

CSE555: Introduction to Pattern Recognition
Spring, 2007

Mid-Term Exam (with solutions)
(100 points, Closed book/notes)

The last page contains some formulas that might be useful.

1. Part(i) (10 pts)

Suppose a bank classifies customers as either good or bad credit risks. On the basis of extensive historical data, the bank has observed that 1% of good credit risks and 10% of bad credit risks overdraw their account in any given month. A new customer opens a checking account at this bank. On the basis of a check with a credit bureau, the bank believe that there is a 70% chance the customer will turn out to be a good credit risk.

- (a) (5 pts) Suppose that this customer's account is overdrawn in the first month. How does this alter the bank's opinion of this customer's creditworthiness?

Answer:

Let G and O represent the following events:

G : customer is considered to be a good credit risk.

O : customer overdraws checking account.

From the bank's historical data, we have $P(O|G) = 0.01$ and $P(O|\tilde{G}) = 0.1$, and we also know that the bank's initial opinion about the customer's creditworthiness is given by $P(G) = 0.7$.

Given the information that the customer has an overdraft in the first month, the bank's revised opinion about the customer's creditworthiness is given by the conditional probability $P(G|O)$. Using Bayes' theorem and the law of total probability,

$$\begin{aligned} P(G|O) &= \frac{P(O|G)P(G)}{P(O)} \\ &= \frac{P(O|G)P(G)}{P(O|G)P(G) + P(O|\tilde{G})P(\tilde{G})} \\ &= \frac{0.01 \times 0.7}{0.01 \times 0.7 + 0.1 \times 0.3} \\ &\approx 0.189 \end{aligned}$$

- (b) (5 pts) Given (a), what would be the bank's opinion of the customer's creditworthiness at the end of the second month if there was not an overdraft in the second month?

Answer:

According to (a), at the beginning of the second month, the prior probability of the customer being a good risk is 0.189. Since the customer does not have an overdraft

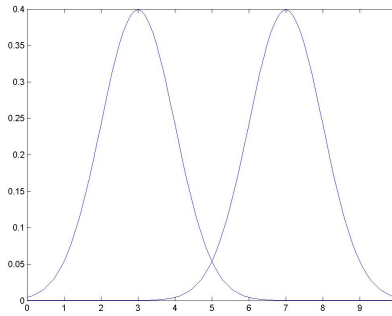
in the second month, the bank's opinion about the customer's creditworthiness at the end of the month is given by

$$\begin{aligned}
 P(G|\tilde{O}) &= \frac{P(\tilde{O}|G)P(G)}{P(\tilde{O})} \\
 &= \frac{P(\tilde{O}|G)P(G)}{P(\tilde{O}|G)P(G) + P(\tilde{O}|\tilde{G})P(\tilde{G})} \\
 &= \frac{0.99 \times 0.189}{0.99 \times 0.189 + 0.9 \times 0.811} \\
 &\approx 0.204
 \end{aligned}$$

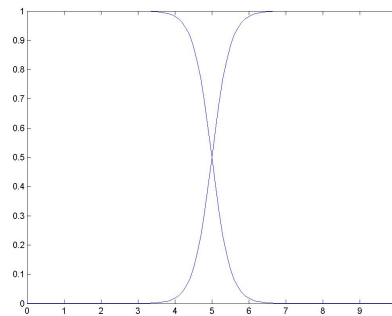
Part(ii) (10 pts)

Consider a two-category classification problem with one-dimensional Gaussian distributions $p(x|w_i) \sim N(\mu_i, \sigma^2), i = 1, 2$ (i.e. they have same variance but different means).

(a) (3 pts) Sketch the two densities $p(x|w_1)$ and $p(x|w_2)$ in one figure. **Answer:**



(b) (7 pts) Sketch the two posterior probabilities $P(w_1|x)$ and $P(w_2|x)$ in one figure, assuming same prior probabilities. **Answer:**



Part(iii) (5 pts)

