

Multivariate Methods



Berlin Chen

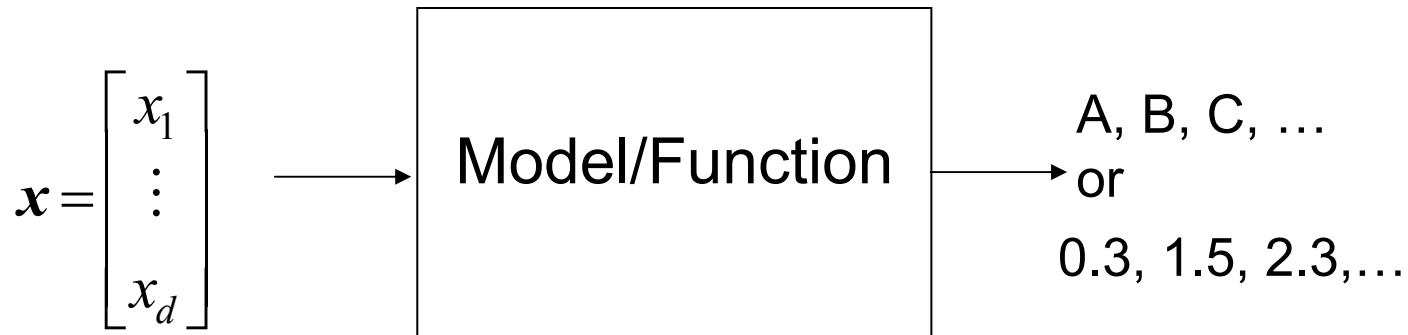
Graduate Institute of Computer Science & Information Engineering
National Taiwan Normal University

References:

1. *Introduction to Machine Learning*, Chapter 5

Multivariate Methods

- Input: a data sample with multiple features (variables/inputs)
- Output (of Prediction)
 - Classification: class code (discrete variable)
 - Regression: real number (continuous variable)



- Supervised learning
 - Model/function to be trained with labeled training samples

Multivariate Data (1/2)

- Each data sample is represented by an observation vector with d dimensions
 - Each dimension of the vector is termed input/feature/attribute
 - Dimensions with different types and value domains
- The whole data set of size N can be viewed as a data matrix

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

samples/instances/items

Inputs/features/attributes

Reduced to a naïve Bayes' classifier

- Features are usually assumed correlated ! $P(\mathbf{x}|C_i) = P(x_1, x_2, \dots, x_d | C_i) \stackrel{?}{=} \prod_{j=1}^d P(x_j | C_i)$
 - Otherwise, there is no need for a multivariate analysis

Multivariate Data (2/2)

- Motivations for multivariate data analysis
 - Simplification
 - Assume that the large body of data can be well summarized by means of relatively few parameters
 - Exploration
 - Predict the value of one variable from the values of other variables
 - Multivariate classification (Discrete)
 - Multivariate regression (Numeric)

Parameter Estimation (1/4)

- Mean of data samples

$$E[\mathbf{X}] = \int \mathbf{x} \cdot P(\mathbf{x}) d\mathbf{x} = \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix}$$

← Mean of Column 1 of matrix \mathbf{X}
← Mean of Column d of matrix \mathbf{X}

- Covariance matrix of data samples

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2d} \\ \vdots & \vdots & \vdots & \vdots \\ \sigma_{d1} & \cdots & \cdots & \sigma_d^2 \end{bmatrix}$$

σ_i^2 : variance of a variable X_i

σ_{ij} : covariace of two variables X_i and X_j

$$E[X_i X_j] = \int_{x_i} \int_{x_j} x_i x_j P(x_i, x_j) dx_i dx_j$$

$$\begin{aligned} \sigma_{ij} &\equiv \text{Cov}(X_i, X_j) \\ &= E[(X_i - \mu_i)(X_j - \mu_j)] = E[X_i X_j] - \mu_i \mu_j \end{aligned}$$

if X_i and X_j are independent

$$E[X_i X_j] = E[X_i] E[X_j] = \mu_i \mu_j$$

$$\begin{aligned} \sigma_{ii} &\equiv \text{Cov}(X_i, X_i) \\ &= E[(X_i - \mu_i)(X_i - \mu_i)] = E[(X_i)^2] - (\mu_i)^2 \\ &= \sigma_i^2 \end{aligned}$$

