

Probabilistic Modeling

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

NetDB-ML, Spring 2015

1 More About Probabilistic Modeling

2 MAP and Bayesian Estimation

3 The Bias/Variance Dilemma

4 Generative Methods

- Univariate Classification
- Maximum Likelihood Estimation
- Multivariate Classification
- Tuning the Model Complexity

Summary of Supervised Learning Models

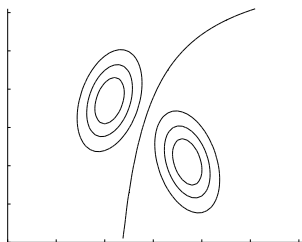
- Three main categories (either parametric or non-parametric):
 - 1 Those learning the **discriminant functions** f 's (no probability interpretation)
 - E.g., perceptron, k NN, etc.
 - 2 Those based on probability and learn $p(r|\mathbf{x})$ directly
 - E.g., linear regression, logistic regression, etc.
 - $p(r|\mathbf{x};\theta)$ with θ (constant) estimated from \mathcal{X}
 - Methods in 1 and 2 are called **discriminative methods**
 - 3 Those learn $p(r|\mathbf{x})$ indirectly from $p(\mathbf{x}|r)p(r)$
 - To be discussed later
 - These are called **generative methods**, as $p(\mathbf{x}|r)p(r)$ explains how \mathbf{x} (and \mathcal{X}) is generated

Probabilistic Modeling

- By assuming the target follows some probability distribution
- Pros and cons?

Probabilistic Modeling

- By assuming the target follows some probability distribution
- Pros and cons?
- Perform well only when the assumption holds
- Essentially solves a problem (i.e., distribution estimation) harder than discrimination
 - E.g., in generative models, if we let $p(\mathbf{x}|r) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then we can plot the contour of each class in addition to the decision boundary
 - Less efficient; but more descriptive



More About Probabilistic Modeling (1)

- The roles of θ in the prediction function $p(r'|\mathbf{x}')$:
 - Constant, from ML estimation of θ :
 - $\theta_{ML} = \arg \max_{\theta} p(\mathcal{X}|\theta)$
 - $p(r|\mathbf{x}') := p(r|\mathbf{x}'; \theta_{ML})$
 - Constant, from MAP estimation of θ :
 - $\theta_{MAP} = \arg \max_{\theta} p(\theta|\mathcal{X}) = \arg \max_{\theta} p(\mathcal{X}|\theta)p(\theta)$
 - $p(r|\mathbf{x}) := p(r|\mathbf{x}; \theta_{MAP})$
 - **Random variable**, for full Bayesian treatment:
 - $p(y|\mathbf{x}, \mathcal{X}) = \int p(y, \theta|\mathbf{x}', \mathcal{X}) d\theta$

More About Probabilistic Modeling (2)

- Can we analyze the generation performance more easily with the aid of distribution assumption?

More About Probabilistic Modeling (3)

- Generative models

The Roles of θ

- The roles of θ in the discrimination function $p(r'|\mathbf{x}')$:
 - Constant, from ML estimation of θ :
 - $\theta_{ML} = \arg \max_{\theta} p(\mathcal{X}|\theta)$
 - $p(r'|\mathbf{x}') := p(r|\mathbf{x}; \theta_{ML})$
 - Constant, from MAP estimation of θ :
 - $\theta_{MAP} = \arg \max_{\theta} p(\theta|\mathcal{X}) = \arg \max_{\theta} p(\mathcal{X}|\theta)p(\theta)$
 - $p(r'|\mathbf{x}') := p(r'|\mathbf{x}'; \theta_{MAP})$
 - **Random variable**, for full Bayesian treatment of r' :
 - $p(r'|\mathbf{x}', \mathcal{X}) = \int p(r', \theta|\mathbf{x}', \mathcal{X}) d\theta$

ML Estimator for θ

- The estimators we discussed so far (e.g., ρ_i , \mathbf{m}_i , and \mathbf{S}_i in classification and \mathbf{w} in regression) are called the **Maximum Likelihood (ML) estimators** since they are derived from

$$\theta_{ML} = \arg_{\theta} \max p(\mathcal{X}|\theta)$$

- E.g., in linear regression where $\theta = \mathbf{w}$, given a new instance \mathbf{x}' , the prediction can be made by

$$y' = \arg_y \max p(y|\mathbf{x}'; \mathbf{w}_{ML}) = \arg_y \max \mathcal{N}(y|\mathbf{w}_{ML}^T \mathbf{x}', \beta^{-1}) = \mathbf{w}_{ML}^T \mathbf{x}'$$

MAP Estimator for θ

- If we have the prior knowledge about θ (i.e., $P(\theta)$), we can obtain the **Maximum A Posteriori (MAP) estimators** based on

$$\theta_{MAP} = \arg_{\theta} \max P(\theta|\mathcal{X}) = \arg_{\theta} \max p(\mathcal{X}|\theta)P(\theta)$$

