

# CS181 Practice Questions: Model Selection and Linear Classification

## 1. Coin Flipping

You have 2 biased coins: one that comes up heads with probability 0.8 and another that comes up with heads with probability 0.2. You select a coin randomly and you observe 5 heads out of 10 coin flips. Calculate  $p(D)$  for the models that represent the first and second coin respectively. Which coin are you more likely to have selected?

$$\text{For coin 1: } p(D) = \binom{10}{5} (.8)^5 (.2)^5$$

$$\text{For coin 2: } p(D) = \binom{10}{5} (.8)^5 (.2)^5$$

The terms are equal - you are uncertain about which coin you have.

## 2. Model Selection

Suppose you had three models,  $M_1$ ,  $M_2$ ,  $M_3$ , each increasing in complexity. For example, you could imagine that the models represented unregularized polynomial regression, with  $M_1$  linear regression,  $M_2$  quadratic regression, and  $M_3$  cubic regression. Within the context of Bayesian model selection, come up with a way to penalize the complexity of a model so you did not always choose  $M_3$ . Additionally, explain why, in many cases, Bayesian model selection will recover the simplest model to explain the data without explicit penalization.

Choose the prior over the model space to assign higher probability to the simpler models. This is the same thing that we did to encourage less complex parameters in Bayesian linear regression; we put a prior over the weight vector assigning the highest probability to weight vectors that lie close to zero. This is a phenomenon known as Bayesian Occam's Razor and is discussed on in Bishop (pg. 164, fig. 3.13). Briefly, since simpler models can explain a smaller subset of data than more complex models, and each model has to integrate to 1, the simpler models can put higher probabilities on the data that they do describe than the complex models, which must put some of their mass on the data that the simpler model cannot describe.

### 3. Convex Hulls and Linear Separability

Define the convex hull of a set of data points  $(\{x_i\})$  as the set

$$\left\{ \sum_i \alpha_i x_i \text{ such that } \alpha_i \geq 0 \text{ and } \sum_i \alpha_i = 1 \right\}$$

Additionally, we say that two sets of points  $\{x_i\}$  and  $\{x'_j\}$  are linearly separable if there exists a vector  $w$  and  $w_0$  such that  $w^\top x_i + w_0 > 0$  for all points in the first set and  $w^\top x'_j + w_0 < 0$  for all points in the second set. Show that if two sets of points  $\{x_i\}$  and  $\{x'_j\}$  are linearly separable, their convex hulls do not intersect.

By definition, we have a  $w$  and  $w_0$  such that  $w^\top x_i + w_0 > 0$  and  $w^\top x'_k + w_0 < 0$ . Assume for the sake of contradiction that the convex hulls do intersect. Thus, there exists some set of  $\{\alpha_i\}$  and  $\{\beta_j\}$  with  $\alpha_i > 0$  and  $\beta_j > 0$  such that  $\sum_i \alpha_i x_i = \sum_j \beta_j x'_j$ . Now, we have that

$$\begin{aligned} w^\top x_i + w_0 &> 0 \implies \\ \sum_i \alpha_i (w^\top x_i + w_0) &> 0 \implies \\ \sum_i \alpha_i w^\top x_i + w_0 &> 0 \implies \\ w^\top \sum_i \alpha_i x_i + w_0 &> 0 \implies \\ w^\top \sum_j \beta_j x'_j + w_0 &> 0 \implies \\ \sum_j \beta_j w^\top x'_j + w_0 &> 0 \implies \\ \sum_j \beta_j (w^\top x'_j + w_0) &> 0. \end{aligned}$$

But this last step is a contradiction because  $w^\top y_i + w_0 < 0$  and  $\beta_i \geq 0$ .

#### 4. Perceptron Algorithm

Consider the perceptron algorithm, which is a binary classification algorithm that finds the best linear hyperplane to separate the basis-transformed input values. The error function that is minimized is 0 when the algorithm correctly labels a data point and otherwise:

$$E_p(\mathbf{w}) = - \sum_{n \in M} \mathbf{w}^\top \phi(\mathbf{x}_n) t_n, \quad (1)$$

where we sum over the mislabeled values and  $t_n = 1$  if the correct classification is  $\mathcal{C}_1$  and  $t_n = -1$  if the correct classification is  $\mathcal{C}_2$ . Derive the stochastic gradient descent relation to optimize the weight vector for this error function.

We note that the SGD is given by  $\mathbf{w}^{(n+1)} = \mathbf{w}^{(n)} - \eta \nabla E_p(\mathbf{w}) = \mathbf{w}^{(n)} + \eta \phi(\mathbf{x}_n) t_n$ . Using this, we can iteratively improve the initial  $\mathbf{w}^{(0)}$  until convergence.

