

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.

16 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]

5

- (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]

5

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

[illegible]

150

- (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]

[illegible]

150

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

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

has a unique minimum.

[5 marks]

5

(f) Explain the main advantage and disadvantage of stochastic gradient descent (SGD) compared to full gradient descent.

[5 marks]

5

- (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, $\text{Dir}(\mathbf{p}|\boldsymbol{\alpha})$ is a conjugate prior to the multinomial likelihood $\text{Multi}(\mathbf{k}|n, \mathbf{p})$ where

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

where $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_d)$ 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 $\boldsymbol{\alpha}$. [5 marks]

5

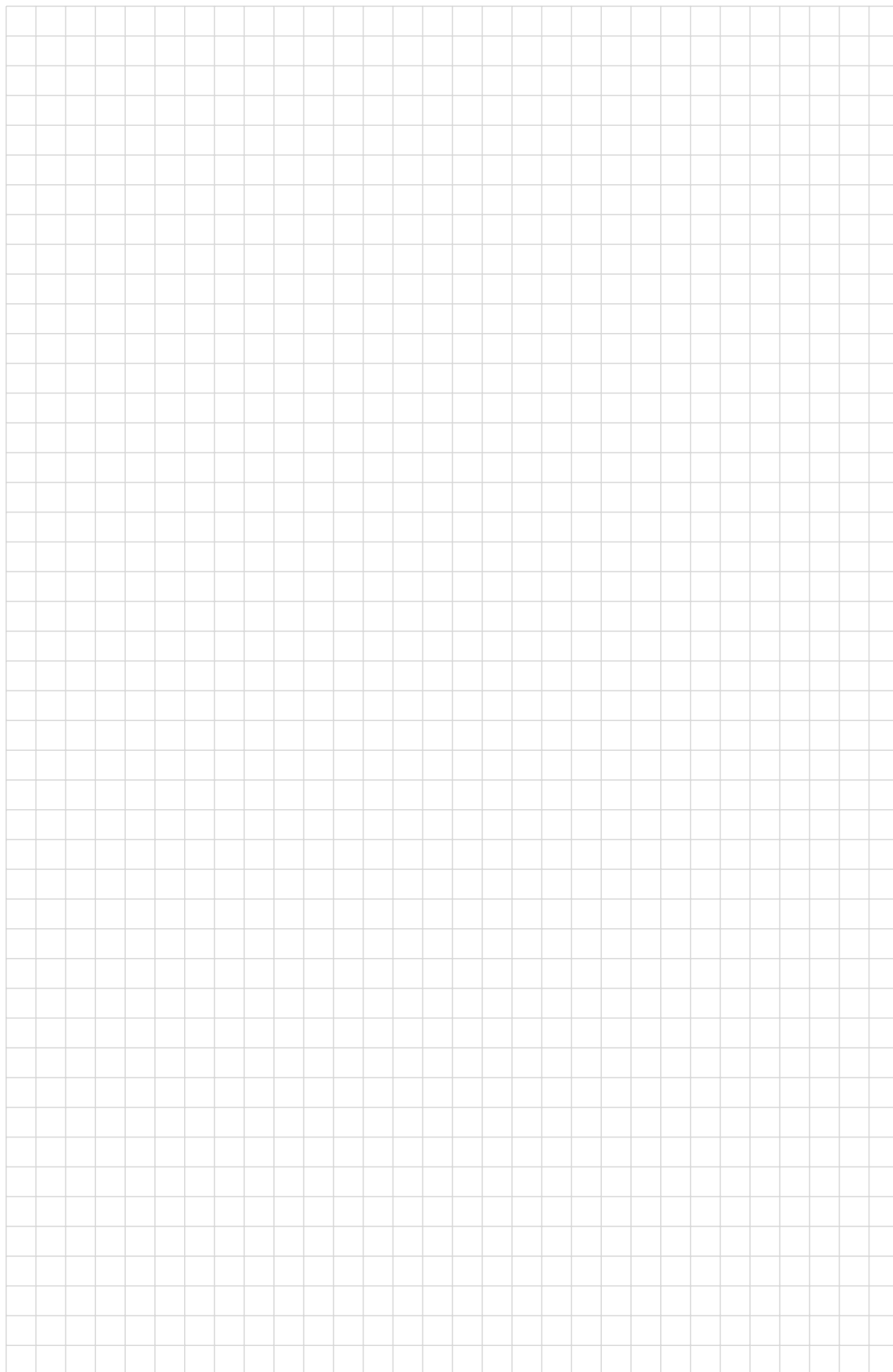
- (h) Describe how the minimum description length formalism is used for model selection. [5 marks]

5

End of question A1

(a)	$\frac{1}{5}$	(b)	$\frac{1}{5}$	(c)	$\frac{1}{5}$	(d)	$\frac{1}{5}$	(e)	$\frac{1}{5}$	(f)	$\frac{1}{5}$	(g)	$\frac{1}{5}$	(h)	$\frac{1}{5}$	Total	$\frac{1}{40}$
-----	---------------	-----	---------------	-----	---------------	-----	---------------	-----	---------------	-----	---------------	-----	---------------	-----	---------------	-------	----------------

Additional space. Do not use unless necessary. Clearly mark corresponding question.