- If we assume that $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha^{-1}\mathbf{I})$ in linear regression, we have

$$\log p(\mathbf{w}|\mathcal{X}) = \log p(\mathcal{X}|\mathbf{w}) + \log p(\mathbf{w}) \propto$$

$$-\frac{\beta}{2} \sum_{t=1}^N (r^{(t)} - \mathbf{w}^T \mathbf{x}^{(t)})^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \quad [\text{Proof}]$$

- We effectively find \mathbf{w}_{MAP} that minimizes

$$\sum_{t=1}^N (r^{(t)} - \mathbf{w}^T \mathbf{x}^{(t)})^2 + \lambda \mathbf{w}^T \mathbf{w}, \text{ where } \lambda = \alpha/\beta$$

- In addition to minimizing the SSE, we **regularize** the norm of \mathbf{w} to prevent a highly complex model, thereby reducing the generalization error
- $y' = \arg_y \max p(y|\mathbf{x}'; \mathbf{w}_{MAP}) = \arg_y \max \mathcal{N}(y|\mathbf{w}_{MAP}^T \mathbf{x}', \beta^{-1}) = \mathbf{w}_{MAP}^T \mathbf{x}'$

Bayesian Estimator for r'

- The above methods treat θ as a deterministic value when making predictions
- Another technique, called the *Bayesian estimation* of r' , treats θ as a random variable, and considers all possible values of θ when estimating r' :
 - $y' = \arg_y \max p(y|\mathbf{x}', \mathcal{X}) = \int p(y, \theta|\mathbf{x}', \mathcal{X}) d\theta$
 - E.g., in linear regression,
 $y' = \arg_y \max p(y|\mathbf{x}', \mathcal{X}) = \arg_y \max \int p(y, \mathbf{w}|\mathbf{x}', \mathcal{X}) d\mathbf{w}$
 - No separated estimation phase for θ
- We will discuss how to solve y' in the lecture of graphical models

Regression Revisited (1)

- Given $\mathcal{X} = \{\mathbf{x}^{(t)}, r^{(t)}\}_{t=1}^N$, where $r^{(t)} \in \mathbb{R}$. Assume
 - $(\mathbf{x}^{(t)}, r^{(t)})$ are i.i.d samples drawn from some joint distribution of \mathbf{x} and r (otherwise can never learn r from \mathbf{x})
 - In particular, $r^{(t)} = f(\mathbf{x}^{(t)}; \theta) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \beta^{-1})$ for some **hyperparameter** (i.e., constant fixed during the objective solving) β
 - The marginal distribution $p(r|\mathbf{x})$ follows: $p(r|\mathbf{x}) = p_{N_{h(\mathbf{x}; \theta), \beta^{-1}}}(r)$
- We want to estimate f using \mathcal{X}
 - Hypothesis: $h(\mathbf{x}; w_0, w_1, \dots, w_d) = w_0 + w_1 x_1 + \dots + w_d x_d$, a line
 - Once getting w_0, w_1, \dots, w_d , we can predict the unknown r' of a new instance \mathbf{x}' by
$$y' = \arg_y \max p(y|\mathbf{x}') = \arg_y \max p_{N_{h(\mathbf{x}'; \theta), \beta^{-1}}}(y) = h(\mathbf{x}'; \theta)$$
 - Note that we don't need to know β to make prediction

Regression Revisited (2)

- How to obtain the estimate h of f ? How to obtain θ ?
- We can pick θ maximizing $p(\theta|\mathcal{X})$, the **posterior** probability
- Or, by Baye's theorem, θ maximizing the **likelihood** $p(\mathcal{X}|\theta)$ (if we assume $p(\theta)$ remains the same for all θ)
- Or, θ maximizing the **log likelihood** $\log p(\mathcal{X}|\theta) = \log \left(\prod_{t=1}^N p(\mathbf{x}^{(t)}, r^{(t)}|\theta) \right) = \log \left(\prod_{t=1}^N p(r^{(t)}|\mathbf{x}^{(t)}, \theta) p(\mathbf{x}^{(t)}|\theta) \right) = \log \left(\prod_{t=1}^N p(h(\mathbf{x}^{(t)}; \theta) + \epsilon|\mathbf{x}^{(t)}, \theta) p(\mathbf{x}^{(t)}|\theta) \right)$
- Ignoring $p(\mathbf{x}^{(t)}|\theta) = p(\mathbf{x}^{(t)})$ (since it is irrelevant to θ) and constants we have $\log p(\mathcal{X}|\theta) \propto -N \log \left(\sqrt{\frac{2\pi}{\beta}} \right) - \frac{\beta}{2} \sum_{t=1}^N (r^{(t)} - h(\mathbf{x}^{(t)}; \theta))^2$
- Dropping the first term and constants we have $\log p(\mathcal{X}|\theta) \propto -\sum_{t=1}^N (r^{(t)} - h(\mathbf{x}^{(t)}; \theta))^2$; that is, we seek for θ minimizing the SSE (sum of square errors)