## 5. Thresholded Discriminant Functions

Suppose we have the discriminant function  $y(x) = \mathbf{w}^\top \mathbf{x} + w_0$ , but that rather than assigning  $x$  to  $\mathcal{C}_1$  when  $y(x) \geq 0$  and to  $\mathcal{C}_2$  otherwise (as in Bishop 4.1.1), we instead assign  $x$  to  $\mathcal{C}_1$  when  $y(x) \geq \eta$  for some  $\eta$  and to  $\mathcal{C}_2$  otherwise. Do we gain any generality by moving to this thresholded decision rule? Why or why not?

No, we don't. In particular, for any  $x$ , we have that  $y(x) \geq \eta$  iff  $y(x) - \eta \geq 0$ , and so by setting  $\hat{w}_0 = w_0 - \eta$ , we can equivalently predict  $\mathcal{C}_1$  if  $\mathbf{w}^\top \mathbf{x} + \hat{w}_0 \geq 0$ , and  $\mathcal{C}_2$  otherwise.

## 6. Maximizing Separation Between Classes (Bishop 4.4)

Suppose, as in Fisher's Discriminant Analysis, that we want to find the vector  $\mathbf{w}$  that maximizes the distance between the means of two classes  $\mathcal{C}_1, \mathcal{C}_2$  that are projected onto it. That is, we want to maximize

$$\mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1), \quad (\text{Bishop 4.2.2})$$

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

- (a) Show that by maximizing the criterion above subject to the constraint that  $\mathbf{w}^\top \mathbf{w} = 1$ , we find that  $\mathbf{w}_{\max} \propto (\mathbf{m}_2 - \mathbf{m}_1)$ . That is,  $\mathbf{w}_{\max} = \alpha(\mathbf{m}_2 - \mathbf{m}_1)$  for some  $\alpha$ .
- (b) Geometrically, what is the interpretation of  $\mathbf{w}_{\max}$ ?

- (a) We can write down the Lagrangian as  $L = \mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1) + \lambda(\mathbf{w}^\top \mathbf{w} - 1)$ . Taking gradients wrt  $\mathbf{w}$ , we get  $\nabla_{\mathbf{w}} L = \mathbf{m}_2 - \mathbf{m}_1 + 2\lambda \mathbf{w}$ . Setting the gradient to zero and solving gives the result.
- (b) Geometrically,  $\mathbf{w}$  is parallel to a line drawn between the means  $\mathbf{m}_1$  and  $\mathbf{m}_2$ .

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

The Fisher Criterion is defined as

$$J(\mathbf{w}) = \frac{(\mathbf{m}_2 - \mathbf{m}_1)^2}{s_1^2 + s_2^2}, \quad (\text{Bishop 4.2.5})$$

where

$$\begin{aligned} m_k &= \mathbf{w}^\top \mathbf{m}_k \\ \mathbf{m}_k &= \frac{1}{N_k} \sum_{\mathbf{x} \in \mathcal{C}_k} \mathbf{x} \\ s_k^2 &= \sum_{\mathbf{x} \in \mathcal{C}_k} (\mathbf{w}^\top \mathbf{x} - m_k)^2 \end{aligned}$$

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

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

For the numerator, we have

$$\begin{aligned} (\mathbf{m}_2 - \mathbf{m}_1)^2 &= (\mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1))^2 \\ &= \mathbf{w}^\top (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^\top \mathbf{w} \\ &= \mathbf{w}^\top \mathbf{S}_B \mathbf{w} \end{aligned}$$

For the denominator, we have

$$\begin{aligned} s_1^2 + s_2^2 &= \sum_{x \in \mathcal{C}_1} (\boldsymbol{w}^\top \boldsymbol{x} - m_1)^2 + \sum_{x \in \mathcal{C}_2} (\boldsymbol{w}^\top \boldsymbol{x} - m_2)^2 \\ &= \sum_{x \in \mathcal{C}_1} (\boldsymbol{w}^\top (\boldsymbol{x} - \boldsymbol{m}_1))^2 + \sum_{x \in \mathcal{C}_2} (\boldsymbol{w}^\top (\boldsymbol{x} - \boldsymbol{m}_2))^2 \\ &= \sum_{x \in \mathcal{C}_1} \boldsymbol{w}^\top (\boldsymbol{x} - \boldsymbol{m}_1) (\boldsymbol{x} - \boldsymbol{m}_1)^\top \boldsymbol{w} + \sum_{x \in \mathcal{C}_2} \boldsymbol{w}^\top (\boldsymbol{x} - \boldsymbol{m}_2) (\boldsymbol{x} - \boldsymbol{m}_2)^\top \boldsymbol{w} \\ &= \boldsymbol{w}^\top \boldsymbol{S}_W \boldsymbol{w} \end{aligned}$$