

## SEMESTER 2 EXAMINATION 2022/23

## ADVANCED MACHINE LEARNING

Duration 120 mins (2 hours)

This paper is a WRITE-ON examination paper.

You **must** write your Student ID on this Page and must not write your name anywhere on the paper.

All answers should be written within the designated boxes in this examination paper and sufficient space is provided for each question.

If, for some reason, space is required to complete or correct an answer to a question, use the "Additional Space" provided on the facing or adjacent page to the question. Clearly indicate which question the answer corresponds to.

No credit will be given for answers presented elsewhere and without clear indication of to what question they correspond. Blue answer books may be used for scratch; they will be discarded without being looked at.

Answer all parts of the question in section A (40 marks)  
and ALL three questions from section B (20 marks each)

Student ID:

Question	Mark	Arithmetic checked	Double Marked
A1	/40		
B2	/20		
B3	/20		
B4	/20		
Total:	/100		

University approved calculators MAY be used.

A foreign language translation dictionary (paper version) is permitted provided it contains no notes, additions or annotations.

**10 page examination paper**

## Section A

### A 1

- (a) Explain why in boosting it is important to use as decorrelated trees as possible and explain how this is achieved in random forest. [5 marks]

---

**The idea of boosting is to reduce the variance in the prediction of the machine by averaging over many machines. This only works if the machines are decorrelated. In random forest this is achieved both by using booststapping and by building decision trees that use only a small randomly sampled subset of the features.**

---

- (b) Multilayer perceptrons (MLPs) are famously universal approximators (they can approximate any smooth function up to arbitrary precision). This should make them susceptible to massively overfitting a training set. Explain why, despite this, MLPs often have good generalisation performance. [5 marks]

---

**Although the class of MLPs (considering any architecture) is universal approximators, a particular MLP (with a fixed number of hidden units) will not be a universal approximator. In practice we choose the architecture of an MLP to fit the complexity of the function we are trying to learn (depending on the size of the training set). This prevents overfitting.**

---

- (c) Show that if  $\lambda > 0$  is an eigenvalue of  $C = X X^T$  then it is also an eigenvalue of  $D = X^T X$ , where  $X$  is a matrix. [5 marks]

---

*(Algebraic questions requiring a good grasp of linear algebra that underlies are lot of duality in machine learning.)*

**If  $\lambda \neq 0$  is an eigenvalue of  $C$  then there exists some eigenvector  $v (\neq 0)$  such that  $C v = \lambda v$ . Multiplying on the left by  $X^T$  then**

$$X^T C v = \lambda X^T v$$

**but  $X^T C = X^T (X X^T) = (X^T X) X^T = D X^T$  (where we have used the associativity of matrix multiplication and the definitions of  $C$  and  $D$ ). Thus**

$$D X^T v = \lambda X^T v$$

**or  $D u = \lambda u$  where  $u = X^T v \neq 0$  is an eigenvector of  $D$  with eigenvalue  $\lambda$ .**

---

- (d) If  $\|x\|$  is a proper norm, use the triangular inequality ( $\|x + y\| \leq \|x\| + \|y\|$ ), linearity of a norm ( $\|a x\| = a \|x\|$ ) and the definition of convexity, to show that the norm is convex. [5 marks]

**(The students have not seen this. The algebra is simple, but the students need a good understanding to do the proof.)**

**For any vectors,  $x$ , and  $y$  and any scalar  $a \in [0, 1]$  then**

$$\|a x + (1 - a) y\| \leq \|a x\| + \|(1 - a) y\| = a \|x\| + (1 - a) \|y\|$$

**where the first inequality follows from the triangular inequality and the second equality from the linearity of the norm. However,**

$$\|a x + (1 - a) y\| \leq a \|x\| + (1 - a) \|y\|$$

**is the defining equation of convexity.**

- (e) Use the fact norms are convex to argue that an elastic net with a loss function

$$L(w) = \sum_{i=1}^m (w^T x_i - y_i)^2 + \alpha \|w\|_{L_1} + \beta \|w\|_{L_2}^2$$

has a unique minimum.

