



Lecture Slides for

INTRODUCTION TO

Machine Learning

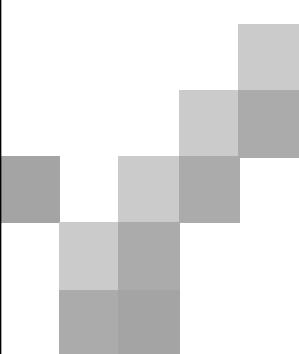
ETHEM ALPAYDIN

© The MIT Press, 2004

Edited for CS 536 Fall 2005 – Rutgers University
Ahmed Elgammal

alpaydin@boun.edu.tr

http://www.cmpe.boun.edu.tr/~ethem/i2ml



CHAPTER 4:

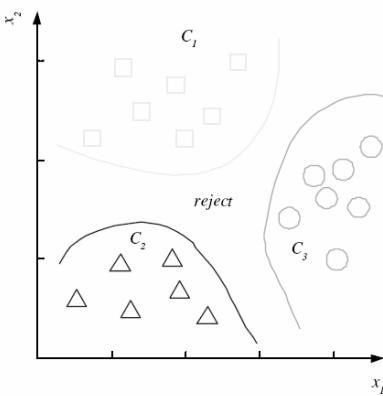
Parametric Methods

CHAPTER 5:

Multivariate Methods

Parametric Classification

- We would like to estimate $p(x/C_i)$ and $p(C_i)$ to be able to estimate the posterior $p(C_i/x)$



3

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Parametric Estimation

- $X = \{x^t\}_t$ where $x^t \sim p(x)$
- X is independent and identically distributed (iid) sample
- Parametric estimation:
 - Assume a form for $p(x | \theta)$ and estimate θ , its sufficient statistics, using X
 - e.g., $N(\mu, \sigma^2)$ where $\theta = \{\mu, \sigma^2\}$

4

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Maximum Likelihood Estimation

- Likelihood of θ given the sample X

$$I(\theta|X) = p(X|\theta) = \prod_t p(x^t|\theta)$$

- Log likelihood

$$\mathcal{L}(\theta|X) = \log I(\theta|X) = \sum_t \log p(x^t|\theta)$$

- Maximum likelihood estimator (MLE)

$$\theta^* = \operatorname{argmax}_{\theta} \mathcal{L}(\theta|X)$$

What is the parameter(s) of the distribution that maximizes the likelihood of the data sample

5

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Examples: Bernoulli/Multinomial

- Bernoulli: Two states, failure/success, x in {0,1}

$$P(x) = p_o^x (1 - p_o)^{1-x}$$

$$\mathcal{L}(p_o|X) = \log \prod_t p_o^{x^t} (1 - p_o)^{1-x^t}$$

$$\text{MLE: } p_o = \frac{\sum_t x^t}{N}$$

- Multinomial: $K > 2$ states, x_i in {0,1}

$$P(x_1, x_2, \dots, x_K) = \prod_i p_i^{x_i}$$

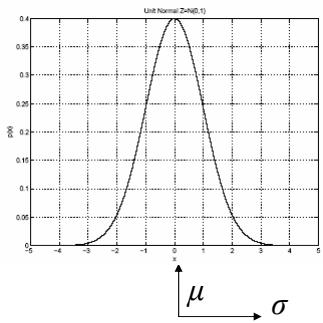
$$\mathcal{L}(p_1, p_2, \dots, p_K|X) = \log \prod_t \prod_i p_i^{x_i^t}$$

$$\text{MLE: } p_i = \frac{\sum_t x_i^t}{N}$$

6

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Gaussian (Normal) Distribution



- $p(x) = N(\mu, \sigma^2)$

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

- MLE for μ and σ^2 :

$$m = \frac{\sum x^t}{N}$$

$$s^2 = \frac{\sum (x^t - m)^2}{N}$$

7

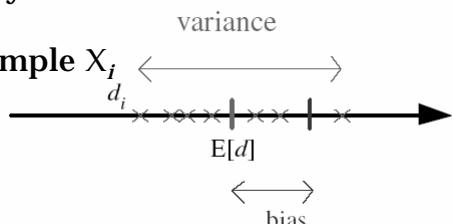
Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Bias and Variance

How to evaluate the quality of an estimator?