The Bias/Variance Dilemma (1/4)

- The likelihood-based classification and regression share the same idea that the estimators $h(x; \theta_x)$ are obtained by $\theta_x = \arg_{\theta} \max p(\mathcal{X}|\theta)$
 - In classification, $h(x; \theta_x)$ estimates the discriminant of a class; in regression, $h(x; \theta_x)$ estimates f
- Given a new instance x' where r' is unknown, the expected square error (over the joint distribution of (x, r)) of our prediction can be written as

$$\begin{aligned} E[(r - h(x'; \theta_x))^2 | x'] &= \int (r - h(x'; \theta_x))^2 p(r | x') dr \\ &= \int [(r - E[r | x']) + (E[r | x'] - h(x'; \theta_x))]^2 p(r | x') dr \\ &= \int (r - E[r | x'])^2 p(r | x') dr + (E[r | x'] - h(x'; \theta_x))^2 \int p(r | x') dr - 2 \cdot 0 \\ &= E[(r - E[r | x'])^2 | x'] + (E[r | x'] - h(x'; \theta_x))^2 \end{aligned}$$

- The first term does not depend on h but the assumption of the joint distribution of (x, r)

The Bias/Variance Dilemma (2/4)

- The second term changes as we vary our hypothesis h and its complexity
- Note that in regression,
 $E[r|x'] = E[f(x') + \epsilon|x'] = f(x') + E[\epsilon|x'] = f(x')$ so the second term measures how our estimator h is difference from its target f
 - The similar argument applies to the case of classification
- Recall that we can measure how good the estimator h is by using the mean square error $E_{\mathcal{X}}[(h - f)^2]$ **over all possible \mathcal{X} of the same size**¹
- Since h and f are functions, we can rewrite the mean square error as follows given an instance x' :

$$\begin{aligned} E_{\mathcal{X}}[(h(x'; \theta_{\mathcal{X}}) - E[r|x'])^2|x'] &= \text{bias}^2 + \text{variance} \\ &= (E_{\mathcal{X}}[h(x'; \theta_{\mathcal{X}})] - E[r|x'])^2 + E_{\mathcal{X}}[(h(x'; \theta_{\mathcal{X}}) - E_{\mathcal{X}}[h(x'; \theta_{\mathcal{X}})])^2] \end{aligned}$$

¹Here we distinguish $E_{\mathcal{X}}$ (over \mathcal{X}) from E (over the joint distribution of (x, r))

The Bias/Variance Dilemma (3/4)

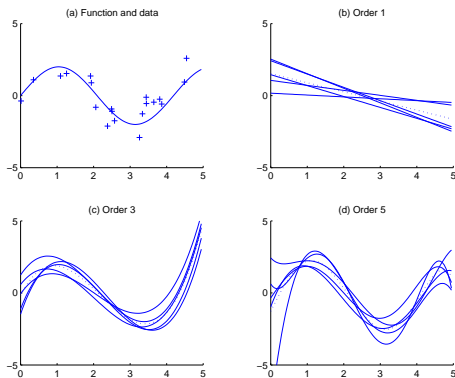


Figure : (a) A function $f(x) = 2\sin(1.5x)$ and a noisy training set ($\epsilon =_{s.t.} N_{0,1}$) consisting of 20 examples. There are totally 5 training sets \mathcal{X}_i , $1 \leq i \leq 5$, generated to calculate $E_{\mathcal{X}}$. (b), (c), and (d) are 5 polynomial fits, namely $h(x; \theta_{\mathcal{X}_i})$ of order 1, 3, and 5 respectively. For each case, the dotted line shows the average of the 5 fits, namely $E_{\mathcal{X}_i} [h(x; \theta_{\mathcal{X}_i})]$.

The Bias/Variance Dilemma (4/4)

- As we can see, a complex (i.e., high order) hypothesis h has
 - Low bias, as the average of the 5 fits is close to f
 - But high variance, as its shape is affected by noise
 - The variance decreases as N increase, since when N is large the different training sets \mathcal{X}_i look similar
- This is a mathematical way to justify: generalization error \propto empirical error + (model complexity / N)
 - Empirical error corresponds to the bias
 - The second term corresponds to the variance

Model Selection

- The right order of h can be determined using the cross validation technique
- Given the validation results at right (the dotted line), which order should we take?

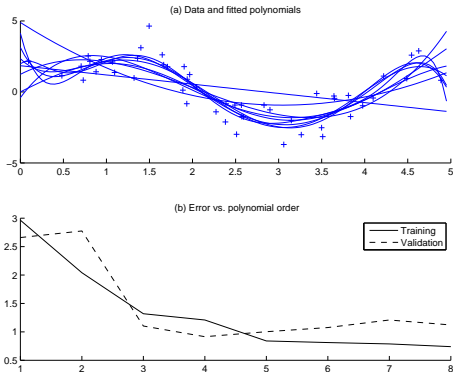


Figure : Cross validation results of 8 hypotheses with orders 1 to 8. Both the training and cross validation sets contain 50 instances.