[5 marks]

**The quadratic minima and the two norms are all convex (up) functions. The sum of convex functions is convex. We note that the  $L_2$  norm,  $\|w\|_{L_2}^2$ , is strongly convex which ensures that the minimum is unique.**

- (f) Explain the main advantage and disadvantage of stochastic gradient descent (SGD) compared to full gradient descent. [5 marks]

**The overwhelming advantage of SGD is that it is far faster to compute the gradient with respect to a mini-batch than computing the gradient using the whole training set.**

**A significant disadvantage is that the gradient estimates can have very high added noise. Thus the algorithm will typically wander around. It will also never converge.**

- (g) The multinomial distribution  $\text{Multi}(\mathbf{k}|n, \mathbf{p})$  describes the likelihood of observed counts  $\mathbf{k} = (k_1, k_2, \dots, k_d)$  for  $d$  possible outcomes, where the total number of counts is  $\sum_{i=1}^d k_i = n$ . The vector  $\mathbf{p} = (p_1, p_2, \dots, p_d)$  is a vector of probabilities (summing to 1) where  $p_i$  is the probability of outcome  $i$  occurring. Show that the Dirichlet distribution,  $\text{Dir}(\mathbf{p}|\alpha)$  is a conjugate prior to the multinomial likelihood  $\text{Multi}(\mathbf{k}|n, \mathbf{p})$  where

$$\text{Dir}(\mathbf{p}|\alpha) = \Gamma(\alpha_0) \prod_{i=1}^d \frac{p_i^{\alpha_i-1}}{\Gamma(\alpha_i)} \quad \text{Multi}(\mathbf{k}|n, \mathbf{p}) = n! \prod_{i=1}^d \frac{p_i^{k_i}}{k_i!}$$

where  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_d)$  is a vector of parameters that controls the Dirichlet distribution with  $\alpha_0 = \sum_{i=1}^d \alpha_i$ . Derive the update equation for the parameters  $\alpha$ . [5 marks]

**We only need to consider the functional form with respect to  $p_i$ . Thus the posterior is proportional to**

$$f(\mathbf{p}|\mathbf{k}) \propto \prod_{i=1}^d \frac{p_i^{k_i}}{k_i!} \prod_{i=1}^d \frac{p_i^{\alpha_i-1}}{\Gamma(\alpha_i)} \propto \prod_{i=1}^d p_i^{\alpha_i+k_i-1} \propto \text{Dir}(\mathbf{p}|\alpha + \mathbf{k}).$$

**The update equation is thus  $\alpha' = \alpha + \mathbf{k}$ .**

- (h) Describe how the minimum description length formalism is used for model selection. [5 marks]

**In the minimum description length formalism we choose a model that generates a prediction that minimises the length of the model plus the length of the residual errors between the prediction of the model and the true data. The description length of model is a minimum encoding of the parameters of the model, while the length message needed to communicate the residual errors is usually taken to be the Shannon bound.**

End of question A1

## Section B

### B 2

(a) Show that the expected generalisation given by

$$\mathbb{E}_{\mathcal{D}}[E(\mathcal{D})] = \mathbb{E}_{\mathcal{D}} \left[ \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right]$$

can be written as the sum of a bias term,  $B$ , and variance term  $V$  where

$$B = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right)^2, \quad V = \mathbb{E}_{\mathcal{D}} \left[ \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right]$$

where  $\hat{f}_m(\mathbf{x}) = \mathbb{E}_{\mathcal{D}}[\hat{f}(\mathbf{x}|\mathcal{D})]$  is the prediction made by averaging over all machines.

[10 marks]

---

**(Students find this derivation difficult because it requires understanding expectations. Being rather conceptual students often struggle with this. Two marks for: (i) writing the first line, (ii) using the trick of adding and subtracting the prediction of the mean machine, (iii) multiplying out the square, (iv) arguing the cross term vanishes (v) explaining the final answer.)**

**The first step is to subtract and add the response of the mean machine**

$$\begin{aligned} \bar{E}_G &= \mathbb{E}_{\mathcal{D}}[E_G(\mathcal{D})] = \mathbb{E}_{\mathcal{D}} \left[ \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right] \\ &= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - f(\mathbf{x}) \right)^2 \right] \\ &= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}_{\mathcal{D}} \left[ \left( \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right) + \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \right)^2 \right] \\ &= \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \left( \mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right)^2 \right] + \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right. \\ &\quad \left. + \mathbb{E}_{\mathcal{D}} \left[ 2 \left( \hat{f}(\mathbf{x}|\mathcal{D}) - \hat{f}_m(\mathbf{x}) \right) \left( \hat{f}_m(\mathbf{x}) - f(\mathbf{x}) \right) \right] \right) \end{aligned}$$

**The last term vanishes on taking the expectation. This leaves us with the variance term plus the bias.**

---

- (b) Explain in words (1) the *bias* and (2) the *variance* terms and (3) the dilemma. [6 marks]