Unknown parameter θ

Estimator $d_i = d(X_i)$ on sample X_i



Bias: $b_\theta(d) = E[d] - \theta$

Variance: $E[(d-E[d])^2]$

Mean square error:

$$\begin{aligned} r(d, \theta) &= E[(d-\theta)^2] \\ &= (E[d] - \theta)^2 + E[(d-E[d])^2] \\ &= \text{Bias}^2 + \text{Variance} \end{aligned}$$

8

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Bayes' Estimator

- Treat θ as a random var with prior $p(\theta)$
- Bayes' rule: $p(\theta|X) = p(X|\theta) p(\theta) / p(X)$
- Full: $p(x|X) = \int p(x|\theta) p(\theta|X) d\theta$
- Maximum a Posteriori (MAP): $\theta_{\text{MAP}} = \operatorname{argmax}_{\theta} p(\theta|X)$
- Maximum Likelihood (ML): $\theta_{\text{ML}} = \operatorname{argmax}_{\theta} p(X|\theta)$
- Bayes': $\theta_{\text{Bayes'}} = E[\theta|X] = \int \theta p(\theta|X) d\theta$

9

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Bayes' Estimator: Example

- $x^t \sim N(\theta, \sigma_0^2)$ and $\theta \sim N(\mu, \sigma^2)$
- $\theta_{\text{ML}} = m$
- $\theta_{\text{MAP}} = \theta_{\text{Bayes'}} =$

$$E[\theta | X] = \frac{N/\sigma_0^2}{N/\sigma_0^2 + 1/\sigma^2} m + \frac{1/\sigma^2}{N/\sigma_0^2 + 1/\sigma^2} \mu$$

10

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Parametric Classification

$$g_i(x) = p(x | C_i)P(C_i)$$

or equivalently

$$g_i(x) = \log p(x | C_i) + \log P(C_i)$$

$$p(x | C_i) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left[-\frac{(x - \mu_i)^2}{2\sigma_i^2}\right]$$
$$g_i(x) = -\frac{1}{2} \log 2\pi - \log \sigma_i - \frac{(x - \mu_i)^2}{2\sigma_i^2} + \log P(C_i)$$

11

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

- Given the sample $\mathcal{X} = \{\mathbf{x}^t, \mathbf{r}^t\}_{t=1}^N$

$$\mathbf{x} \in \mathfrak{R} \quad \mathbf{r}_i^t = \begin{cases} \mathbf{1} & \text{if } \mathbf{x}^t \in C_i \\ \mathbf{0} & \text{if } \mathbf{x}^t \in C_j, j \neq i \end{cases}$$