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

1	<hr/> <hr/> <hr/> <hr/>
2	<hr/> <hr/> <hr/> <hr/>
3	<hr/> <hr/> <hr/> <hr/>

6

- (c) Consider a classification problem where $\hat{f}_c(\mathbf{x}|\boldsymbol{\theta}_{\mathcal{D}})$ is the softmax output of a learning machine trained on a dataset \mathcal{D} for class c . Let $\hat{m}_c(\mathbf{x}) = \mathbb{E}_{\mathcal{D}}[\hat{f}_c(\mathbf{x}|\boldsymbol{\theta}_{\mathcal{D}})]$ 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

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

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]

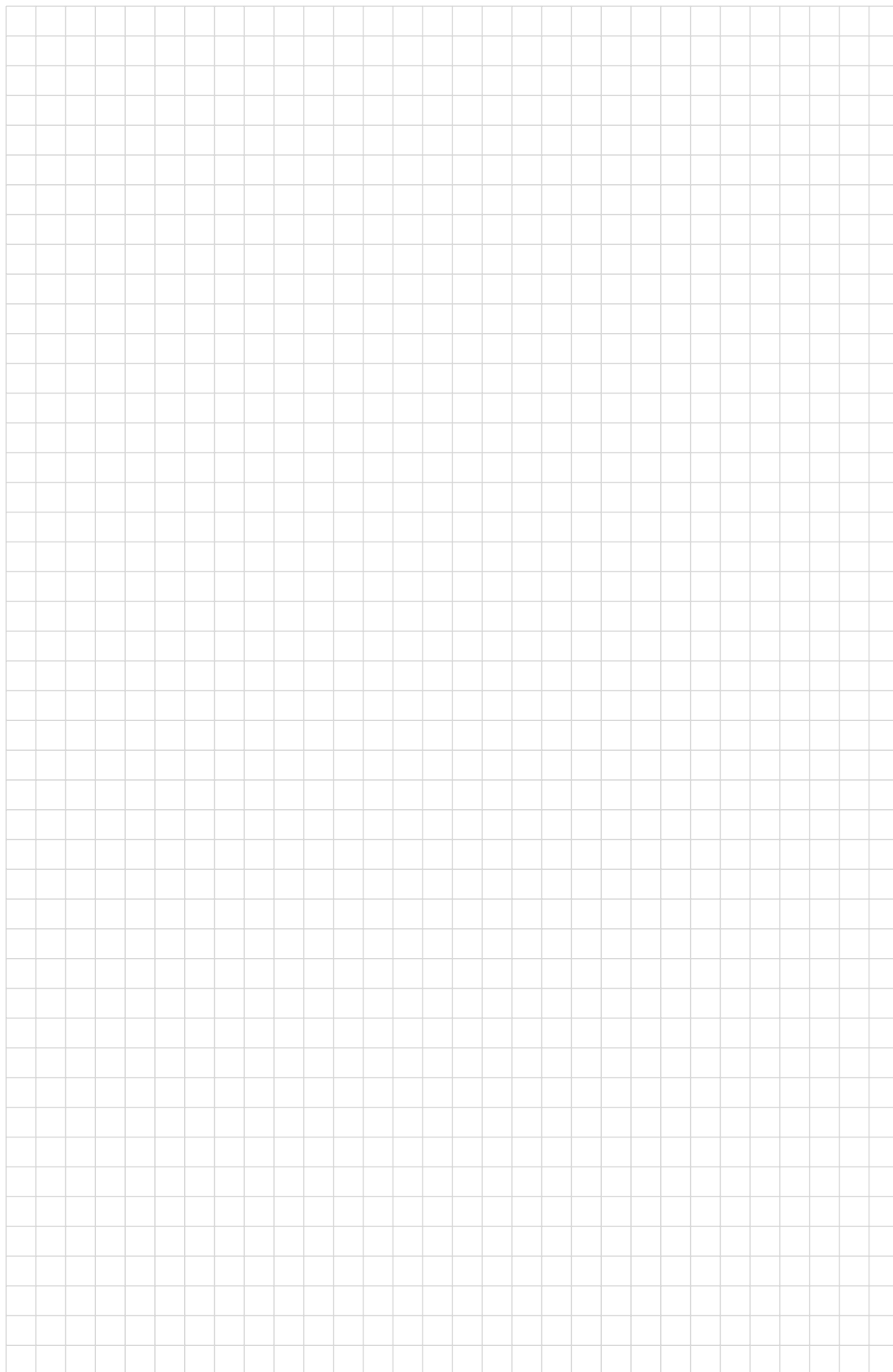
This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and extend across the width of the page. There is no handwriting or other markings on the paper.

