameterized: the basic assumptions are:
  - The form of P(x |  $\theta$ ) is assumed known, but the value of  $\theta$  is not known exactly
  - Our knowledge about θ is assumed to be contained in a known prior density P(θ)
  - The rest of our knowledge θ is contained in a set D of n random variables x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub> that follows P(x)

The basic problem is: "Compute the posterior density P( $\theta \mid D$ )" then "Derive P(x | D)", where  $p(\mathbf{x} \mid D) = \int p(\mathbf{x} \mid \theta) p(\theta \mid D) d\theta$ 

Using Bayes formula, we have:

$$P(\boldsymbol{\theta} \mid \mathsf{D}) = \frac{P(\mathsf{D} \mid \boldsymbol{\theta})P(\boldsymbol{\theta})}{\int P(\mathsf{D} \mid \boldsymbol{\theta})P(\boldsymbol{\theta})d\boldsymbol{\theta}}$$

And by the independence assumption:

$$p(\mathsf{D} | \mathbf{\theta}) = \prod_{k=1}^{n} p(\mathbf{x}_k | \mathbf{\theta})$$

Pattern Classification, Chapter 3

## Convergence (from notes)

- Problems of Dimensionality
  - Problems involving 50 or 100 features are common (usually binary valued)
  - Note: microarray data might entail ~20000 real-valued features
  - Classification accuracy dependant on
    - dimensionality
    - amount of training data
    - discrete vs continuous

 Case of two class multivariate normal with the same covariance

- $P(\mathbf{x}|\omega_j) \sim N(\mu_j, \Sigma), j=1,2$
- Statistically independent features
- If the priors are equal then:

$$P(error) = \frac{1}{\sqrt{2\pi}} \int_{r/2}^{\infty} e^{-u^{2}/2} du \qquad \text{(Bayes error)}$$
where :  $r^{2} = (\mu_{1} - \mu_{2})^{t} \Sigma^{-1} (\mu_{1} - \mu_{2})$ 
 $r^{2}$  is the squared Mahalanobi s distance
$$\lim_{r \to \infty} P(error) = 0$$

#### • If features are *conditionally* independent then:

$$\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2, ..., \sigma_d^2)$$
$$r^2 = \sum_{i=1}^{i=d} \left(\frac{\mu_{i1} - \mu_{i2}}{\sigma_i}\right)^2$$

- Do we remember what conditional independence is?
- Example for binary features:

Let  $p_i = \Pr[x_i=1|\omega_1]$  then  $P(x|\omega_1)$  is the product of the  $p_i$ 

- Most useful features are the ones for which the difference between the means is large relative to the standard deviation
  - Doesn't require independence
- Adding independent features helps increase  $r \rightarrow$  reduce error
- Caution: adding features increases cost & complexity of feature extractor and classifier
- It has frequently been observed in practice that, beyond a certain point, the inclusion of additional features leads to worse rather than better performance:
  - we have the wrong model !
  - we don't have enough training data to support the additional dimensions



**FIGURE 3.3.** Two three-dimensional distributions have nonoverlapping densities, and thus in three dimensions the Bayes error vanishes. When projected to a subspace—here, the two-dimensional  $x_1 - x_2$  subspace or a one-dimensional  $x_1$  subspace—there can be greater overlap of the projected distributions, and hence greater Bayes error. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Computational Complexity

 Our design methodology is affected by the computational difficulty

"big oh" notation
f(x) = O(h(x)) "big oh of h(x)"
If: ∃(c<sub>0</sub>,x<sub>0</sub>) ∈ ℜ<sup>2</sup>; |f(x)| ≤ c<sub>0</sub>|h(x)|
(An upper bound on f(x) grows no worse than h(x) for sufficiently large x!)