- ML estimates are

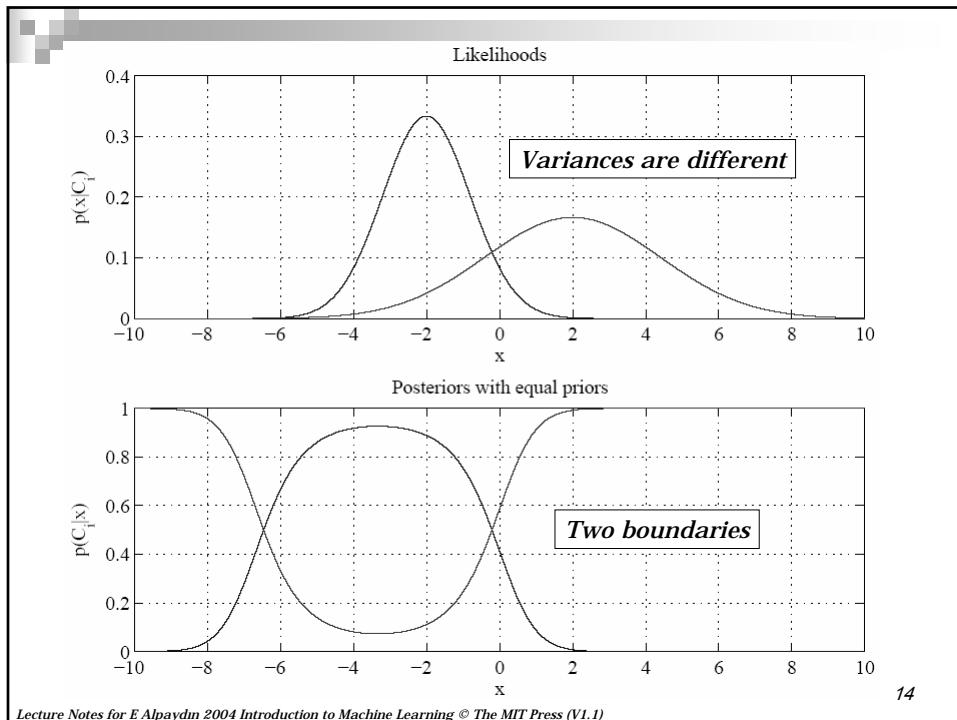
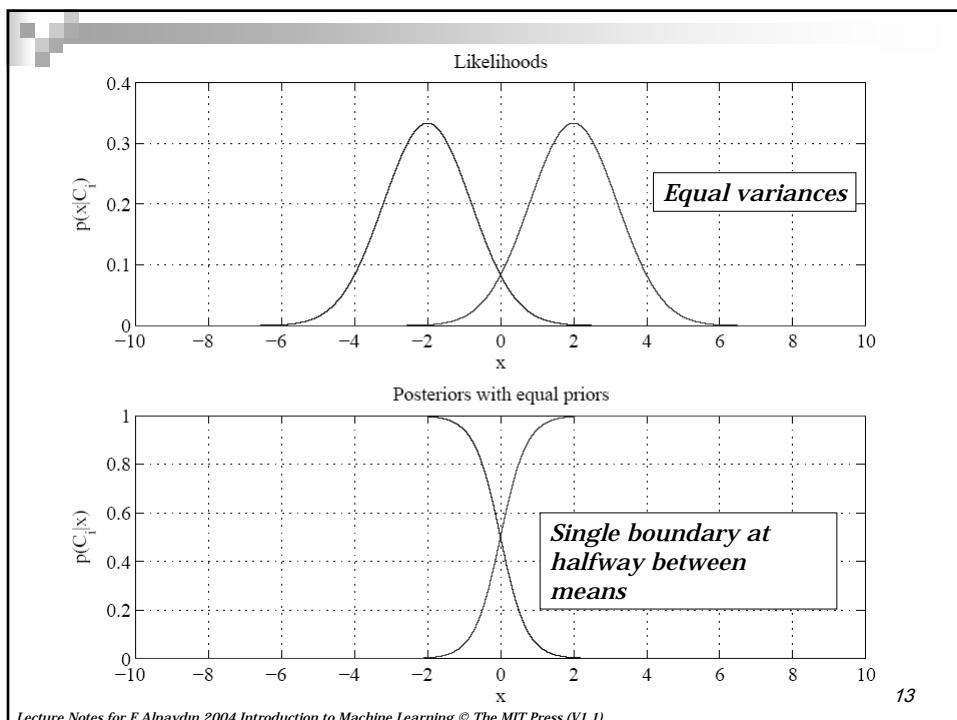
$$\hat{P}(C_i) = \frac{\sum r_i^t}{N} \quad \mathbf{m}_i = \frac{\sum \mathbf{x}^t r_i^t}{\sum r_i^t} \quad s_i^2 = \frac{\sum (\mathbf{x}^t - \mathbf{m}_i)^2 r_i^t}{\sum r_i^t}$$

- Discriminant becomes

$$g_i(x) = -\frac{1}{2} \log 2\pi - \log s_i - \frac{(x - m_i)^2}{2s_i^2} + \log \hat{P}(C_i)$$

12

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)



Multivariate Data

- Multiple measurements (sensors)
- d inputs/features/attributes: d -variate
- N instances/observations/examples

$$\mathbf{X} = \begin{bmatrix} X_1^1 & X_2^1 & \cdots & X_d^1 \\ X_1^2 & X_2^2 & \cdots & X_d^2 \\ \vdots & & & \\ X_1^N & X_2^N & \cdots & X_d^N \end{bmatrix}$$