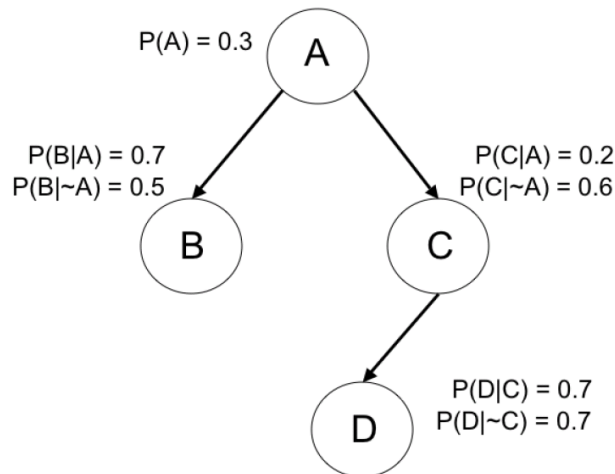
Suppose there is a two-class one-dimensional problem with the Gaussian distributions: $p(x|w_1) \sim N(-1, 1)$, and $p(x|w_2) \sim N(4, 1)$ and equal prior probabilities. In order to find a good classifier, you are given as much training data as you would like and you are free to pick any method. What is the best error on test data that any learning algorithm can attain and why? (Try to keep your reasoning to two sentences or fewer.)

- (a) 0% error
- (b) more than 0% but less than 10%
- (c) more than 20%
- (d) can't tell from this information

Answer: b

The best error can be achieved by using Bayes decision rule. The Bayes decision boundary is $x = 1.5$, which means that $|x - \mu| > 2\sigma$. Therefore the error rate should be less than 5%.

2. (15 pts) Compute the following probabilities from the Bayes network below.



(a) $P(A, B, C, D)$

Answer:

$$\begin{aligned}
 P(A, B, C, D) &= P(A)P(B|A)P(C|A)P(D|C) \\
 &= 0.3 \times 0.7 \times 0.2 \times 0.7 \\
 &= 0.0294
 \end{aligned}$$

(b) $P(A|B)$

Answer:

$$\begin{aligned}P(A|B) &= \frac{P(B|A)P(A)}{P(B)} \\&= \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|\tilde{A})P(\tilde{A})} \\&= \frac{0.7 \times 0.3}{0.7 \times 0.3 + 0.5 \times 0.7} \\&= 0.375\end{aligned}$$

(c) $P(C|B)$

Answer:

$$\begin{aligned}P(C|B) &= P(C|A)P(A|B) + P(C|\tilde{A})P(\tilde{A}|B) \\&= 0.2 \times 0.375 + 0.6 \times 0.625 \\&= 0.45\end{aligned}$$

(d) $P(B|D)$ (Hint: there is a shortcut for this one.)

Answer: since $P(D|C) = P(D|\tilde{C}) = 0.7$, D is independent of C and therefore D is independent from the network. Thus

$$P(B|D) = P(B) = P(B|A)P(A) + P(B|\tilde{A})P(\tilde{A}) = 0.56$$

3. (20 pts) Assume a two-class problem with equal a priori class probabilities and Gaussian class-conditional densities as follows:

$$p(x|\omega_1) = N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} a & c \\ c & b \end{bmatrix}\right) \quad p(x|\omega_2) = N\left(\begin{bmatrix} d \\ e \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

where $a \times b - c \times c = 1$.

(a) (7 pts) Find the equation of the decision boundary between these two classes in terms of the given parameters, after choosing a logarithmic discriminant function.

Answer:

Given two normal density and equal prior probabilities, we choose the discriminant function as

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mu_i)^t \Sigma_i^{-1}(\mathbf{x} - \mu_i) - \frac{1}{2} \ln |\Sigma_i|$$

then the decision boundary is given by $g_1(\mathbf{x}) = g_2(\mathbf{x})$, and after simplified the function, we have

$$(b-1)x_1^2 + (a-1)x_2^2 - 2cx_1x_2 + 2dx_1 + 2ex_2 - d^2 - e^2 = 0$$

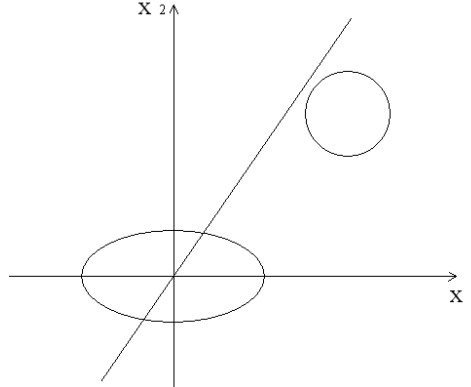
(b) (5 pts) Determine the constraints on the values of a, b, c, d and e , such that the resulting discriminant function results with a linear decision boundary.

Answer:

To get a linear decision boundary, we need $\Sigma_1 = \Sigma_2$. Thus we have $a = b = 1$ and $c = 0$.

(c) (8 pts) Let $a = 2, b = 1, c = 0, d = 4, e = 4$. Draw typical equal probability contours for both densities and determine the projection line for the Fisher's discriminant, such that these two classes are separated in an optimal way.