- 
- (i) The *bias* is the generalisation performance of the mean machine (i.e. the prediction made by taking the mean response averaged over machines trained with every possible training set of a given size).
- (ii) The *variance* measures how the responses of the individual machines vary from machine to machine.
- (iii) To get good generalisation we want to both reduce the bias and variance. To reduce the bias we need a complex machine, but that will typically lead to an increase of the variance. Conversely a simple machine is likely to have a small variance, but a high bias.
- 

- (c) Consider a classification problem where  $\hat{f}_c(\mathbf{x}|\boldsymbol{\theta}_D)$  is the softmax output of a learning machine trained on a dataset  $D$  for class  $c$ . Let  $\hat{m}_c(\mathbf{x}) = \mathbb{E}_D[\hat{f}_c(\mathbf{x}|\boldsymbol{\theta}_D)]$  be the output of the mean machine for class  $c$  (i.e. the outputs averaged over machines trained on all possible datasets). Consider the cross-entropy loss

$$L(\mathbf{x}, y, \boldsymbol{\theta}) = - \sum_{(\mathbf{x}, y) \in \mathcal{D}} \sum_{c \in \mathcal{C}} \mathbb{I}[y = c] \log(\hat{f}_c(\mathbf{x}|\boldsymbol{\theta}_D))$$

where  $\mathbb{I}[y = c]$  is an indicator function equal to 1 if the target  $y$  is equal to class  $c$ , and 0 otherwise. Show that the expected loss over all training sets can be written as the expected loss for the mean machine (a bias) plus an additional term (a variance). Use Jensen's inequality ( $\mathbb{E}[\log(X)] \leq \log(\mathbb{E}[X])$ ) to show the second term is non-negative. [4 marks]

---

**We can add and subtract the loss for the mean machine**

$$\begin{aligned} \bar{L} &= -\mathbb{E}_{(\mathbf{x}, y)} \left[ \mathbb{E}_D \left[ \sum_{c \in \mathcal{C}} \mathbb{I}[y = c] \log(\hat{f}_c(\mathbf{x}|\boldsymbol{\theta}_D)) \right] \right] \\ &= -\mathbb{E}_{(\mathbf{x}, y)} \left[ \sum_{c \in \mathcal{C}} \mathbb{I}[y = c] \log(\hat{m}_c(\mathbf{x})) \right] - \mathbb{E}_{(\mathbf{x}, y)} \left[ \mathbb{E}_D \left[ \sum_{c \in \mathcal{C}} \mathbb{I}[y = c] \log\left(\frac{\hat{f}_c(\mathbf{x}|\boldsymbol{\theta}_D)}{\hat{m}_c(\mathbf{x})}\right) \right] \right] \end{aligned}$$

**The first term acts like a bias. The second (variance-like) term is**

$$V = -\mathbb{E}_{(\mathbf{x}, y)} \left[ \sum_{c \in \mathcal{C}} \mathbb{I}[y = c] \mathbb{E}_D \left[ \log(\hat{f}_c(\mathbf{x}|\boldsymbol{\theta}_D)) \right] \right] + \mathbb{E}_{(\mathbf{x}, y)} \left[ \sum_{c \in \mathcal{C}} \mathbb{I}[y = c] \log(\hat{m}_c(\mathbf{x})) \right]$$

**But using Jensen's inequality**

$$\mathbb{E}_D \left[ \log(\hat{f}_c(\mathbf{x}|\boldsymbol{\theta}_D)) \right] \leq \log(\mathbb{E}_D[\hat{f}_c(\mathbf{x}|\boldsymbol{\theta}_D)]) = \log(\hat{m}_c(\mathbf{x})). \text{ Thus}$$

$$V \geq -\mathbb{E}_{(\mathbf{x}, y)} \left[ \sum_{c \in \mathcal{C}} \mathbb{I}[y = c] \log(\hat{m}_c(\mathbf{x})) \right] + \mathbb{E}_{(\mathbf{x}, y)} \left[ \sum_{c \in \mathcal{C}} \mathbb{I}[y = c] \log(\hat{m}_c(\mathbf{x})) \right] = 0$$