Model Selection

- The right order of h can be determined using the cross validation technique
- Given the validation results at right (the dotted line), which order should we take?
 - Why not 4? Occam's razor tells us that we should choose the simplest hypothesis provided that its error is comparable
 - Note the validation results may not be as V-shaped as we might expect when N is large

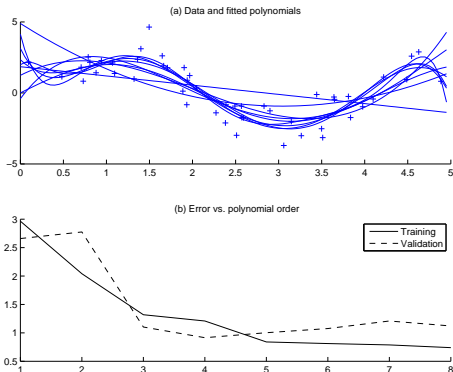


Figure : Cross validation results of 8 hypotheses with orders 1 to 8. Both the training and cross validation sets contain 50 instances.

Outline

1 More About Probabilistic Modeling

2 MAP and Bayesian Estimation

3 The Bias/Variance Dilemma

4 Generative Methods

- Univariate Classification
- Maximum Likelihood Estimation
- Multivariate Classification
- Tuning the Model Complexity

Univariate Classification

- Given a training set $\mathcal{X} = \{x^{(t)}, r^{(t)}\}_{t=1}^N$, where $r_i^{(t)} = 1$ if $x^{(t)} \in C_i$ and 0 otherwise, we find the discriminant $f_i(x) = P(C_i|x)$ for each class C_i , and then classify a new instance x' as $C_{y'}$ if $y' = \arg_i \max P(C_i|x)$
- Based on the generative assumption and Bayes' rule, we pick C_i such that $f_i(x') = \log(p(x'|C_i)P(C_i)) = \log p(x'|C_i) + \log P(C_i)$ is maximized
- To be able to make prediction given all possible x'
 - We estimate the prior $P(C_i)$ by $\hat{P}[C_i] = \frac{\sum_{t=1}^N r_i^{(t)}}{N}$
 - By assuming that instances of the same class are normally distributed, we estimate the likelihood $p(x|C_i)$ by $\hat{p}(x|C_i) = \frac{1}{\sqrt{2\pi s_i^2}} \exp\left(\frac{-(x-m_i)^2}{2s_i^2}\right)$, where $m_i = \frac{\sum_{t=1}^N x^{(t)} r_i^{(t)}}{\sum_{t=1}^N r_i^{(t)}}$ and $s_i^2 = \frac{\sum_{t=1}^N (x^{(t)} - m_i)^2 r_i^{(t)}}{\sum_{t=1}^N r_i^{(t)} - 1}$

Outline

1 More About Probabilistic Modeling

2 MAP and Bayesian Estimation

3 The Bias/Variance Dilemma

4 Generative Methods

- Univariate Classification
- **Maximum Likelihood Estimation**
- Multivariate Classification
- Tuning the Model Complexity

Maximum Likelihood Estimation

- Why $\hat{P}(C_i) = \frac{\sum_{t=1}^N r_i^{(t)}}{N}$ and $\hat{p}(x|C_i) = \frac{1}{\sqrt{2\pi s_i^2}} \exp\left(\frac{-(x-m_i)^2}{2s_i^2}\right)$ are good choices?

Maximum Likelihood Estimation

- Why $\hat{P}(C_i) = \frac{\sum_{t=1}^N r_i^{(t)}}{N}$ and $\hat{p}(x|C_i) = \frac{1}{\sqrt{2\pi s_i^2}} \exp\left(\frac{-(x-m_i)^2}{2s_i^2}\right)$ are good choices?
 - It turns out that each of these estimators maximizes the likelihood $p(\mathcal{X}|\theta)$, where θ is the parameters of the distribution used to model the target probability ($P(C_i)$ and $p(x|C_i)$ respectively)
- When we talk about the likelihood-based classification, ***the “likelihood” actually refers to the one ($p(\mathcal{X}|\theta)$) of θ given \mathcal{X} rather than that ($p(x'|C_i)$) of C_i given x'***

ML Estimation for $P(C_i)$ (1/2)

- To estimate $P(C_i)$, we first assume that $P(C_i)$ has the Bernoulli distribution parametrized by $\theta = \rho_i$ and can be written as $P(C_i) = P(C_i; \theta)$
 - Let X_i be a random variable where $X_i = 1$ if the event “the outcome of a toss is C_i and $X_i = 0$ if “the outcome is not C_i ”
 - Let ρ_i be the probability that $X_i = 1$, we have $P(X_i = c; \theta) = \rho_i^c (1 - \rho_i)^{1-c}$, $c \in \{0, 1\}$
- Now the problem estimating $P(C_i | \theta) = P(X_i = 1; \theta) = \rho_i$ can be reduced to estimating $\theta = \rho_i$