4

End of question B2

(a) $\frac{\quad}{10}$ (b) $\frac{\quad}{6}$ (c) $\frac{\quad}{4}$ Total $\frac{\quad}{20}$

Additional space. Do not use unless necessary. Clearly mark corresponding question.



B 3

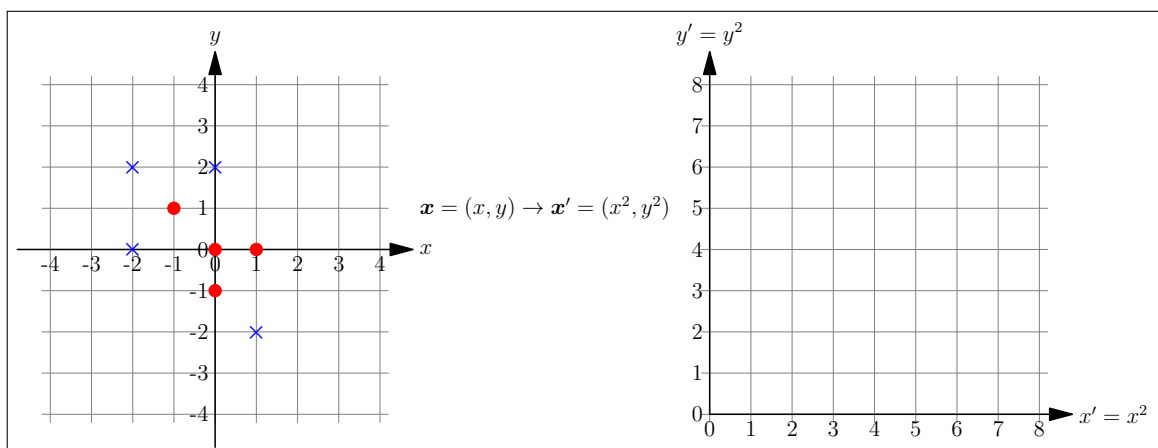
(a) Show that for the mapping

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

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]

5

(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]



5

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

[6 marks]

1	_____

2	_____

3	_____

6

- (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.

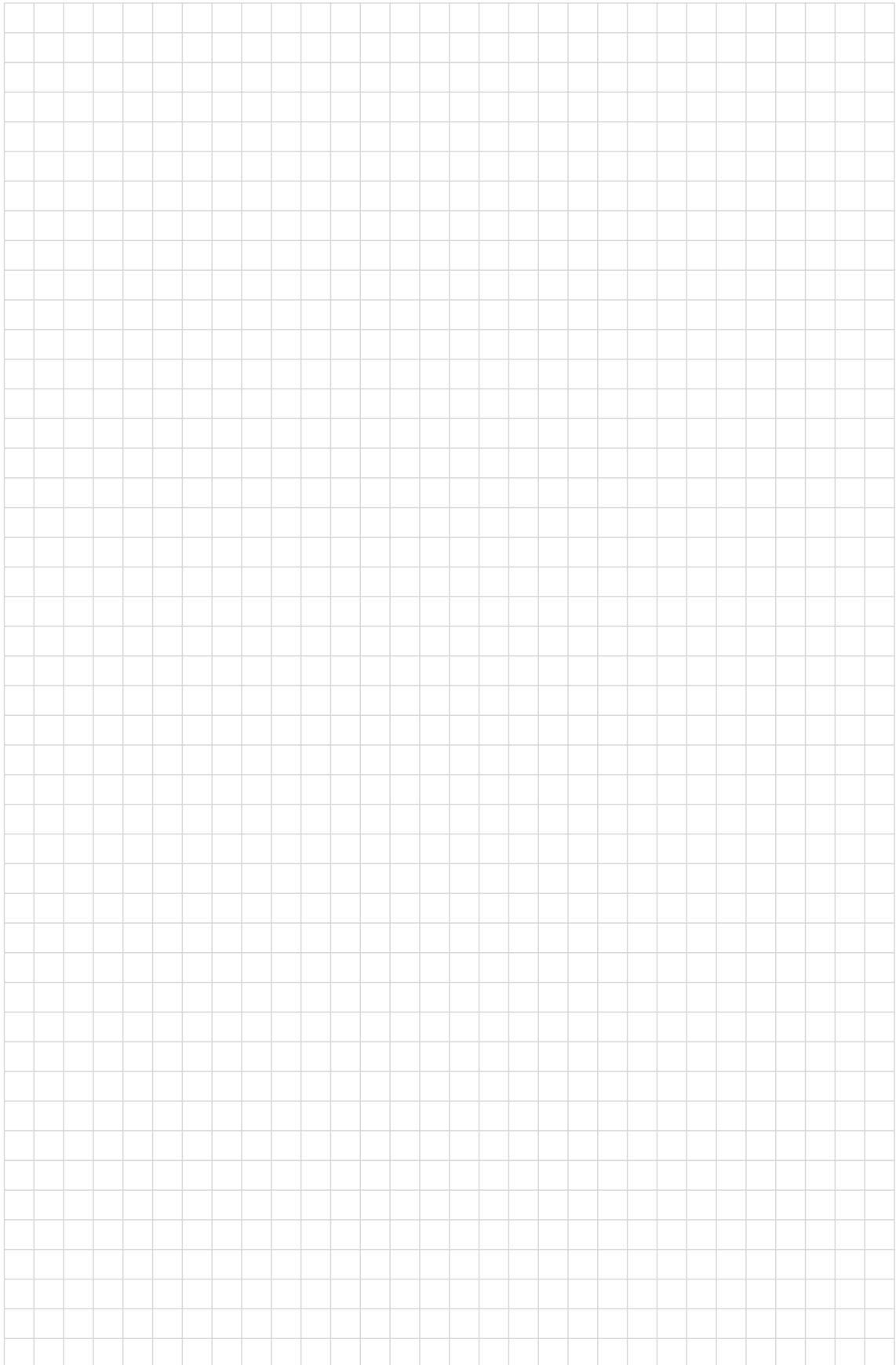
[4 marks]

