

## Non-Linear Signal Processing: Quiz Exercise 6

C.M. Bishop: Pattern Recognition and Machine Learning,  
Sections 1.5, 4.2, 4.3.4, 5.1-5.4

### Questions

1. Consider the probability of class 1,  $C_1$ , given a data point,  $x$ ,  $P(C_1|x)$  in a binary classification setting. How would the probability  $P(C_1|x)$  change if the probability of the data point given class 2,  $P(x|C_2)$ , were to increase? Assume all other quantities  $P(C_1)$ ,  $P(C_2)$ , and  $P(x|C_1)$  remain the same.
  - (a)  $P(C_1|x)$  would not be affected.
  - (b)  $P(C_1|x)$  would increase.
  - (c)  $P(C_1|x)$  would decrease
  - (d)  $P(C_1|x)$  would decrease, but so would  $P(C_2|x)$  so that the ratio between the two would remain the same.
2. Do the error functions of logistic regression and multinomial regression have global, unique minima?
  - (a) logistic regression: yes, multinomial regression: yes
  - (b) logistic regression: yes, multinomial regression: no
  - (c) logistic regression: no, multinomial regression: yes
  - (d) logistic regression: no, multinomial regression: no
3. Which assumptions on class distributions lead to linear decision boundaries in classification problems?
  - (a) Class distributions are Gaussian
  - (b) Class distributions have the same covariance matrix
  - (c) Class distributions have different covariance matrices
  - (d) Class distributions have the same covariance matrix and are Gaussian
4. The standard form for linear equations with  $n$  variables is  $a_1x_1 + a_2x_2 + \dots + a_nx_n = b$ , where  $a_1, a_2, \dots, a_n$ , and  $b$  are constants. Let  $\mathbf{a} = (a_1, a_2, \dots, a_n)^T$  be the column vector containing the coefficients  $a_1, a_2, \dots, a_n$  of the variables. Then the standard form can be written  $\mathbf{a}^T \mathbf{x} = b$ , where  $\mathbf{x}$  is the column vector containing the variables.

Consider equation (4.65) in Bishop, which states  $P(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$ . This can be re-arranged to make it clear that it is a linear equation when  $P(C_1|\mathbf{x})$  is constant, as is the case on decision boundaries. Re-arranging and identifying the resulting terms with those in the standard form for linear equations gives

- (a)  $\mathbf{w}$  corresponds to  $\mathbf{a}$  in the standard form and  $w_0 + \sigma^{-1}(P(C_1|\mathbf{x}))$  to  $b$
  - (b)  $\mathbf{w}$  corresponds to  $\mathbf{a}$  in the standard form and  $\sigma^{-1}(P(C_1|\mathbf{x})) - w_0$  to  $b$
  - (c)  $\mathbf{x}$  corresponds to  $\mathbf{a}$  in the standard form and  $w_0 + \sigma^{-1}(P(C_1|\mathbf{x}))$  to  $b$
  - (d)  $\mathbf{x}$  corresponds to  $\mathbf{a}$  in the standard form and  $\sigma^{-1}(P(C_1|\mathbf{x})) - w_0$  to  $b$
5. Consider the cancer treatment loss function in Figure 1.25 in Bishop. The hospital has a probabilistic classifier  $p(C_k|\mathbf{x})$  which takes the biomarker measurement  $\mathbf{x}$  and predicts probability for  $C_1 = \text{cancer}$  and  $C_2 = \text{normal}$ . For a given patient the classifier gives  $p(\text{cancer}|\mathbf{x}) = 1/1000$ . What should we do?
- (a) Treat the patient as if the patient is normal because that has the lowest associated expected loss.
  - (b) Treat the patient as if the patient has cancer because that has the lowest associated expected loss.
  - (c) The expected losses of normal and cancer are so close that we cannot make a decision.
  - (d) We cannot use the classifier to make decisions because its predictions are not 100 % certain.
6. In a neural network for binary classification the output

$$y(\mathbf{x}|\mathbf{w}) = \frac{1}{1 + \exp(-a(\mathbf{x}; \mathbf{w}))} = \sigma(a(\mathbf{x}; \mathbf{w}))$$

gives the probability for class 1. What is  $a(\mathbf{x}; \mathbf{w})$ ?

- (a)  $\mathbf{w} \cdot \mathbf{x}$
- (b)  $\sigma(\mathbf{w}^{(2)} \cdot \mathbf{z})$
- (c)  $\tanh(\mathbf{w}^{(1)} \cdot \mathbf{x})$
- (d)  $\mathbf{w}^{(2)} \cdot \mathbf{z},$

where  $\mathbf{z}$  is shorthand for the output of the hidden unit.