15

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Multivariate Parameters

Mean : $E[\mathbf{x}] = \boldsymbol{\mu} = [\mu_1, \dots, \mu_d]^T$

Covariance : $\sigma_{ij} \equiv \text{Cov}(X_i, X_j)$

Correlation : $\text{Corr}(X_i, X_j) \equiv \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$

$$\Sigma \equiv \text{Cov}(\mathbf{X}) = E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2d} \\ \vdots & & & \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_d^2 \end{bmatrix}$$

16

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Parameter Estimation

$$\text{Sample mean } \mathbf{m} : m_i = \frac{\sum_{t=1}^N x_i^t}{N}, i = 1, \dots, d$$

$$\text{Covariance matrix } \mathbf{S} : s_{ij} = \frac{\sum_{t=1}^N (x_i^t - m_i)(x_j^t - m_j)}{N}$$

$$\text{Correlation matrix } \mathbf{R} : r_{ij} = \frac{s_{ij}}{s_i s_j}$$

17

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

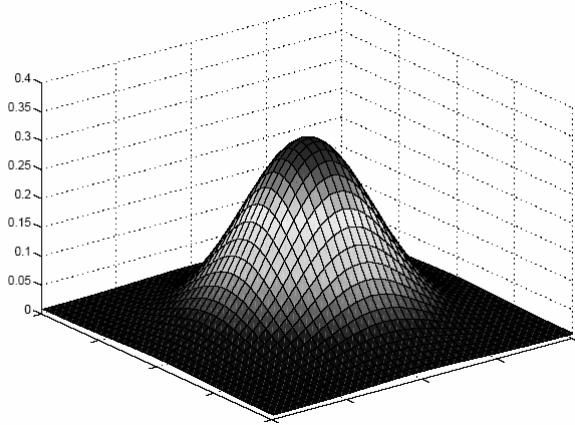
Estimation of Missing Values

- What to do if certain instances have missing attributes?
- Ignore those instances: not a good idea if the sample is small
- Use ‘missing’ as an attribute: may give information
- Imputation: Fill in the missing value
 - Mean imputation: Use the most likely value (e.g., mean)
 - Imputation by regression: Predict based on other attributes

18

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Multivariate Normal Distribution



$$\mathbf{x} \sim N_d(\boldsymbol{\mu}, \Sigma)$$

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

19

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Multivariate Normal Distribution

- Mahalanobis distance: $(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$
measures the distance from \mathbf{x} to $\boldsymbol{\mu}$ in terms of
(normalizes for difference in variances and
correlations)
- Bivariate: $d = 2$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

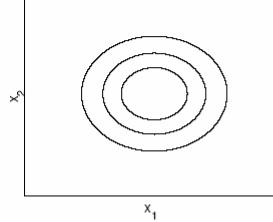
$$p(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left[-\frac{1}{2(1-\rho^2)} (z_1^2 - 2\rho z_1 z_2 + z_2^2) \right]$$
$$z_i = (\mathbf{x}_i - \boldsymbol{\mu}_i) / \sigma_i$$

20

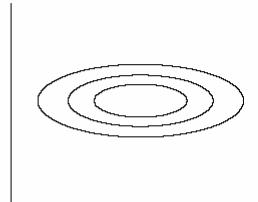
Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Bivariate Normal

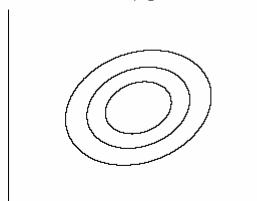
$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$



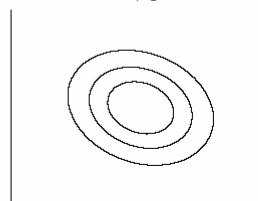
$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) > \text{Var}(x_2)$