Parameter Estimation (2/4)

- Covariance matrix is symmetric
 - Diagonal terms: variances
 - Off-diagonal terms: covariances

$$\Sigma = \text{Cov}(\mathbf{X}) = E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = E[\mathbf{X}\mathbf{X}^T] - \boldsymbol{\mu}\boldsymbol{\mu}^T$$

- Correlation between two variables X_i and X_j

$$\text{Corr}(X_i, X_j) = \frac{\sigma_{ij}}{\sigma_i \sigma_j}, \quad -1 \leq \text{Corr}(X_i, X_j) \leq 1$$

independent \rightarrow uncorrelated

- Two variables independent \rightarrow covariance = correlation = 0
- But covariance = correlation = 0 does not imply two variables are independent (nonlinear dependence)

Parameter Estimation (3/4)

- If X_i and X_j are linear dependent

$$X_j = a \cdot X_i + b \quad (a > 0)$$

$$\Rightarrow E[X_j] = \mu_j = a \cdot [X_i] + b = a \cdot \mu_i + b$$

$$\sigma_j^2 = \text{var}(X_j) = E[(X_j - \mu_j)^2] = a^2 \cdot \text{var}(X_i) = a^2 \cdot \sigma_i^2$$

$$\sigma_{ij} = E[(X_i - \mu_i)(X_j - \mu_j)] = E[(X_i - \mu_i)((a \cdot X_i + b) - (a \cdot \mu_i + b))] = a \cdot E[(X_i - \mu_i)^2] = a \cdot \sigma_i^2$$

$$\Rightarrow \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j} = \frac{a \cdot \sigma_i^2}{\sigma_i \cdot (a \cdot \sigma_i)} = 1$$

Similarly, for

$$X_j = -a \cdot X_i + b \quad (a > 0)$$

$$\rho_{ij} = -1$$

Parameter Estimation (4/4)

- Maximum Likelihood Estimators

- Sample mean as an estimator for mean

$$\bar{\boldsymbol{\mu}} = \mathbf{m} = \frac{\sum_{t=1}^N \mathbf{x}^t}{N} \quad \text{with } m_i = \frac{\sum_{t=1}^N x_i^t}{N}, \quad i=1, \dots, d$$

- Sample covariance matrix as an estimator for covariance matrix

$$s_i^2 = \frac{\sum_{t=1}^N (x_i^t - m_i)^2}{N}$$
$$s_{ij} = \frac{\sum_{t=1}^N (x_i^t - m_i)(x_j^t - m_j)}{N}$$
$$\bar{\boldsymbol{\Sigma}} = \mathbf{C} = \begin{bmatrix} s_1^2 & s_{12} & \cdots & s_{1d} \\ s_{21} & s_2^2 & \cdots & s_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ s_{d1} & \cdots & \cdots & s_d^2 \end{bmatrix}$$

- Sample correlation coefficients

$$r_{ij} = \frac{s_{ij}}{s_i s_j}$$

Multivariate Normal Distribution (1/4)

- A random variable \mathbf{x} that is d -dimensional and normal distributed, is denoted as $\mathbf{x} \sim N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

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

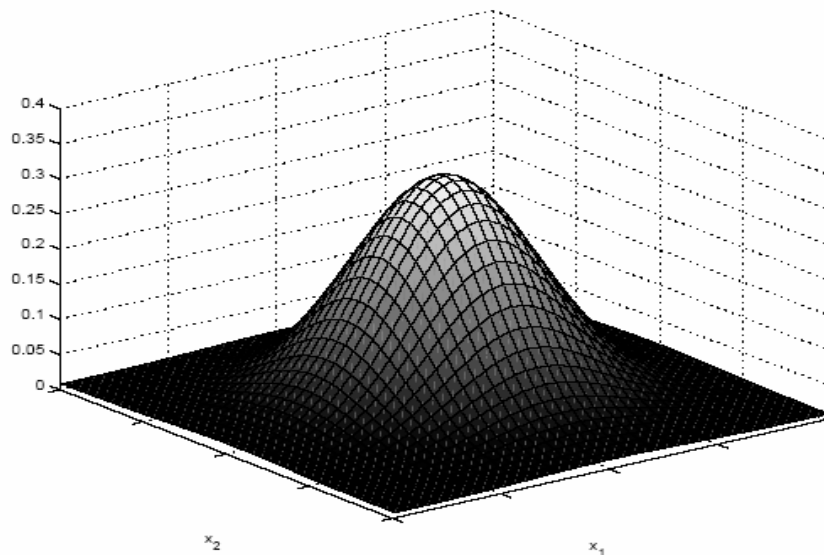
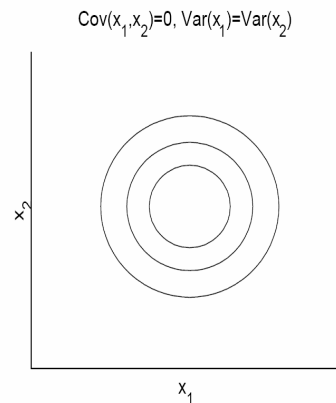