Answer:



$$\mathbf{S}_w = \mathbf{S}_1 + \mathbf{S}_2 = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix}$$

therefore

$$\begin{aligned} \mathbf{w} &= \mathbf{S}_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2) \\ &= \frac{1}{6} \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} -4 \\ -4 \end{bmatrix} \\ &= \left(-\frac{2}{3}\right) \begin{bmatrix} 2 \\ 3 \end{bmatrix} \end{aligned}$$

4. (i) (5 pts) When is maximum a posteriori (MAP) estimation equivalent to maximum-likelihood (ML) estimation? (Try to keep your answer to two sentences or fewer.)

Answer: MAP estimation is equivalent to ML estimation we when the prior probability is uniformly distributed over the parameters. Alternatively when there is an infinite amount of data.

- (ii) (15 pts) Suppose that n samples x_1, x_2, \dots, x_n are drawn independently according to the Erlang distribution with the following probability density function

$$p(x|\theta) = \theta^2 x e^{(-\theta x)} u(x)$$

where $u(x)$ is the unit step function as follows

$$u(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \end{cases}$$

find the maximum likelihood estimate of the parameter θ .

Answer:

The log-likelihood function is

$$l(\theta) = \sum_{k=1}^n \ln p(x_k | \theta) = \sum_{k=1}^n \ln[\theta^2 x_k e^{-\theta x_k}] = \sum_{k=1}^n (2 \ln \theta + \ln x_k - \theta x_k) \quad \text{for } x_k > 0$$

we solve $\nabla_{\theta} l(\theta) = 0$ to find $\hat{\theta}$

$$\nabla_{\theta} l(\theta) = \sum_{k=1}^n \nabla_{\theta} (2 \ln \theta + \ln x_k - \theta x_k) = \sum_{k=1}^n \left(\frac{2}{\theta} - x_k \right) = 0$$

thus we have

$$\hat{\theta} = \frac{2n}{\sum_{k=1}^n x_k}$$

5. (20 pts) Suppose we have a HMM with N hidden states and M visible states, and a particular sequence of visible states V^T is observed, where T is the length of the sequence.

(a) (13 pts) Describe a smart algorithm to compute the probability that V^T was generated by this model, i.e. using either *HMM Forward algorithm* or *HMM Backward algorithm*. You need to give the definition of $\alpha_j(t)$ or $\beta_j(t)$ and explain the meaning of it in words.

Answer:

To compute $P(V^T)$ recursively, we use HMM Forward algorithm by defining

$$\alpha_j(t) = \left[\sum_i \alpha_i(t-1) a_{ij} \right] b_{jk} v(t)$$

where a_{ij} and b_{jk} are the transition and emission probabilities.

Since $\alpha_j(t)$ is the probability that the HMM is in hidden state ω_i at step t having generated the first t elements of V^T , therefore $P(V^T)$ can be represented by $\alpha_0(T)$ at the final state (i.e. the absorber state with a unique null visible symbol v_0).

Check the textbook (DHS) for the complete algorithm.

(b) (3 pts) What is the computational complexity of the algorithm above and why?

Answer:

The computational complexity of this algorithm is $O(N^2T)$, because at each step t we need to compute $\alpha_j(t)$ for all $j = 1$ to N and each $\alpha_j(t)$ involves a summation over all $\alpha_i(t - 1)$.

(c) (4 pts) How to compute the probability above without using this smart algorithm? What is the computational complexity and why?

Answer:

In general, $P(V^T)$ can be computed by

$$P(V^T) = \sum_{r=1}^{R_{max}} P(V^T|\omega_r^T)P(\omega_r^T)$$

where $R_{max} = N^T$ for a HMM with N hidden states.

Therefore the computational complexity will be $O(N^T T)$.

Appendix: Useful formulas.

- For a 2×2 matrix,

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

the matrix inverse is

$$A^{-1} = \frac{1}{|A|} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

- The *scatter matrices* \mathbf{S}_i are defined as

$$\mathbf{S}_i = \sum_{X \in D_i} (X - \mathbf{m}_i)(X - \mathbf{m}_i)^t$$

where \mathbf{m}_i is the d -dimensional sample mean.

The *within-class scatter matrix* is defined as

$$\mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$$

The *between-class scatter matrix* is defined as

$$\mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^t$$

The solution for the \mathbf{w} that optimizes $J(\mathbf{w}) = \frac{\mathbf{w}^t \mathbf{S}_B \mathbf{w}}{\mathbf{w}^t \mathbf{S}_W \mathbf{w}}$ is

$$\mathbf{w} = \mathbf{S}_W^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$