$\text{Cov}(x_1, x_2) > 0$



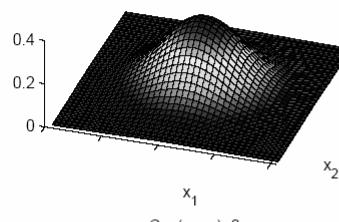
$\text{Cov}(x_1, x_2) < 0$



21

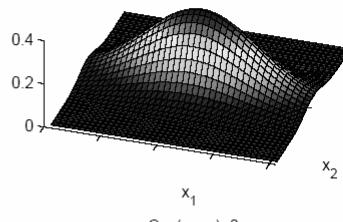
Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$

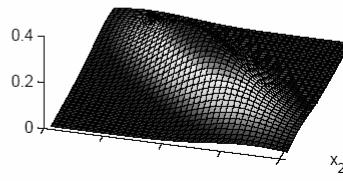
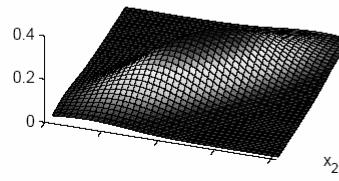


$\text{Cov}(x_1, x_2) > 0$

$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) > \text{Var}(x_2)$



$\text{Cov}(x_1, x_2) < 0$



22

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Independent Inputs: Naive Bayes

- If x_i are independent, offdiagonals of Σ are 0, Mahalanobis distance reduces to weighted (by $1/\sigma_i$) Euclidean distance:

$$p(\mathbf{x}) = \prod_{i=1}^d p_i(x_i) = \frac{1}{(2\pi)^{d/2} \prod_{i=1}^d \sigma_i} \exp \left[-\frac{1}{2} \sum_{i=1}^d \left(\frac{\mathbf{x}_i - \mu_i}{\sigma_i} \right)^2 \right]$$

- If variances are also equal, reduces to Euclidean distance

23

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Parametric Classification

- If $p(\mathbf{x} | C_i) \sim N(\mu_i, \Sigma_i)$

$$p(\mathbf{x} | C_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu_i)^T \Sigma_i^{-1} (\mathbf{x} - \mu_i) \right]$$

- Discriminant functions are

$$\begin{aligned} g_i(\mathbf{x}) &= \log p(\mathbf{x} | C_i) + \log P(C_i) \\ &= -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\mathbf{x} - \mu_i)^T \Sigma_i^{-1} (\mathbf{x} - \mu_i) + \log P(C_i) \end{aligned}$$

24

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Estimation of Parameters

$$\begin{aligned}\hat{P}(C_i) &= \frac{\sum_t r_i^t}{N} \\ \mathbf{m}_i &= \frac{\sum_t r_i^t \mathbf{x}^t}{\sum_t r_i^t} \\ \mathbf{s}_i &= \frac{\sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T}{\sum_t r_i^t}\end{aligned}$$

$$g_i(\mathbf{x}) = -\frac{1}{2} \log |\mathbf{s}_i| - \frac{1}{2} (\mathbf{x} - \mathbf{m}_i)^T \mathbf{s}_i^{-1} (\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i)$$