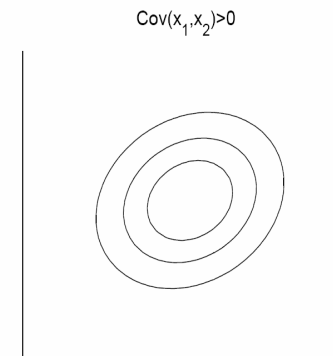
Figure 5.1: Bivariate normal distribution.

Multivariate Normal Distribution (2/4)

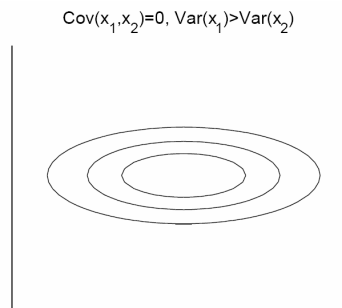
- If each dimension is independent of one another and with same variance value



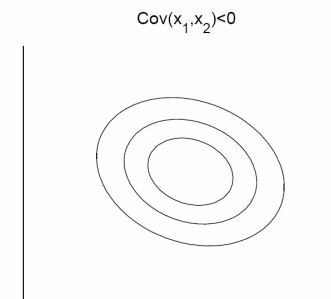
- If each dimension is dependent of one another with $\text{Corr}(X_1, X_2) > 0$



- If each dimension is independent of one another and with different variance value



- If each dimension is dependent of one another with $\text{Corr}(X_1, X_2) < 0$



Multivariate Normal Distribution (3/4)

- If the components of random variable \mathbf{x} are independent

$$\begin{aligned} P(\mathbf{x}) &= P(x_1, \dots, x_j, \dots, x_d) = \prod_{j=1}^d P(x_j) \\ &= \frac{1}{(2\pi)^{d/2} \prod_{j=1}^d \sigma_j} \exp \left[-\frac{1}{2} \sum_{j=1}^d \left(\frac{x_j - \mu_j}{\sigma_j} \right)^2 \right] \end{aligned}$$

Also note that
$$P(x_j) = \frac{1}{(2\pi)^{1/2} \sigma_j} \exp \left[-\frac{1}{2} \left(\frac{x_j - \mu_j}{\sigma_j} \right)^2 \right]$$

Multivariate Normal Distribution (4/4)

- Recall: The projection of a d -dimensional normal distribution on a vector \mathbf{w} is univariate normal (suppose that $\|\mathbf{w}\|=1$)

$$\mathbf{w}^T \mathbf{x} = w_1 x_1 + w_2 x_2 + \dots + w_d x_d \sim N(\mathbf{w}^T \boldsymbol{\mu}, \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w})$$

$$E[\mathbf{w}^T \mathbf{x}] = \mathbf{w}^T E[\mathbf{x}] = \underbrace{\mathbf{w}^T \boldsymbol{\mu}}_{\text{scalar}}$$

$$\begin{aligned} \text{Var}(\mathbf{w}^T \mathbf{x}) &= E\left[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})^T\right] \\ &= E\left[\mathbf{w}^T (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{w}\right] = \mathbf{w}^T E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \mathbf{w} = \underbrace{\mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w}}_{\text{scalar}} \end{aligned}$$

