

---

## ENSEMBLE LEARNING PROBLEM SHEET

---

**1** These questions have both appeared in past examinations.

(a) If  $\{X_i | i = 1, 2, \dots, n\}$  is a set of correlated random variables such that

$$\mathbb{E}[X_i] = \mu \quad \mathbb{E}[(X_i - \mu)(X_j - \mu)] = \begin{cases} \sigma^2 & \text{if } i = j \\ \rho\sigma^2 & \text{if } i \neq j \end{cases}$$

show

$$\mathbb{E} \left[ \left( \frac{1}{n} \sum_{i=1}^n X_i - \mu \right)^2 \right] = \rho\sigma^2 + \frac{(1-\rho)\sigma^2}{n}$$

[10 marks]

---

$$\begin{aligned} \mathbb{E} \left[ \left( \frac{1}{n} \sum_{i=1}^n X_i - \mu \right)^2 \right] &= \mathbb{E} \left[ \left( \frac{1}{n} \sum_{i=1}^n (X_i - \mu) \right)^2 \right] = \frac{1}{n^2} \mathbb{E} \left[ \left( \sum_{i=1}^n (X_i - \mu) \right)^2 \right] \\ &= \frac{1}{n^2} \mathbb{E} \left[ \sum_{i=1}^n \sum_{j=1}^n (X_i - \mu)(X_j - \mu) \right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[(X_i - \mu)(X_j - \mu)] \\ &= \frac{1}{n^2} \sum_{i=1}^n \sigma^2 + \frac{1}{n^2} \sum_{\substack{i,j=1 \\ i \neq j}}^n \rho\sigma^2 = \frac{1}{n^2} n\sigma^2 + \frac{1}{n^2} n(n-1)\rho\sigma^2 \\ &= \rho\sigma^2 + \frac{1-\rho}{n}\sigma^2 \end{aligned}$$

---

- (b) We consider a regression problem where the data  $(x, y)$  is distributed according to  $\gamma(x, y)$ . We consider a learning machine that makes a prediction  $\hat{f}(x|\theta)$ , where the parameters,  $\theta$  are trained using a stochastic algorithm that returns parameters distributed according to a probability density  $\rho(\theta)$ . We can define the mean machine as  $\hat{m}(x) = \mathbb{E}_{\theta \sim \rho}[\hat{f}(x|\theta)]$ . We assume that

$$\mathbb{E}_{(x,y) \sim \gamma}[(\hat{m}(x) - y)^2] = B, \quad \mathbb{E}_{(x,y) \sim \gamma} \left[ \mathbb{E}_{\theta \sim \rho} \left[ \left( \hat{f}(x|\theta) - \hat{m}(x) \right)^2 \right] \right] = V.$$

That is, we can define a bias  $B$  and variance  $V$ . We now consider ensembling  $n$  machines

$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n \hat{f}(x|\theta_i)$$

where  $\theta_i$  are drawn independently from  $\rho(\theta)$ . Compute the expected generalisation error of  $\hat{f}_n(x)$ . (Note this is different from the usual bias-variance calculation because we are averaging the performance of  $n$  machines). [10 marks]

We note that  $\mathbb{E}_{\theta}[\hat{f}_n(x)] = \hat{m}(x)$  so the bias variance dilemma calculation still holds

$$\begin{aligned} \mathcal{E}_n &= \mathbb{E}_{\theta} \left[ \mathbb{E}_{(x,y)} \left[ \left( \hat{f}_n(x) - y \right)^2 \right] \right] = \mathbb{E}_{\theta} \left[ \mathbb{E}_{(x,y)} \left[ \left( (\hat{f}_n(x) - \hat{m}(x)) + (\hat{m}(x) - y) \right)^2 \right] \right] \\ &= \mathbb{E}_{(x,y)} \left[ (\hat{m}(x) - y)^2 \right] + \mathbb{E}_{\theta} \left[ \mathbb{E}_{(x,y)} \left[ \left( \hat{f}_n(x) - \hat{m}(x) \right)^2 \right] \right] \end{aligned}$$

where the cross term vanishes as

$$\mathbb{E}_{\theta} \left[ \left( \hat{f}_n(x) - \hat{m}(x) \right) (y - \hat{m}(x)) \right] = (y - \hat{m}(x)) \mathbb{E}_{\theta} \left[ \hat{f}_n(x) - \hat{m}(x) \right] = 0.$$

Thus  $\mathcal{E}_n$  consists of a bias  $B$  and a new variance  $V_n$  where

$$\begin{aligned} V_n &= \mathbb{E}_{\theta} \left[ \mathbb{E}_{(x,y)} \left[ \left( \hat{f}_n(x) - \hat{m}(x) \right)^2 \right] \right] = \mathbb{E}_{\theta} \left[ \mathbb{E}_{(x,y)} \left[ \left( \frac{1}{n} \sum_{i=1}^n (\hat{f}(x|\theta_i) - \hat{m}(x)) \right)^2 \right] \right] \\ &= \mathbb{E}_{\theta} \left[ \mathbb{E}_{(x,y)} \left[ \frac{1}{n^2} \sum_{i=1}^n (\hat{f}(x|\theta_i) - \hat{m}(x))^2 + \frac{1}{n^2} \sum_{i,j=1, j \neq i}^n (\hat{f}(x|\theta_i) - \hat{m}(x)) (\hat{f}(x|\theta_j) - \hat{m}(x)) \right] \right] \\ &= \mathbb{E}_{\theta} \left[ \mathbb{E}_{(x,y)} \left[ \frac{1}{n^2} \sum_{i=1}^n (\hat{f}(x|\theta_i) - \hat{m}(x))^2 \right] \right] = \frac{V}{n} \end{aligned}$$

since

$$\mathbb{E}_{\theta} \left[ (\hat{f}(x|\theta_i) - \hat{m}(x)) (\hat{f}(x|\theta_j) - \hat{m}(x)) \right] = 0$$

as  $\theta_i$  and  $\theta_j$  are independent. Thus the expected generalisation performance is equal to

$$\mathcal{E}_n = B + \frac{V}{n}.$$

That is, by averaging over  $n$  different machines we reduce the variance by  $n$ . Note that we assume the predictions of the machine where independent so

$$\mathbb{E}_{\theta} \left[ \left( \hat{f}(\mathbf{x}|\theta_i) - \hat{m}(\mathbf{x}) \right) \left( \hat{f}(\mathbf{x}|\theta_j) - \hat{m}(\mathbf{x}) \right) \right] = 0$$

If the predictions of the machines are correlated then we don't do so well.

---

End of question 1

**END OF PAPER**