7. Which of the following statements are correct?
- I) A discriminative function can be used to simulate data.
  - II) Posterior probabilities of class membership are found both when using discriminative models and generative models.
  - III) Discriminant functions map explanatory variables directly onto class labels without using probabilities.
- (a) I)
  - (b) I) and II)
  - (c) II) and III)
  - (d) I), II), and III)

8. When using the squared loss as cost function in regression problems, what is the best prediction?
  - (a) The conditional mean, conditioning on target variables used during training
  - (b) The conditional mean, conditioning on explanatory variables used during training
  - (c) The conditional mean, conditioning on the observed explanatory variables for which a prediction is desired
  - (d) The conditional mean, conditioning on the observed explanatory variables for which a prediction is desired, plus the intrinsic variance of the data
9. When using the squared loss as cost function in regression problems, what is the minimal attainable error?
  - (a) Zero
  - (b) The precision of the noise on the target variables
  - (c) The standard deviation of the noise on the target variables
  - (d) The variance of the noise on the target variables
10. In multi-class logistic regression, we have that  $P(C_k|\phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$ , where  $a_k = \mathbf{w}_k^T \phi$ , where  $\mathbf{w}_k$  is the vector of coefficients and  $\phi$  the vector of basis functions of explanatory variables. What is  $\frac{\partial y_k}{\partial a_j}$ ?
  - (a)  $y_k(1 - y_k)$  when  $k = j$  and  $-y_k y_j$  when  $k \neq j$
  - (b)  $\exp(a_k)\phi$  when  $k = j$  and  $\exp(a_j)\phi$  when  $k \neq j$
  - (c)  $\exp(a_k)\mathbf{w}$  when  $k = j$  and  $\exp(a_j)\mathbf{w}$  when  $k \neq j$
  - (d)  $y_k^2(1 - y_k)$  when  $k = j$  and  $-y_k^2 y_j$  when  $k \neq j$

## Hint - reading for each question

1. Section 4.2
2. Section 4.3.4
3. Section 4.2.1
4. Section 4.2.1
5. Section 1.5.2
6. Section 5.1
7. Section 1.5.4
8. Section 1.5.5
9. Section 1.5.5
10. Section 4.3.4

## Correct answers

1. (c) Using Bayes theorem as given in equation (4.57) in the textbook, we write up the probability of interest:

$$P(C_1|\mathbf{x}) = \frac{P(\mathbf{x}|C_1)P(C_1)}{P(\mathbf{x}|C_1)P(C_1) + P(\mathbf{x}|C_2)P(C_2)}.$$

If  $P(\mathbf{x}|C_2)$  were to increase, the numerator would not be affected while the denominator would increase. Hence  $P(C_1|\mathbf{x})$  would decrease. Intuitively, this could be understood by noting that when the probability of that particular data point being observed in class two increases, seeing that data point does not point as strongly to class one being the generating class as before. Hence the probability that class one was the generating class decreases.

2. (a) Since logistic regression can be seen as multinomial regression (another term for multi-class logistic regression) with only two classes, the answer for logistic regression will be the same as that for multinomial regression. As stated at the end of section 4.3.4 in the textbook, the error function for multinomial regression has a unique minimum.
3. (d) When deriving the posterior probability for each class given a data point as a function of the data point  $x$ , a second order equation in  $x$  is obtained when the class distributions are assumed Gaussian. When the covariance matrices for all classes are equal, all quadratic terms in  $x$  are cancelled. Hence linear decision boundaries are obtained.
4. (b) To get  $P(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$  in the normal form, we isolate  $\mathbf{w}^T \mathbf{x}$  in the following way:

$$\begin{aligned} P(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) &\Leftrightarrow \sigma(P(C_1|\mathbf{x}))^{-1} = \mathbf{w}^T \mathbf{x} + w_0 \\ &\Leftrightarrow \mathbf{w}^T \mathbf{x} = \sigma(P(C_1|\mathbf{x}))^{-1} - w_0 \end{aligned}$$

When  $P(C_1|\mathbf{x})$  is constant, the right-hand side is clearly constant, showing  $P(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$  is linear in  $\mathbf{x}$  for constant  $P(C_1|\mathbf{x})$ , i.e. that decision boundaries are linear.