- The projection of a d -dimensional normal distribution to a k -dimensional space is k -variate normal

$$\mathbf{W}^T \mathbf{x} \sim N_k(\underbrace{\mathbf{W}^T \boldsymbol{\mu}}_{\substack{k\text{-dim} \\ \text{vector}}}, \underbrace{\mathbf{W}^T \boldsymbol{\Sigma} \mathbf{W}}_{\substack{k \times k \\ \text{matrix}}}) \quad \mathbf{W} = [\mathbf{w}_1 \quad \mathbf{w}_2 \quad \dots \quad \mathbf{w}_k]$$

Multivariate Classification (1/2)

- Normal density $N_d(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ as the class-conditional probability $P(\mathbf{x}|C_i)$ of random variable $\mathbf{x} \in \mathbb{R}^d$

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

- Define the discriminant function as

$$\begin{aligned} g_i(\mathbf{x}) &= \log P(\mathbf{x}|C_i) + \log P(C_i) \\ &= -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \log P(C_i) \end{aligned}$$

Multivariate Classification (2/2)

- Maximum likelihood (ML) training of classifiers
 - Given a set of labeled samples

$$\mathbf{X} = \{\mathbf{x}^t, \mathbf{r}^t\}, \quad r_i^t = 1 \text{ if } \mathbf{x}^t \in C_i \text{ and } 0 \text{ otherwise}$$

$$\hat{P}(C_i) = \frac{\sum_{t=1}^N r_i^t}{N}$$

$$\hat{\boldsymbol{\mu}}_i = \mathbf{m}_i \text{ (sample mean)} = \frac{\sum_{t=1}^N r_i^t \mathbf{x}^t}{\sum_{t=1}^N r_i^t}$$

$$\hat{\boldsymbol{\Sigma}}_i = \mathbf{S}_i \text{ (sample covariance matrix)} = \frac{\sum_{t=1}^N r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T}{\sum_{t=1}^N r_i^t}$$

Quadratic Discriminant (1/2)

- The discriminant function with normal class-conditional density can be expressed as a **quadratic discriminant**

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

quadratic discriminant

$$\mathbf{W}_i = -\frac{1}{2} \mathbf{S}_i^{-1}$$

$$\mathbf{w}_i = \mathbf{S}_i^{-1} \mathbf{m}_i$$

$$w_{i0} = -\frac{d}{2} \log 2\pi - \frac{1}{2} \mathbf{m}_i^T \mathbf{S}_i^{-1} \mathbf{m}_i - \frac{1}{2} \log |\mathbf{S}_i| + \log P(C_i)$$

Recall : second - order polynomial for scalar variable

$$g_i(x) = w_{i2}x^2 + w_{i1}x + w_{i0}$$

Quadratic Discriminant (2/2)

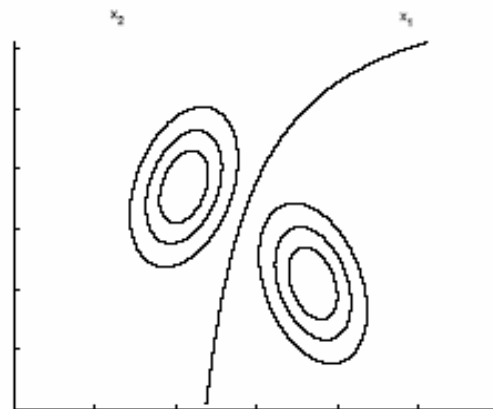


likelihood



posterior

$$\begin{aligned}
 g_1(\mathbf{x}) &= \mathbf{x}^T \mathbf{W}_1 \mathbf{x} + \mathbf{w}_1^T \mathbf{x} + w_{10} \\
 g_2(\mathbf{x}) &= \mathbf{x}^T \mathbf{W}_2 \mathbf{x} + \mathbf{w}_2^T \mathbf{x} + w_{20} \\
 \tilde{g}(\mathbf{x}) &= g_1(\mathbf{x}) - g_2(\mathbf{x}) \\
 &= \mathbf{x}^T (\mathbf{W}_1 - \mathbf{W}_2) \mathbf{x} + (\mathbf{w}_1^T - \mathbf{w}_2^T) \mathbf{x} \\
 &\quad + (w_{10} - w_{20}) \\
 &= \mathbf{x}^T \mathbf{W}' \mathbf{x} + \mathbf{w}'^T \mathbf{x} + w'_0
 \end{aligned}$$



likelihood

Figure 5.3: Classes have different covariance matrices.