4

End of question B3

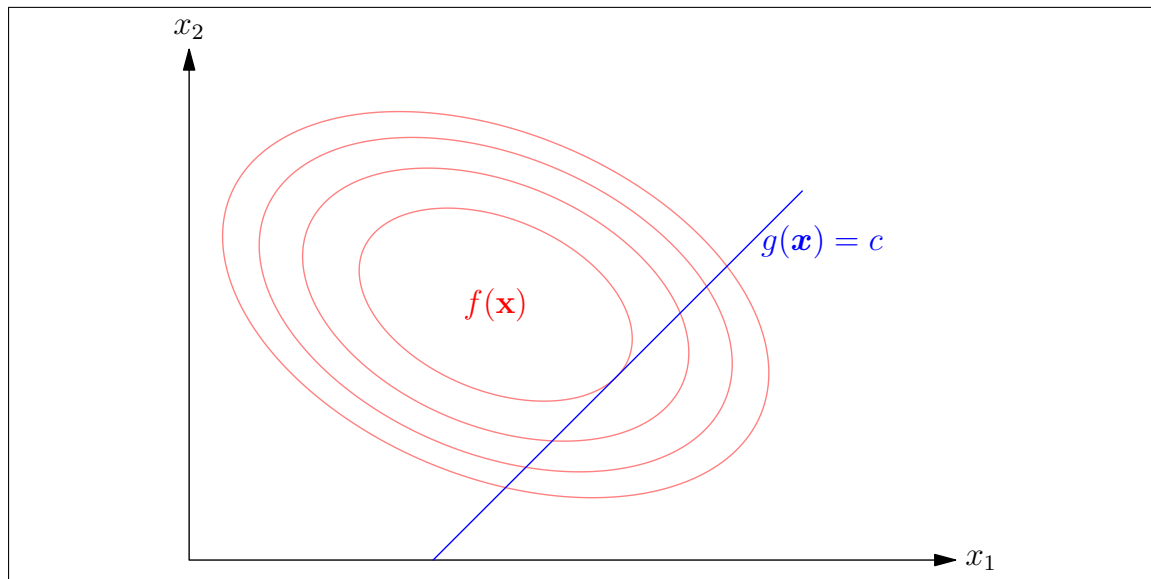
(a) $\frac{1}{5}$	(b) $\frac{1}{5}$	(c) $\frac{1}{6}$	(d) $\frac{1}{4}$	Total $\frac{1}{20}$
-------------------	-------------------	-------------------	-------------------	----------------------

Additional space. Do not use unless necessary. Clearly mark corresponding question.



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]



$\overline{10}$

- (b) Consider for a dataset $\{\mathbf{x}_i | i = 1, 2, \dots, n\}$. Subtracting the mean and projecting onto a vector \mathbf{v} gives us a number $z_i = \mathbf{v}^T(\mathbf{x}_i - \boldsymbol{\mu})$. Show that the direction \mathbf{v} , with $\|\mathbf{v}\|^2 = 1$, 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]

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and extend across the width of the page. There is no text or other markings on the paper.

10

End of question B4

(a) $\frac{\quad}{10}$	(b) $\frac{\quad}{10}$	Total $\frac{\quad}{20}$
------------------------	------------------------	--------------------------

Additional space. Do not use unless necessary. Clearly mark corresponding question.