5. (b) From Figure 1.25 in Bishop, we find that the losses for correct classification are zero, whereas the loss for treating the patient as having cancer if the patient is in fact normal is  $L(cancer|normal) = 1$  and the loss for treating the patient as normal if the patient does in fact have cancer is  $L(normal|cancer) = 1000$ . The expected loss if we treat the patient as having cancer is  $p(cancer|\mathbf{x})L(cancer|normal) + p(cancer|\mathbf{x})L(cancer|cancer) = 1/1000 \cdot 1 + 1/1000 \cdot 0 = 1/1000$ . The expected loss if we treat the patient as being normal is  $p(normal|\mathbf{x})L(normal|normal) + p(normal|\mathbf{x})L(normal|cancer) = (1 - 1/1000) \cdot 0 + (1 - 1/1000) \cdot 1000 = 999$ . Hence the expected loss is least if we treat the patient as having cancer.
6. (d) Since the output is  $\sigma(a(\mathbf{x}; \mathbf{w}))$ , the input to the output layer must be  $a(\mathbf{x}; \mathbf{w})$ . Since inputs from the hidden layer are linearly combined inside the activation function in the output layer using the weights from the hidden layer to the output layer,  $\mathbf{w}^{(2)}$ , this must be the operation represented by  $a(\mathbf{x}; \mathbf{w})$ . That is,  $a(\mathbf{x}; \mathbf{w}) = \mathbf{w}^{(2)} \cdot \mathbf{z}$ .

7. (c) A discriminative function directly gives a prediction given the explanatory variables. Discriminative models directly give the posterior probabilities, while generative models calculate the posterior probability using Bayes' theorem. Since discriminative functions do not have an associated generating model, data cannot be simulated based on discriminative functions.
8. (c) When a prediction is desired given some explanatory variables, those explanatory variables are used to calculate the function value of the fitted regression function. This gives the expectation of the new target variable conditioned on the observed explanatory variables.
9. (d) If the fitted regression function is equal to the true function  $y(x)$ , then the first term in equation (1.90) in the textbook becomes zero and the second term becomes equal to the variance of the noise on the target variables.
10. (a) We have  $P(C_k|\phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$ . The definition of  $a_k$  is irrelevant. To find the derivative, we use the quotient rule  $\frac{d \frac{g(x)}{h(x)}}{dx} = \frac{h(x) \frac{dg(x)}{dx} - g(x) \frac{dh(x)}{dx}}{(h(x))^2}$ . We get:

$$\begin{aligned} \frac{\partial y_k}{\partial a_j} &= \frac{\partial \frac{\exp(a_k)}{\sum_i \exp(a_i)}}{\partial a_j} \\ &= \frac{(\sum_i \exp(a_i)) \frac{\partial \exp(a_k)}{\partial a_j} - \exp(a_k) \frac{\partial \sum_i \exp(a_i)}{\partial a_j}}{(\sum_i \exp(a_i))^2}. \end{aligned}$$

For  $k \neq j$ , this is:

$$\begin{aligned} \frac{(\sum_i \exp(a_i)) \frac{\partial \exp(a_k)}{\partial a_j} - \exp(a_k) \frac{\partial \sum_i \exp(a_i)}{\partial a_j}}{(\sum_i \exp(a_i))^2} &= \frac{(\sum_i \exp(a_i)) \cdot 0 - \exp(a_k) \exp(a_j)}{(\sum_i \exp(a_i))^2} \\ &= -\frac{\exp(a_k)}{\sum_i \exp(a_i)} \frac{\exp(a_j)}{\sum_i \exp(a_i)} = -y_k y_j \end{aligned}$$

For  $k = j$ , we get:

$$\begin{aligned} \frac{(\sum_i \exp(a_i)) \frac{\partial \exp(a_k)}{\partial a_j} - \exp(a_k) \frac{\partial \sum_i \exp(a_i)}{\partial a_j}}{(\sum_i \exp(a_i))^2} &= \frac{(\sum_i \exp(a_i)) \cdot \exp(a_k) - \exp(a_k) \exp(a_k)}{(\sum_i \exp(a_i))^2} \\ &= \frac{\exp(a_k)}{\sum_i \exp(a_i)} \cdot \left( \frac{\sum_i \exp(a_i) - \exp(a_k)}{\sum_i \exp(a_i)} \right) \\ &= y_k \cdot \left( 1 - \frac{\exp(a_k)}{\sum_i \exp(a_i)} \right) = y_k(1 - y_k). \end{aligned}$$