Linear Discriminant (1/2)

- The discriminant function with normal class-conditional density **sharing the same covariance matrix** can be expressed as a **linear discriminant**

$$\mathbf{S} = \sum_{i=1}^K \hat{P}(C_i) \cdot \mathbf{S}_i \quad \text{All classes share the same covariance matrix}$$

$$g_i(\mathbf{x}) = \log P(\mathbf{x}|C_i) + \log P(C_i)$$

$$= -\frac{1}{2}(\cancel{\mathbf{x}^T \mathbf{S}^{-1} \mathbf{x}} - 2\mathbf{x}^T \mathbf{S}^{-1} \mathbf{m}_i + \mathbf{m}_i^T \mathbf{S}^{-1} \mathbf{m}_i) - \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\mathbf{S}| + \log P(C_i)$$

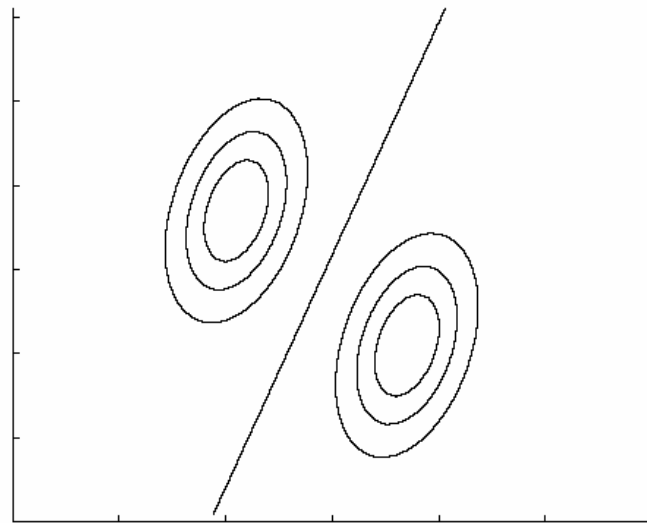
$$\Rightarrow \tilde{g}_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

$$\mathbf{w}_i = \mathbf{S}^{-1} \mathbf{m}_i$$

$$w_{i0} = -\frac{d}{2} \log 2\pi - \frac{1}{2} \mathbf{m}_i^T \mathbf{S}^{-1} \mathbf{m}_i - \frac{1}{2} \log |\mathbf{S}| + \log P(C_i)$$

Linear Discriminant (2/2)

covariances (off-diagonal terms of covariance matrix) are not equal to zero



$$\begin{aligned}g_1(\mathbf{x}) &= \mathbf{w}_1^T \mathbf{x} + w_{10} \\g_2(\mathbf{x}) &= \mathbf{w}_2^T \mathbf{x} + w_{20} \\ \tilde{g}(\mathbf{x}) &= g_1(\mathbf{x}) - g_2(\mathbf{x}) \\ &= (\mathbf{w}_1^T - \mathbf{w}_2^T) \mathbf{x} + (w_{10} - w_{20}) \\ &= \mathbf{w}'^T \mathbf{x} + w'_0\end{aligned}$$

likelihood

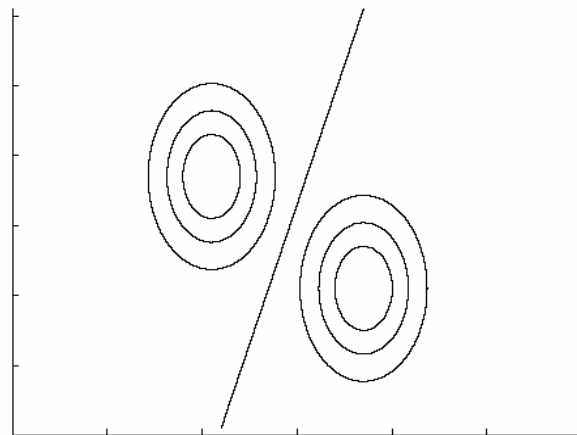
Figure 5.4: Covariances may be arbitrary but shared by both classes.