25

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

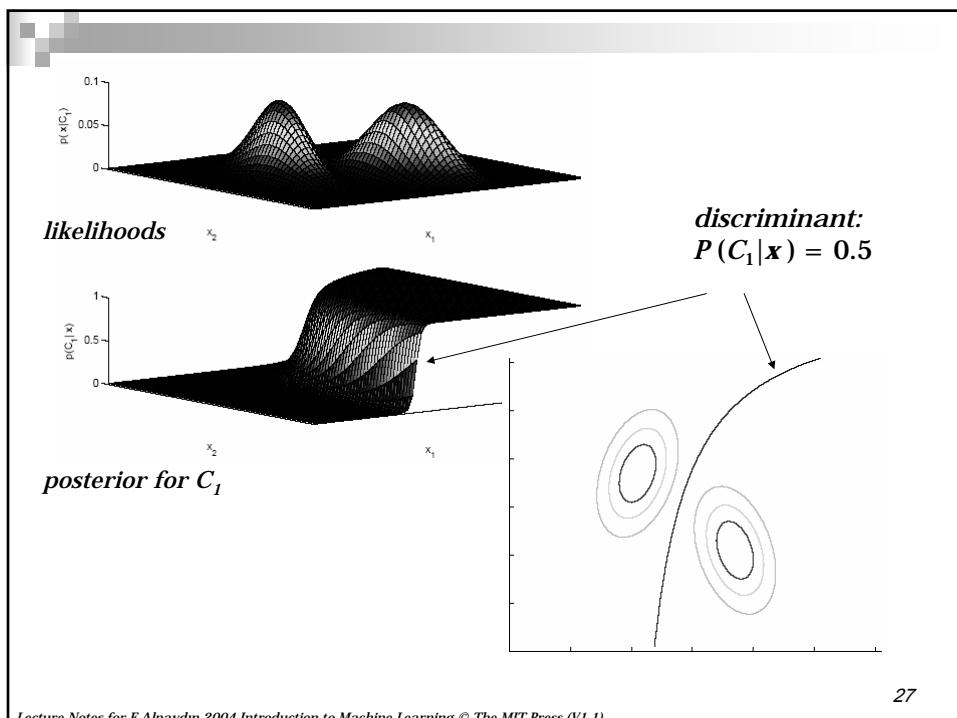
Different \mathbf{S}_i

■ Quadratic discriminant

$$\begin{aligned}g_i(\mathbf{x}) &= -\frac{1}{2} \log |\mathbf{s}_i| - \frac{1}{2} (\mathbf{x}^T \mathbf{s}_i^{-1} \mathbf{x} - 2 \mathbf{x}^T \mathbf{s}_i^{-1} \mathbf{m}_i + \mathbf{m}_i^T \mathbf{s}_i^{-1} \mathbf{m}_i) + \log \hat{P}(C_i) \\ &= \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0} \\ \text{where} \\ \mathbf{W}_i &= -\frac{1}{2} \mathbf{s}_i^{-1} \\ \mathbf{w}_i &= \mathbf{s}_i^{-1} \mathbf{m}_i \\ w_{i0} &= -\frac{1}{2} \mathbf{m}_i^T \mathbf{s}_i^{-1} \mathbf{m}_i - \frac{1}{2} \log |\mathbf{s}_i| + \log \hat{P}(C_i)\end{aligned}$$

26

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)



27

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Common Covariance Matrix \mathbf{S}

- Shared common sample covariance \mathbf{S}

$$\mathbf{S} = \sum_i \hat{P}(C_i) \mathbf{S}_i$$

- Discriminant reduces to

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}^{-1}(\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i)$$

which is a linear discriminant

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

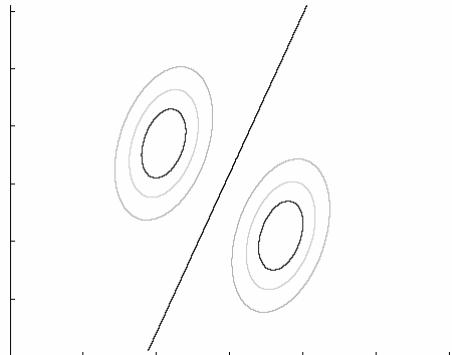
where

$$\mathbf{w}_i = \mathbf{S}^{-1} \mathbf{m}_i \quad w_{i0} = -\frac{1}{2} \mathbf{m}_i^T \mathbf{S}^{-1} \mathbf{m}_i + \log \hat{P}(C_i)$$

28

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Common Covariance Matrix S



29

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Diagonal S

- When $x_j, j = 1..d$, are independent, S is diagonal
 $p(\mathbf{x}|C_i) = \prod_j p(x_j|C_i)$ (Naive Bayes' assumption)