ML Estimation for $P(C_i)$ (2/2)

- Given the training set \mathcal{X} , a good estimate of θ is the one that maximizes $P(\theta|\mathcal{X})$
 - From Bayes' rule, we can instead pick $\hat{\theta}$ maximizing $P(\mathcal{X}|\theta)$ if we don't have prior reason to favor certain θ
 - Equivalently, we pick $\hat{\theta}$ maximizing $\log P(\mathcal{X}|\theta)$
 - We have $\log P(\mathcal{X}|\theta) = \log \left(\prod_{t=1}^N \rho_i^{r_i^{(t)}} (1 - \rho_i)^{1-r_i^{(t)}} \right)$
 - Solving $\frac{d(\log P(\mathcal{X}|\theta))}{d\rho} = 0$ we obtain the **Maximum Likelihood (ML) estimator** $\hat{\rho}_i = \frac{\sum_{t=1}^N r_i^{(t)}}{N}$ [Proof]
- $\hat{P}[C_i] = P(C_i|\hat{\theta}) = \hat{\rho}_i$
- Note we can also consider all classes together and assume that $P(C_i)$ follows the Multinomial distribution parametrized by $\theta = (\rho_1, \dots, \rho_K)$ with constrains $\sum_{i=1}^K \rho_i = 1$
 - The ML estimator for each ρ_i will be the same as the above [Homework]

ML Estimation for $p(x|C_i)$ (1/2)

- We assume that $p(x|C_i)$ is normal and can be written as $p(x|C_i) = p(x|C_i; \theta)$ with some $\theta = (\mu_i, \sigma_i)$
 - $p(x|C_i; \theta) = p_{N_{\mu_i, \sigma_i^2}}(x) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(\frac{-(x-\mu_i)^2}{2\sigma_i^2}\right)$
- Now the problem estimating $p(x|C_i; \theta)$ can be reduced to estimating $\theta = (\mu_i, \sigma_i)$
- Given the training set \mathcal{X} , a good estimate of θ is the one that maximizes $\log p(\mathcal{X}|\theta)$

- We have $\log p(\mathcal{X}|\theta) = \log \left(\prod_{t=1}^N \left(\frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(\frac{-(x^{(t)}-\mu_i)^2}{2\sigma_i^2}\right) \right)^{r_i^{(t)}} \right)$
- Taking the partial derivatives of $\log p(\mathcal{X}|\theta)$ in terms of μ_i and σ_i and setting them equal to 0 we obtain the estimators $m_i = \frac{\sum_{t=1}^N x^{(t)} r_i^{(t)}}{\sum_{t=1}^N r_i^{(t)}}$ and

$$s_i^2 = \frac{\sum_{t=1}^N (x^{(t)} - m_i)^2 r_i^{(t)}}{\sum_{t=1}^N r_i^{(t)}} \text{ respectively [Proof]}$$

ML Estimation for $p(x|C_i)$ (2/2)

- Recall that $s_i^2 = \frac{\sum_{t=1}^N (x^{(t)} - m_i)^2 r_i^{(t)}}{\sum_{t=1}^N r_i^{(t)}}$ is a bias estimator, we can replace the denominator with $\sum_{t=1}^N r_i^{(t)} - 1$
 - This step is optional
 - The difference, actually, is negligible when N is large
- $\hat{p}(x|C_i) = p(x|C_i, \hat{\theta}) = \frac{1}{\sqrt{2\pi s_i^2}} \exp\left(\frac{-(x - m_i)^2}{2s_i^2}\right)$

Outline

1 More About Probabilistic Modeling

2 MAP and Bayesian Estimation

3 The Bias/Variance Dilemma

4 Generative Methods

- Univariate Classification
- Maximum Likelihood Estimation
- **Multivariate Classification**
- Tuning the Model Complexity

Multivariate Data

- Let's go back to a higher dimensional feature space
 - We are given a training set $\mathcal{X} = \{\mathbf{x}^{(t)}, \mathbf{r}^{(t)}\}_{t=1}^N$ where $\mathbf{x}^{(t)} \in \mathbb{R}^d$ and $(\mathbf{x}^{(t)}, \mathbf{r}^{(t)})$ are i.i.d. samples drawn from some unknown (multivariate) distribution
 - Typically, the features of $\mathbf{x}^{(t)}$ are correlated (otherwise we can discuss each attribute individually using the univariate methods)
- It might be a good idea to review the multivariate distributions now

Multivariate Classification

- The idea remains the same: given a new instance $\mathbf{x}' \in \mathbb{R}^d$, we make prediction by picking the class C_i if its discriminant $f_i(\mathbf{x}') = P(C_i|\mathbf{x}')$ is maximized
 - Generative assumption: pick C_i if $f_i(\mathbf{x}') = \log p(\mathbf{x}'|C_i) + \log P(C_i)$ is maximized