Naïve Bayes' Classifier (1/2)

- The discriminant function with a normal class-conditional density **sharing the same diagonal covariance matrix** can be expressed as a **naïve Bayes' classifier**

$$g_i(\mathbf{x}) = -\frac{1}{2} \sum_{j=1}^d \left(\frac{x_j - m_{ij}}{s_j} \right)^2 + \log P(C_i)$$

variances (diagonal terms of covariance matrix) are not equal



likelihood

Figure 5.5: All classes have equal, diagonal covariance matrices but variances are not equal.

Naïve Bayes' Classifier (2/2)

- If the variances (diagonal terms) of the naïve Bayes' classifier are further set to equal for all dimensions

$$g_i(\mathbf{x}) = -\frac{1}{2s^2} \sum_{j=1}^d (x_j - m_{ij})^2 + \log P(C_i)$$

The Mahalanobis distance is reduced to Euclidean distance. (all variables/features have the same variance and are independent)

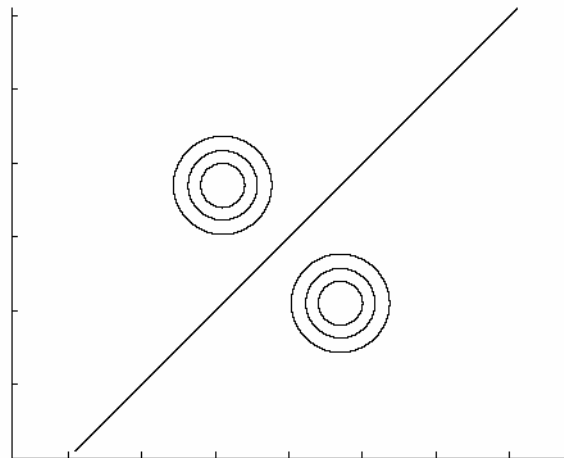


Figure 5.6: All classes have equal, diagonal covariance matrices of equal variances on both dimensions.

Nearest Mean Classifier

- Assign the data sample to the class of the nearest mean
 - If the priors $P(C_i)$ are further set to equal

$$\begin{aligned}g_i(\mathbf{x}) &= -\sum_{j=1}^d (x_j - m_{ij})^2 = -\|\mathbf{x} - \mathbf{m}_i\|^2 \\ &= -(\mathbf{x} - \mathbf{m}_i)^T (\mathbf{x} - \mathbf{m}_i) \\ &= -(\cancel{\mathbf{x}^T \mathbf{x}} - 2\mathbf{m}_i^T \mathbf{x} + \mathbf{m}_i^T \mathbf{m}_i)\end{aligned}$$

$$\begin{aligned}\Rightarrow \bar{g}_i(\mathbf{x}) &= \mathbf{m}_i^T \mathbf{x} - \frac{1}{2} \mathbf{m}_i^T \mathbf{m}_i \\ &= \mathbf{w}_i^T \mathbf{x} + w_{i0}\end{aligned}$$

where $\mathbf{w}_i = \mathbf{m}_i$

$$w_{i0} = -\frac{1}{2} \mathbf{m}_i^T \mathbf{m}_i = -\frac{1}{2} \|\mathbf{m}_i\|^2$$

Tuning Complexity

- Tradeoff between the bias and variance of an estimator

Table 5.1 Reducing variance through simplifying assumptions.

Assumption	Covariance matrix	No. of parameters
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	$K \cdot (d(d+1)/2)$

- Simplifying covariance matrix \rightarrow decreasing number of parameters \rightarrow introducing estimation bias
- Arbitrary covariance matrix \rightarrow much more data is needed \rightarrow introducing estimation variance
- Regularized Discriminant Analysis (RDA, 1989)
 - Use a weighted average of three special cases of covariance matrix

$$\mathbf{S}' = \alpha \cdot \underline{\sigma^2 \mathbf{I}} + \beta \cdot \underline{\mathbf{S}} + (1 - \alpha - \beta) \cdot \underline{\mathbf{S}_i}$$

a shared diagonal covariance matrix (with equal variance) for all classes

a shared covariance matrix for all classes

a specific covariance matrix for each class

Discrete Features (1/4)