**that is, the second term is non-negative.**

---

End of question B2

## B 3

(a) Show that for the mapping

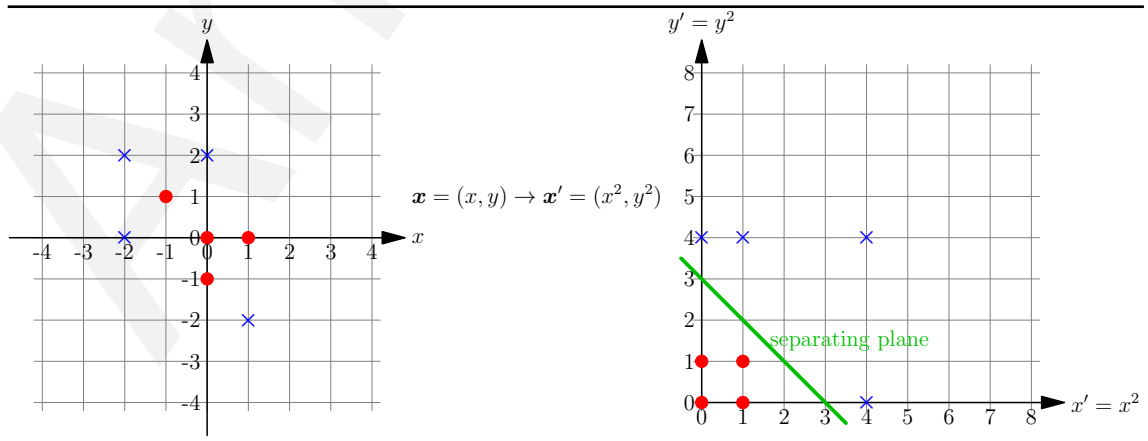
$$\mathbf{x} = (x_1, x_2, x_3)^T \rightarrow \vec{\phi}(\mathbf{x}) = (x_1^2, x_2^2, x_3^2, \sqrt{2} x_1 x_2, \sqrt{2} x_1 x_3, \sqrt{2} x_2 x_3)^T$$

the kernel  $K(\mathbf{x}, \mathbf{y}) = \langle \vec{\phi}(\mathbf{x}), \vec{\phi}(\mathbf{y}) \rangle$  is equal to  $(\mathbf{x}^T \mathbf{y})^2$ . [5 marks]

*(This is an extension of the 2-D example shown in the lecture.)*

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= \langle \vec{\phi}(\mathbf{x}), \vec{\phi}(\mathbf{y}) \rangle \\ &= (x_1^2, x_2^2, x_3^2, \sqrt{2} x_1 x_2, \sqrt{2} x_1 x_3, \sqrt{2} x_2 x_3) \\ &\quad (y_1^2, y_2^2, y_3^2, \sqrt{2} y_1 y_2, \sqrt{2} y_1 y_3, \sqrt{2} y_2 y_3)^T \\ &= x_1^2 y_1^2 + x_2^2 y_2^2 + x_3^2 y_3^2 + 2 x_1 y_1 x_2 y_2 + 2 x_1 y_1 x_3 y_3 + 2 x_2 y_2 x_3 y_3 \\ &= (x_1 y_1 + x_2 y_2 + x_3 y_3)^2 = (\mathbf{x}^T \mathbf{y})^2 \end{aligned}$$

(b) Show how the data points  $\{\mathbf{x}_i = (x_i, y_i) | i = 1, 2, \dots\}$  shown below transform under the mapping  $\mathbf{x} = (x_i, y_i) \rightarrow \mathbf{x}' = (x_i^2, y_i^2)$  and sketch the position of the maximal margin dividing plane in the new feature space. [5 marks]



(c) Give three conditions that any positive semi-definite kernel,  $K(x, y)$ , should satisfy.

[6 marks]

- For any function  $f(x)$

$$\iint f(x) K(x, y) f(y) dx dy \geq 0.$$

- The eigenvalues are non-negative.
- Can be written as an inner-product

$$K(x, y) = \langle \phi(x), \phi(y) \rangle = \sum_i \phi_i(x) \phi_i(y)$$

for some (possibly infinite) vector function  $\phi(x)$ .

(d) Using properties of positive semi-definite kernels to show that

$$K^{(3)}(x, y) = K^{(2)}(x, y) K^{(1)}(x, y)$$

is positive semi-definite if  $K^{(1)}(x, y)$  and  $K^{(2)}(x, y)$  are positive semi-definite.

[4 marks]

As  $K^{(1)}(x, y)$  and  $K^{(2)}(x, y)$  are positive-definite we can write

$$K^{(1)}(x, y) = \sum_i \phi_i^{(1)}(x) \phi_i^{(1)}(y) \quad K^{(2)}(x, y) = \sum_j \phi_j^{(2)}(x) \phi_j^{(2)}(y)$$

and thus we can write

$$\begin{aligned} K^{(3)}(x, y) &= \sum_{i,j} \phi_i^{(1)}(x) \phi_j^{(1)}(x) \phi_i^{(2)}(y) \phi_j^{(2)}(y) \\ &= \sum_{i,j} \phi_{ij}^{(3)}(x) \phi_{ij}^{(3)}(y) = \langle \phi^{(3)}(x), \phi^{(3)}(y) \rangle \end{aligned}$$