- It's common to assume that $p(\mathbf{x}|C_i)$ follows the multivariate normal distribution, i.e.,

$$p(\mathbf{x}|C_i) = \mathcal{PN}_{\mu_i, \Sigma_i}(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \det(\Sigma_i)^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu_i)^\top \Sigma_i^{-1} (\mathbf{x} - \mu_i) \right]$$

- Why?
 - Major reason: analytical simplicity
 - Studies also show that the model is robust to datasets departing from normality

Maximum Likelihood Estimation

- The ML estimators of $P(C_i)$ is $\hat{P}[C_i] = \sum_{t=1}^N r_i^{(t)} / N$

- We have seen this in the univariate cases before

- The ML estimators of $p(\mathbf{x}|C_i)$ is

$$\frac{1}{(2\pi)^{d/2} \det(\mathbf{S}_i)^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mathbf{m}_i)^\top \mathbf{S}_i^{-1} (\mathbf{x} - \mathbf{m}_i) \right], \text{ where}$$

$$\mathbf{m}_i = \frac{\sum_{t=1}^N \mathbf{x}^{(t)} r_i^{(t)}}{\sum_{t=1}^N r_i^{(t)}} \text{ and}$$

$$\mathbf{S}_i = \frac{\sum_{t=1}^N (\mathbf{x}^{(t)} - \mathbf{m}_i)(\mathbf{x}^{(t)} - \mathbf{m}_i)^\top r_i^{(t)}}{\sum_{t=1}^N r_i^{(t)}}$$

- Why?

Maximum Likelihood Estimation

- The ML estimators of $P(C_i)$ is $\hat{P}[C_i] = \sum_{t=1}^N r_i^{(t)} / N$

- We have seen this in the univariate cases before

- The ML estimators of $p(\mathbf{x}|C_i)$ is

$\frac{1}{(2\pi)^{d/2} \det(\mathbf{S}_i)^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mathbf{m}_i)^\top \mathbf{S}_i^{-1} (\mathbf{x} - \mathbf{m}_i) \right]$, where

$$\mathbf{m}_i = \frac{\sum_{t=1}^N \mathbf{x}^{(t)} r_i^{(t)}}{\sum_{t=1}^N r_i^{(t)}} \text{ and}$$

$$\mathbf{S}_i = \frac{\sum_{t=1}^N (\mathbf{x}^{(t)} - \mathbf{m}_i)(\mathbf{x}^{(t)} - \mathbf{m}_i)^\top r_i^{(t)}}{\sum_{t=1}^N r_i^{(t)}}$$

- Why? It's a good idea to review the matrix calculus now

ML Estimator of μ_j

- Let $\theta = (\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$, we have the likelihood $\log p(\mathcal{X}|\theta) =$
$$\log \left(\prod_{t=1}^N \left(\frac{1}{(2\pi)^{d/2} \det(\boldsymbol{\Sigma}_j)^{1/2}} e^{-\frac{1}{2}(\mathbf{x}^{(t)} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_j)} \right)^{r_j^{(t)}} \right) =$$

$$-\frac{N_j d}{2} \log(2\pi) - \frac{N_j}{2} \log(\det(\boldsymbol{\Sigma}_j)) - \frac{1}{2} \sum_{t=1}^N r_j^{(t)} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_j),$$
 where $N_j = \sum_{t=1}^N r_j^{(t)}$
- Recall that for any $\mathbf{a} \in \mathbb{R}^n$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$,
 - $\frac{\partial}{\partial \mathbf{x}} (\mathbf{a}^\top \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^\top \mathbf{a}) = \mathbf{a}^\top$
 - $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) = \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top)$
- Taking the partial derivative of $\log p(\mathcal{X}|\theta)$ with respect to $\boldsymbol{\mu}_j$ and setting it to zero, we get $\sum_{t=1}^N r_j^{(t)} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1} = \mathbf{0}^\top$ [Proof]
- So $\mathbf{m}_j = \frac{\sum_{t=1}^N \mathbf{x}^{(t)} r_j^{(t)}}{\sum_{t=1}^N r_j^{(t)}}$

ML Estimator of Σ_i (1/2)

- $\log p(\mathcal{X}|\theta) = -\frac{N_i d}{2} \log(2\pi) - \frac{N_i}{2} \log(\det(\Sigma_i)) - \frac{1}{2} \sum_{t=1}^N r_i^{(t)} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_i)$
- Note $\log(\det(\Sigma_i^{-1})) = -\log(\det(\Sigma_i))$
- Also, $(\mathbf{x}^{(t)} - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_i) = \text{tr}(\Sigma_i^{-1} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_i)(\mathbf{x}^{(t)} - \boldsymbol{\mu}_i)^\top)$
[Proof]
- We can rewrite the likelihood as $\log p(\mathcal{X}|\theta) = -\frac{N_i d}{2} \log(2\pi) + \frac{N_i}{2} \log(\det(\Sigma_i^{-1})) - \frac{1}{2} \sum_{t=1}^N r_i^{(t)} \text{tr}(\Sigma_i^{-1} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_i)(\mathbf{x}^{(t)} - \boldsymbol{\mu}_i)^\top)$