- Features that take one of n different values
 - E.g. $\text{pixel} \in \{\text{on}, \text{off}\}; \text{color} \in \{\text{red}, \text{blue}, \text{green}, \text{black}\}$
- If each feature is a binary random variable x_j^t ($x_j^t = 1$ or $x_j^t = 0$)
 - Bernoulli distribution

$$P_{ij} = P(x_j^t = 1 | C_i) \quad 1 - P_{ij} = P(x_j^t = 0 | C_i) \quad \Rightarrow P(x_j^t | C_i) = (P_{ij})^{x_j^t} (1 - P_{ij})^{(1-x_j^t)}$$

- If features are further assumed to be independent

- Assumption for Naïve Bayes' classifier

$$\mathbf{x}^t = \{x_1^t, \dots, x_j^t, \dots, x_d^t\}$$

$$P(\mathbf{x}^t | C_i) = P(x_1^t, \dots, x_d^t | C_i) = \prod_{j=1}^d P(x_j^t | C_i) = \prod_{j=1}^d (P_{ij})^{x_j^t} (1 - P_{ij})^{(1-x_j^t)}$$

- Discriminant function

$$g_i(\mathbf{x}^t) = \log P(\mathbf{x}^t | C_i) + \log P(C_i) = \sum_{j=1}^d [x_j^t \log P_{ij} + (1 - x_j^t) \log(1 - P_{ij})] + \log P(C_i)$$

- Linear ?

Discrete Features (2/4)

- Appendix A

$$\begin{aligned}g_i(\mathbf{x}^t) &= \sum_{j=1}^d \left[x_j^t \log P_{ij} + (1-x_j^t) \log(1-P_{ij}) \right] + \log P(C_i) \\ &= \left[\sum_{j=1}^d x_j^t \log P_{ij} \right] - \left[\sum_{j=1}^d x_j^t (1 - \log P_{ij}) \right] + \left[\sum_{j=1}^d (1 - \log P_{ij}) \right] + \log P(C_i) \\ &= \left[\sum_{j=1}^d x_j^t (2 \cdot (\log P_{ij}) - 1) \right] + \left[\sum_{j=1}^d (1 - \log P_{ij}) \right] + \log P(C_i) \\ &= \mathbf{w}_i^T \mathbf{x}^t + D_i \quad \text{Linear}\end{aligned}$$

where

$$\mathbf{w}_i = \begin{bmatrix} 2 \cdot (\log P_{i1}) - 1 \\ 2 \cdot (\log P_{i2}) - 1 \\ \vdots \\ 2 \cdot (\log P_{id}) - 1 \end{bmatrix} \quad D_i = \left[\sum_{j=1}^d (1 - \log P_{ij}) \right] + \log P(C_i)$$

Discrete Features (3/4)

- Maximum likelihood estimation for P_{ij} (for binary variable)

$$\hat{P}_{ij} = \frac{\sum_{t=1}^N x_j^t \cdot r_i^t}{\sum_{t=1}^N r_i^t}$$

- Extension: features are independent **multinomial random variables**

$$x_j^t \in \{v_{j1}, \dots, v_{jk}, \dots, v_{jn_j}\} \quad z_{jk}^t = \begin{cases} 1 & \text{if } x_j^t = v_{jk} \\ 0 & \text{otherwise} \end{cases}$$

- Define the probability that x_j^t belongs to C_i and take v_{jk} :

$$P_{ijk} = P(z_{jk}^t = 1 | C_i) = P(x_j^t = v_{jk} | C_i) \Rightarrow \prod_{j=1}^d P(x_j^t | C_i) = \prod_{k=1}^{n_j} P_{ijk}^{z_{jk}^t}$$

Discrete Features (4/4)

- Assumption for Naïve Bayes' classifier

$$P(\mathbf{x}^t | C_i) = P(x_1^t, \dots, x_d^t | C_i) = \prod_{j=1}^d P(x_j^t | C_i) = \prod_{j=1}^d \prod_{k=1}^{n_j} P_{ijk}^{z_{jk}^t}$$

- Discriminant function

$$g_i(\mathbf{x}^t) = \log P(\mathbf{x}^t | C_i) + \log P(C_i) = \sum_{j=1}^d \sum_{k=1}^{n_j} z_{jk}^t \log P_{ijk} + \log P(C_i)$$

- Linear ?

- Maximum likelihood estimation for P_{ijk}

$$\hat{P}_{ijk} = \frac{\sum_{t=1}^N z_{jk}^t \cdot r_i^t}{\sum_{t=1}^N r_i^t}$$