 $f(x) = 2+3x+4x^2$  $g(x) = x^2$  $f(x) = O(x^2)$ 

• "big oh" is not unique!  $f(x) = O(x^2); f(x) = O(x^3); f(x) = O(x^4)$ 

• "big theta" notation  $f(x) = \theta(h(x))$ If:  $\exists (x_0, c_1, c_2) \in \Re^3; \forall x > x_0$   $0 \le c_1 g(x) \le f(x) \le c_2 g(x)$ 

 $f(x) = \theta(x^2)$  but  $f(x) \neq \theta(x^3)$ 

### Complexity of the ML Estimation

- Gaussian priors in *d* dimensions classifier with *n* training samples for each of *c* classes
- For each category, we have to compute the discriminant function

$$g(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \hat{\vec{\mu}})^t \hat{\Sigma}^{-1} (\mathbf{x} - \hat{\mu}) - \frac{\partial}{\partial t} \ln 2\pi - \frac{1}{2} \ln |\hat{\Sigma}| + \ln P(\omega)$$

Total =  $O(d^2 n)$ Total for *c* classes =  $O(cd^2 n) \cong O(d^2 n)$ 

• Cost increase when *d* and *n* are large!

Pattern Classification, Chapter 3

# Overfitting

Dimensionality of model vs size of training data
Issue: not enough data to support the model
Possible solutions:

Reduce model dimensionality

• Make (possibly incorrect) assumptions to better estimate  $\Sigma$ 

# Overfitting

Estimate better Σ
use data pooled from all classes

normalization issues

use pseudo-Bayesian form λΣ<sub>0</sub> + (1-λ)Σ<sub>n</sub>
"doctor" Σ by thresholding entries

reduces chance correlations

assume statistical independence

zero all off-diagonal elements

# Shrinkage

• Shrinkage: weighted combination of common and individual covariances

$$\Sigma_{i}(\alpha) = \frac{(1-\alpha)n_{i}\Sigma_{i} + \alpha n\Sigma}{(1-\alpha)n_{i} + \alpha n} \quad \text{for } 0 < \alpha < 1$$

• We can also shrink the estimate common covariances toward the identity matrix

$$\Sigma(\beta) = (1 - \beta)\Sigma + \beta \mathbf{I}$$
 for  $0 < \beta < 1$ 

## • Component Analysis and Discriminants

- Combine features in order to reduce the dimension of the feature space
- Linear combinations are simple to compute and tractable
- Project high dimensional data onto a lower dimensional space
- Two classical approaches for finding "optimal" linear transformation
  - PCA (Principal Component Analysis) "Projection that best represents the data in a least- square sense"
  - MDA (Multiple Discriminant Analysis) "Projection that best separates the data in a least-squares sense"

# PCA (from notes)

## • Hidden Markov Models: Markov Chains

- Goal: make a sequence of decisions
  - Processes that unfold in time, states at time t are influenced by a state at time t-1
  - Applications: speech recognition, gesture recognition, parts of speech tagging and DNA sequencing,
  - Any temporal process without memory  $\omega^{\mathsf{T}} = \{\omega(1), \omega(2), \omega(3), \dots, \omega(\mathsf{T})\}$  sequence of states We might have  $\omega^6 = \{\omega 1, \omega 4, \omega 2, \omega 2, \omega 1, \omega 4\}$
  - The system can revisit a state at different steps and not every state need to be visited<sup>Pattern Classification, Chapter 3</sup>

### First-order Markov models

• Our productions of any sequence is described by the transition probabilities

 $\mathsf{P}(\omega_{i}(t+1) \mid \omega_{i}(t)) = \mathsf{a}_{ij}$ 



**FIGURE 3.8.** The discrete states,  $\omega_i$ , in a basic Markov model are represented by nodes, and the transition probabilities,  $a_{ij}$ , are represented by links. In a first-order discrete-time Markov model, at any step *t* the full system is in a particular state  $\omega(t)$ . The state at step t + 1 is a random function that depends solely on the state at step *t* and the transition probabilities. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

$$\begin{aligned} \theta &= (a_{ij}, \, \omega^{\mathsf{T}}) \\ \mathsf{P}(\omega^{\mathsf{T}} \mid \theta) &= a_{14} \, . \, a_{42} \, . \, a_{22} \, . \, a_{21} \, . \, a_{14} \\ \mathsf{P}(\omega(1) = \omega_{i}) \end{aligned}$$

**Example:** speech recognition

"production of spoken words"
Production of the word: "pattern" represented by phonemes
/p/ /a/ /tt/ /er/ /n/ // (// = silent state)
Transitions from /p/ to /a/, /a/ to /tt/, /tt/ to er/, /er/ to /n/ and /n/ to a silent\_state