

# Midterm 1 Practice Questions

## 1. Biased Coins

You have a box full of coins. There are two types of coins,  $\mathcal{C}_1$  and  $\mathcal{C}_2$ . Coins of type  $\mathcal{C}_1$  come up heads with probability 0.8 and coin of type  $\mathcal{C}_2$  come up heads with probability 0.2. There are many more  $\mathcal{C}_1$  coins in the box than  $\mathcal{C}_2$  coins, in fact 90% of the coins are of type  $\mathcal{C}_1$ . You grab a coin at random from inside the box and flip it 10 times, getting five heads and five tails. Compute  $p(D | \mathcal{C}_1)$ ,  $p(D | \mathcal{C}_2)$ . How probable is it that you have a coin of type  $\mathcal{C}_1$ , given these ten flips?

## 2. Redundant Features in Naïve Bayes

Suppose that we use a Naïve Bayes classifier to classify binary data with binary feature vectors  $\mathbf{x}_n \in \{0, 1\}^D$ . We'll classify them into two classes,  $\mathcal{C}_1$  and  $\mathcal{C}_2$ . With Naïve Bayes and binary features, the class conditional distributions will be of the form of a product of Bernoulli distributions:

$$p(\mathbf{x} | \mathcal{C}_k) = \prod_{d=1}^D \mu_{kd}^{x_d} (1 - \mu_{kd})^{(1-x_d)},$$

where  $x_d \in \{0, 1\}$ , and  $\mu_{kd} = p(x_d = 1 | \mathcal{C}_k)$ . Assume also that the class priors are uniform, i.e.,  $p(\mathcal{C}_1) = p(\mathcal{C}_2) = \frac{1}{2}$ .

- (a) If  $D = 1$  (i.e., there is only one feature), use the equations above to write out  $\ln \frac{p(\mathcal{C}_1 | x)}{p(\mathcal{C}_2 | x)}$  for a single binary feature  $x$ .
- (b) Now suppose we change our feature representation so that instead of using just a single feature, we use two redundant features (i.e., two features that always have the same value so that  $x_1 = x_2$ ). Since they are the same, you can assume that  $\mu_{k1} = \mu_{k2}$  also. With this feature representation, let's write  $\hat{x} = x_1 \cdot x_2$ , since there can only be two configurations of the  $x_1, x_2$  pair, instead of four. What is  $\ln \frac{p(\mathcal{C}_1 | \hat{x})}{p(\mathcal{C}_2 | \hat{x})}$  in terms of the value for  $\ln \frac{p(\mathcal{C}_1 | x)}{p(\mathcal{C}_2 | x)}$  you calculated in part (a)?
- (c) Does this seem like a bug or a feature? Why?

### 3. Binomial Regression

You've been hired by a startup to build a ratings system for restaurants. Users rate the restaurants on a scale of 0 to 10 (i.e.,  $t_n \in \{0, 1, \dots, 10\}$ ) and you have a set of real-valued features for each restaurant,  $\mathbf{x}_n \in \mathbb{R}^D$ . Given the range of the  $t_n$ , it seems like a binomial distribution would be a good choice for building a regression model:

$$p(k|\rho) = \binom{10}{k} \rho^k (1-\rho)^{10-k},$$

where  $\rho$  parameterizes the distribution and takes values in  $(0, 1)$ , while  $k$  is the rating. Recall that  $\binom{N}{K}$  is the binomial coefficient, i.e.,  $N!/(K!(N-K)!)$ .

- (a) We cook up some basis functions  $\phi_j(\mathbf{x})$  and we plan to weight them using a set of weights  $\mathbf{w}$  to determine  $\rho$ . However,  $\phi(\mathbf{x})^\top \mathbf{w}$  can be negative and can be greater than one. How can we map it into the right space?
- (b) Having figured out how to get a map into the right space, write down the log likelihood of a set of  $N$  data  $\{t_n, \mathbf{x}_n\}_{n=1}^N$ . You can ignore constants in the sum that don't depend on the inputs or  $\mathbf{w}$ .
- (c) Compute the gradient of the log likelihood in terms of  $\mathbf{w}$ . Hint: the derivative of the logistic function is  $\frac{d}{dz} \sigma(z) = \sigma(z)(1 - \sigma(z))$ .

#### 4. Hyperplanes and Discriminant functions

Suppose we have the discriminant function  $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ , and if  $y(\mathbf{x}) \geq 0$  we assign  $\mathbf{x}$  to  $\mathcal{C}_1$ , and if  $y(\mathbf{x}) < 0$  we assign  $\mathbf{x}$  to  $\mathcal{C}_2$ . Show that for any  $\mathbf{x}_0, \mathbf{x}_1$  on the decision boundary  $(\mathbf{x}_0 - \mathbf{x}_1)$  is perpendicular to the vector  $\mathbf{w}$ .

## 5. Fisher Criterion in Matrix Form (Bishop 4.5)

The Fisher Criterion is defined as

$$J(\mathbf{w}) = \frac{(\mathbf{m}_2 - \mathbf{m}_1)^2}{s_1^2 + s_2^2},$$

where

$$\mathbf{m}_k = \mathbf{w}^\top \mathbf{m}_k$$

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{\mathbf{x} \in \mathcal{C}_k} \mathbf{x}$$

$$s_k^2 = \sum_{\mathbf{x} \in \mathcal{C}_k} (\mathbf{w}^\top \mathbf{x} - \mathbf{m}_k)^2$$

Show that we can write  $J(\mathbf{w})$  in matrix form as

$$J(\mathbf{w}) = \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}},$$

where

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

and

$$\mathbf{S}_W = \sum_{\mathbf{x} \in \mathcal{C}_1} (\mathbf{x} - \mathbf{m}_1)(\mathbf{x} - \mathbf{m}_1)^\top + \sum_{\mathbf{x} \in \mathcal{C}_2} (\mathbf{x} - \mathbf{m}_2)(\mathbf{x} - \mathbf{m}_2)^\top$$

## 6. Classification with Same-Mean Different-Variance Gaussians

Consider the task of recognizing which of two Gaussian distributions a data point  $\mathbf{x} = (x_1, x_2, \dots, x_D)$  comes from. We will assume that the two distributions have exactly the same mean but different variances. Let the probability that  $\mathbf{x}$  is in class  $C_i$  (where  $i \in \{0, 1\}$ ) be given by

$$\Pr(\mathbf{x}|C_i) = \prod_{j=1}^D \mathcal{N}(x_j|\mu_j, \sigma_{ij})$$

Show that  $P(C_0|x)$  can be written in the form

$$P(C_0|x) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{y} + \theta)}$$

where  $y_i$  is an appropriate function of  $x_i$ ,  $y_i = g(x_i)$ , and  $\theta$  is some constant.

## 7. Margin Distances

Consider the hyperplane given by  $\boldsymbol{w}^T \boldsymbol{x} + b = 0$ . For an arbitrary data point  $\boldsymbol{x}$ , what is the distance between  $\boldsymbol{x}$  and the hyperplane, in terms of  $\boldsymbol{w}$  and  $b$ ?