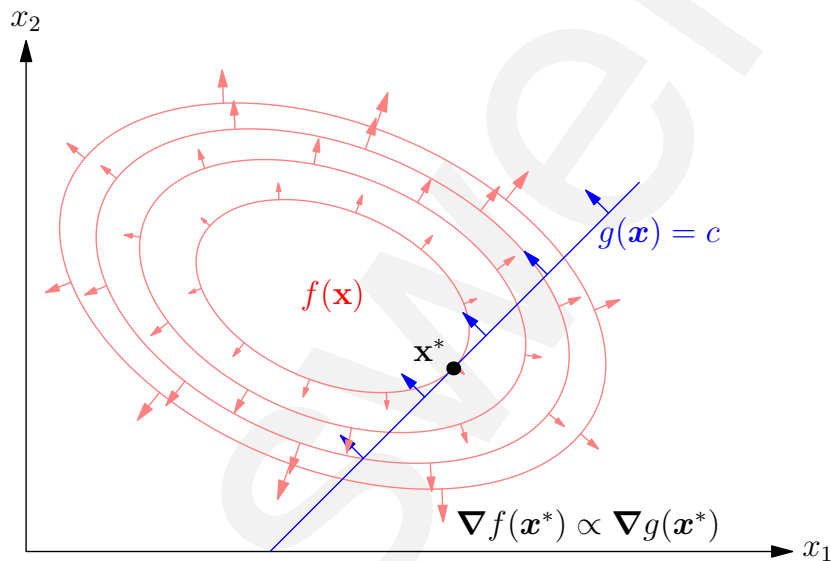
where  $\phi^{(3)}(x)$  is a vector with components  $\phi_{ij}^{(3)}(x) = \phi_i^{(1)}(x) \phi_j^{(2)}(x)$ .  
Therefore  $K^{(3)}(x, y)$  is positive semi-definite.

End of question B3



B 4

- (a) Below we show contour lines for a quadratic minimum  $f(\mathbf{x})$  and a constraint  $g(\mathbf{x}) = x_2 - x_1 = c$ . Plot the gradient  $\nabla f(\mathbf{x})$  at various points along the contour lines and  $\nabla g(\mathbf{x})$  at various points along the constraint. Mark the point that minimises  $f(\mathbf{x})$ , subject to the constraint  $g(\mathbf{x}) = c$ . Write down the condition for the minimum points. [10 marks]



**(2 marks showing gradient for  $f(\mathbf{x})$ , 2 marks showing gradient for constraint  $g(\mathbf{x}) - c = 0$ , 3 marks for identifying  $\mathbf{x}^*$ , 3 marks for noting  $\nabla f(\mathbf{x}^*) \propto \nabla g(\mathbf{x}^*)$ .)**

- (b) Consider for a dataset  $\{x_i | i = 1, 2, \dots, n\}$ . Subtracting the mean and projecting onto a vector  $v$  gives us a number  $z_i = v^T(x_i - \mu)$ . Show that the direction  $v$ , with  $\|v\|^2 = 1$ , that maximises the variance of the set of numbers  $\{z_i | i = 1, 2, \dots, n\}$ , is given by the eigenvector of the covariance matrix with the largest eigenvalue. [10 marks]

**We consider the variance**

$$\begin{aligned}\sigma^2(v) &= \frac{1}{n-1} \sum_{i=1}^n z_i^2 = \frac{1}{n-1} \sum_{i=1}^n (v^T(x_i - \mu))^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n v^T(x_i - \mu)(x_i - \mu)^T v = v^T \left( \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T \right) v = v^T C v\end{aligned}$$

**where**

$$C = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$$

**is the empirical covariance matrix. To compute the direction that maximises this variance we write a Lagrangian**

$$\mathcal{L} = v^T C v - \lambda (\|v\|^2 - 1)$$

**which is extremised when**

$$\nabla \mathcal{L} = 2Cv - 2\lambda v$$

**or  $Cv = \lambda v$ . That is, the direction is an eigenvector of  $C$ . But, for an eigenvector,  $v$  of  $C$**

$$\sigma^2(v) = v^T C v = \lambda v^T v = \lambda.$$

**Thus the direction,  $v$ , that maximises the variance of the data projected onto  $v$  is given by the eigenvector with the largest eigenvalue.**

**(4 marks for showing  $\sigma^2(v) = v^T C v$ , 2 marks for writing down Lagrangian, 2 marks for solving Lagrangian and two marks for arguing that we need to select the eigenvector with the largest eigenvalue.)**

End of question B4