ML Estimator of Σ_j (2/2)

- Given any function $f(x)$, let $g(x) = f(\frac{1}{x})$ for any $x > 0$, then x^* is a stationary point of g iff $\frac{1}{x^*}$ is a stationary point of f
 - The matrix version $g(\mathbf{A}) = f(\mathbf{A}^{-1})$ applies when \mathbf{A} is positive definite
- We can seek for the partial derivative of $\log p(\mathcal{X}|\theta)$ with respect to Σ_j^{-1}
- Recall that $\frac{\partial}{\partial \mathbf{A}} \ln(\det(\mathbf{A})) = (\mathbf{A}^{-1})^\top$, and $\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{A}\mathbf{B}) = \mathbf{B}^\top$
- Taking the partial derivative of $\log p(\mathcal{X}|\theta)$ with respect to Σ_j^{-1} and setting it to zero, we get
$$\frac{N_j}{2} \Sigma_j - \frac{1}{2} \sum_{t=1}^N r_j^{(t)} (\mathbf{x}^{(t)} - \boldsymbol{\mu}_j)(\mathbf{x}^{(t)} - \boldsymbol{\mu}_j)^\top = \mathbf{O} \text{ [Proof]}$$
- Therefore, $\mathbf{S}_j = \frac{\sum_{t=1}^N (\mathbf{x}^{(t)} - \mathbf{m}_j)(\mathbf{x}^{(t)} - \mathbf{m}_j)^\top r_j^{(t)}}{\sum_{t=1}^N r_j^{(t)}}$

Quadratic Discrimination

- Ignoring the constant terms we have the discriminant $f_i(\mathbf{x}) = -\frac{1}{2} \log(\det(\mathbf{S}_i)) - \frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^\top \mathbf{S}_i^{-1}(\mathbf{x} - \mathbf{m}_i) + \log \hat{P}[C_i]$, which can be rewritten as $f_i(\mathbf{x}) = \mathbf{x}^\top \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^\top \mathbf{x} + w_i$, where $\mathbf{W}_i = -\frac{1}{2} \mathbf{S}_i^{-1}$, $\mathbf{w}_i = \mathbf{S}_i^{-1} \mathbf{m}_i$, and $w_i = -\frac{1}{2} \mathbf{m}_i^\top \mathbf{S}_i^{-1} \mathbf{m}_i - \frac{1}{2} \log(\det(\mathbf{S}_i)) + \log \hat{P}[C_i]$ [Proof]
 - The classification is done via **quadratic discrimination**
 - The decision boundary between any two classes is quadratic too [Proof]

Multivariate Classification (3/3)

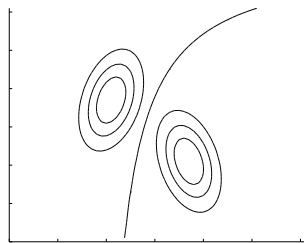
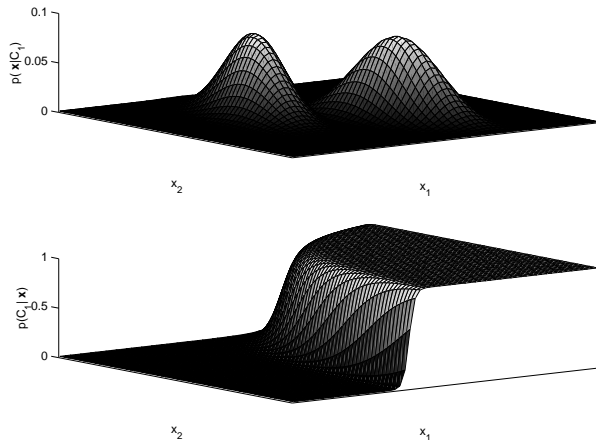


Figure : (a) The graphs of $p(\mathbf{x}|C_i)$ for two classes with different covariance matrices. (b) The graph of posterior $P(C_1|\mathbf{x})$. (c) Level sets of $p(\mathbf{x}|C_i)$ and the decision boundary.

1 More About Probabilistic Modeling

2 MAP and Bayesian Estimation

3 The Bias/Variance Dilemma

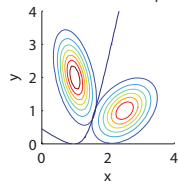
4 Generative Methods

- Univariate Classification
- Maximum Likelihood Estimation
- Multivariate Classification
- **Tuning the Model Complexity**

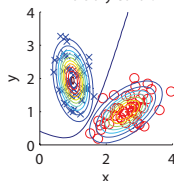
Simplifications (1/2)

