

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?

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.

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.

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.

5. Thresholded Discriminant Functions

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

6. Maximizing Separation Between Classes (Bishop 4.4)

Suppose, as in Fisher's Discriminant Analysis, that we want to find the vector \boldsymbol{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

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

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

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

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

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