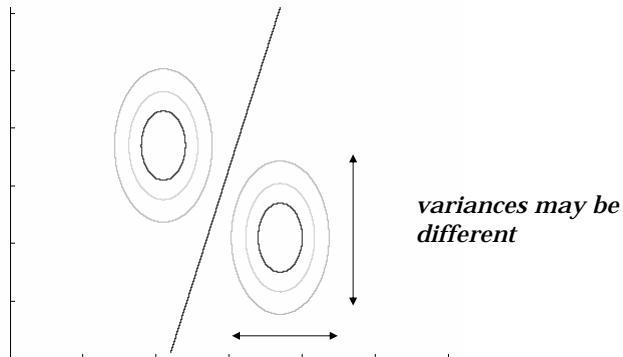
$$g_i(\mathbf{x}) = -\frac{1}{2} \left(\frac{\mathbf{x}_j^t - \mathbf{m}_j}{s_j} \right)^2 + \log \hat{P}(C_i)$$

Classify based on weighted Euclidean distance (in s_j units) to the nearest mean

30

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Diagonal S



31

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Diagonal S , equal variances

- Nearest mean classifier: Classify based on Euclidean distance to the nearest mean

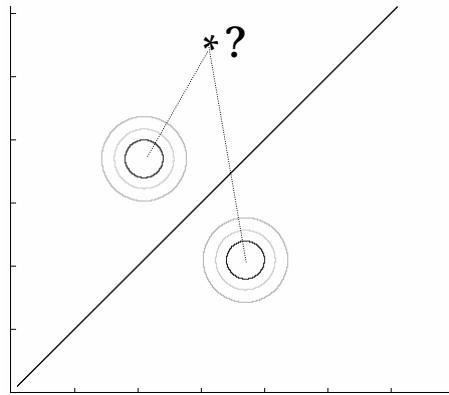
$$\begin{aligned}g_i(\mathbf{x}) &= -\frac{\|\mathbf{x} - \mathbf{m}_i\|^2}{2s^2} + \log \hat{P}(C_i) \\&= -\frac{1}{2s^2} \sum_{j=1}^d (\mathbf{x}_j^t - m_{ij})^2 + \log \hat{P}(C_i)\end{aligned}$$

- Each mean can be considered a prototype or template and this is template matching

32

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Diagonal S , equal variances



33

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Model Selection

<i>Assumption</i>	<i>Covariance matrix</i>	<i>No of parameters</i>
Shared, Hyperspheric	$\mathbf{S}_i = \mathbf{S} = s^2 \mathbf{I}$	1
Shared, Axis- aligned	$\mathbf{S}_i = \mathbf{S}$, with $s_{ij} = 0$	d
Shared, Hyperellipsoidal	$\mathbf{S}_i = \mathbf{S}$	$d(d+1)/2$
Different, Hyperellipsoidal	\mathbf{S}_i	$K d(d+1)/2$

■ As we increase complexity (less restricted S), bias

decreases and variance increases

■ Assume simple models (allow some bias) to control variance (regularization)

34

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Discrete Features

- **Binary features:** $p_{ij} \equiv p(x_j=1 | C_i)$
if x_j are independent (Naive Bayes')

$$p(\mathbf{x} | C_i) = \prod_{j=1}^d p_{ij}^{x_j} (1 - p_{ij})^{1-x_j}$$

the discriminant is linear

$$\begin{aligned} g_i(\mathbf{x}) &= \log p(\mathbf{x} | C_i) + \log P(C_i) \\ &= \sum_j [x_j \log p_{ij} + (1 - x_j) \log (1 - p_{ij})] + \log P(C_i) \end{aligned}$$

Estimated parameters $\hat{p}_{ij} = \frac{\sum_t x_j^t r_i^t}{\sum_t r_i^t}$

35

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Discrete Features

- **Multinomial (1-of- n_j) features:** $x_j \in \{v_1, v_2, \dots, v_{n_j}\}$
 $p_{ijk} \equiv p(z_{jk}=1 | C_i) = p(x_j=v_k | C_i)$
if x_j are independent

$$\begin{aligned} p(\mathbf{x} | C_i) &= \prod_{j=1}^d \prod_{k=1}^{n_j} p_{ijk}^{z_{jk}} \\ g_i(\mathbf{x}) &= \sum_j \sum_k z_{jk} \log p_{ijk} + \log P(C_i) \\ \hat{p}_{ijk} &= \frac{\sum_t z_{jk}^t r_i^t}{\sum_t r_i^t} \end{aligned}$$

36

Lecture Notes for E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)