- Quadratic discrimination:
 - Attributes in different classes have different covariance matrices \mathbf{S}_i ($\{x : p(x|C_i) = c\}$ are ellipsoids)
- **Linear discrimination:**
 - Attributes in different classes share the same correlation $\mathbf{S}_i = \mathbf{S}$ (ellipsoids with the same shape/orientation)
 - Attributes in each classes are independent $\mathbf{S}_i = \mathbf{S} = \mathbf{D}$ (axis-aligned ellipsoids with the same shape/orientation)
 - Attributes in each classes has the same variance $\mathbf{S}_i = \mathbf{S} = s^2 \mathbf{I}$ (equal-sized spheres)

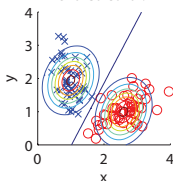
Population likelihoods and posteriors



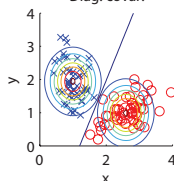
Arbitrary covar.



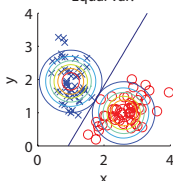
Shared covar.



Diag. covar.



Equal var.



Simplifications (2/2)

- Linear discrimination models seem to be oversimplified, but why are they popular in real applications?

Simplifications (2/2)

- Linear discrimination models seem to be oversimplified, but why are they popular in real applications?
- Quadratic discrimination has lower bias, but higher variance
- Experience tells us that *when we have a small dataset*, it may be better to assume a shared and simplified covariance matrix
 - $\mathbf{S}_i = \mathbf{S}$ can be estimated using all examples in a dataset together
 - $\mathbf{S} = \mathbf{D}$ if we do not have enough data to estimate the covariance between attributes accurately
 - $\mathbf{D} = s^2 \mathbf{I}$ if attributes are z-normalized
- Linear discrimination is *not* necessarily linear
 - We can augment the inputs (e.g., $x_{d+1} = \exp(x_1 + x_4)$) to build a higher dimensional feature space, if we believe this is useful
 - Linear discrimination in the augmented feature space corresponds to a nonlinear model in the original input space
- We can perform the cross validation to decide which assumption is the best

Linear Discrimination ($\mathbf{S}_i = \mathbf{S}$)

- The discriminant for each class is $f_i(\mathbf{x}) = \mathbf{x}^\top \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^\top \mathbf{x} + w_i$, where
$$\mathbf{W}_i = -\frac{1}{2} \mathbf{S}_i^{-1},$$
$$\mathbf{w}_i = \mathbf{S}_i^{-1} \mathbf{m}_i, \text{ and}$$
$$w_i = -\frac{1}{2} \mathbf{m}_i^\top \mathbf{S}_i^{-1} \mathbf{m}_i - \frac{1}{2} \log(\det(\mathbf{S}_i)) + \log \hat{P}[C_i]$$
- We can replace \mathbf{S}_i with \mathbf{S} , the estimator of Σ of all instances in the training set
 - The level sets $\{\mathbf{x} : p(\mathbf{x}|C_i) = c\}$ are ellipsoids with the same shape/orientation
- Ignoring the constant terms, the discriminant now becomes
$$f_i(\mathbf{x}) = \mathbf{w}_i^\top \mathbf{x} + w_i, \text{ where}$$
$$\mathbf{w}_i = \mathbf{S}^{-1} \mathbf{m}_i \text{ and}$$
$$w_i = -\frac{1}{2} \mathbf{m}_i^\top \mathbf{S}^{-1} \mathbf{m}_i + \log \hat{P}[C_i] \text{ [Proof]}$$

Naive Bayes Classifiers ($\mathbf{S}_i = \mathbf{S} = \mathbf{D}$)

- We can further assume that attributes are independent with each

other, i.e., $\mathbf{S}_i = \mathbf{S} = \begin{bmatrix} s_0^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & s_d^2 \end{bmatrix}$ are diagonal

- Likelihood-based classifiers using this strong (naive) independence assumption are called the **naive Bayes' classifiers**
- The level sets $\{\mathbf{x} : p(\mathbf{x}|C_i) = c\}$ are axis-aligned ellipsoids
- $f_i(\mathbf{x}) = -\frac{1}{2} \sum_{j=1}^d \left(\frac{m_{i,j}^2 - 2x_j m_{i,j}}{s_j^2} \right) + \log \hat{P}[C_i]$ [Proof]
- If we further assume that that attributes have the same variance, i.e., $\mathbf{S}_i = \mathbf{S} = sI$
 - The level sets $\{\mathbf{x} : p(\mathbf{x}|C_i) = c\}$ degenerate into spheres
 - $f_i(\mathbf{x}) = -\frac{1}{2s^2} \left(\|\mathbf{m}_i\|^2 - 2\mathbf{m}_i^T \mathbf{x} \right) + \log \hat{P}[C_i]$ (or $f_i(\mathbf{x}) = -\frac{1}{2s^2} \|\mathbf{x} - \mathbf{m}_i\|^2 + \log \hat{P}[C_i]$) [Proof]
 - If we drop $\log \hat{P}[C_i]$, we obtain a **nearest mean classifier**