

Duration 120 mins (2 hours)

Answer all questions. Section A (worth 40 marks) are a series of questions with short answers. Section B (worth 60 marks) involve longer questions

Question	Mark	<i>Arithmetic checked</i>	<i>Double Marked</i>
Total:			

Page 1 of 10

Section A

Question A 1

- (a) Explain why over expressive machines are likely to generalise poorly when the number of training examples is small. (5 marks)

Over expressive machine are like to find spurious rules that are consistent with the training data. They typically have a high variance in that the rules they find vary strongly depending on the training set.

- (b) Explain why CNNs capture the structure of typical image datasets. (5 marks)

CNNs are built from local filters that respond to objects in a translationally invariant way (the same feature maps will be activated by an object irrespective of where it is).

- (c) Explain the major ways (i) gradient descent (ii) Newton's method and (iii) quasi-Newton methods differ in terms of the information they use. (5 marks)

(i) Gradient descent uses just the gradient $x^{(t+1)} = x^{(t)} - r \nabla f(x^{(t)})$

(ii) Newton's method uses both the gradient and the Hessian

$$x^{(t+1)} = x^{(t)} - H^{-1} \nabla f(x^{(t)})$$

(iii) Quasi-Newton methods use some method for approximating the Hessian.

5

5

5

- (d) In stochastic gradient descent (SGD) explain what are mini-batches and their possible advantages and disadvantages. (5 marks)

In SGD we compute the gradient of the loss function for a subset of the training examples (the mini-batch). This is far faster than computing the gradient of the loss function for the whole training set, but it is only an approximation to the true gradient.

- (e) Describe the Karush-Kuhn-Tucker (KKT) conditions for constrained optimisation. (5 marks)

The KKT conditions are used when we have inequality constraints. They state that either the Lagrange multiplier is zero and the point found satisfies the constraint and would be an optimum regardless of the constraint, or the Lagrange multiplier is non-zero and the point found lies on the constraint.

- (f) Show that the Dirichlet distribution given by $\text{Dir}(\mathbf{p}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\prod_{i=1}^d \Gamma(\alpha_i)} \prod_{i=1}^d p_i^{\alpha_i-1}$, where $\mathbf{p} = (p_1, p_2, \dots, p_d)$, $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_d)$ and $\alpha_0 = \sum_{i=1}^d \alpha_i$ is a conjugate prior to the multinomial likelihood $\text{Binom}(\mathbf{k}|n, \mathbf{p}) = n! \prod_{i=1}^d \frac{p_i^{k_i}}{k_i!}$, where \mathbf{k} is a vector of counts (k_1, k_2, \dots, k_d) with $\sum_{i=1}^d k_i = n$. Derive update equations for the parameters $\boldsymbol{\alpha}'$ of the posterior distribution after observing counts \mathbf{k} . (5 marks)

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

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

The updated equation is thus $\boldsymbol{\alpha}' = \boldsymbol{\alpha} + \mathbf{k}$.

- (g) Prove that the set of positive semi-definite matrices is a convex set. (5 marks)

For any two PSD matrices $M_1, M_2 \succeq 0$ the convex combination $M_3 = a M_1 + (1 - a) M_2$ (where $a \in [0, 1]$) has a quadratic form with an arbitrary vector v of

$$v^T M_3 v = v^T M_1 v + v^T M_2 v$$

But as $M_1, M_2 \succeq 0$ then $v^T M_1 v \geq 0$ and $v^T M_2 v \geq 0$. Since $a, (1 - a) \geq 0$ it follows that $v^T M_3 v \geq 0$ or $M_3 \succeq 0$. Thus all convex combinations of positive semi-definite matrices are positive semi-definite so they form a convex set.

- (h) Explain what are the hyper-parameters of a Gaussian process and why they are relatively easy to learn. (5 marks)

The hyper-parameters are the mean function $m(x)$ and the kernel or covariance $k(x, y)$ (often the kernel will depend on a scale ℓ that is a hyper-parameter). It is relatively straightforward to find good hyper-parameters as we can compute the evidence $\mathbb{P}(\mathcal{D})$ in closed form.

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 $\frac{1}{40}$
-------------------	-------------------	-------------------	-------------------	-------------------	-------------------	-------------------	-------------------	----------------------

Section B

Question B 2

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

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

show

$$\left\langle \left(\frac{1}{n} \sum_{i=1}^n X_i - \mu \right)^2 \right\rangle = \rho \sigma^2 + \frac{(1 - \rho) \sigma^2}{n}$$

(10 marks)

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

10

- (b) Using the result derived in part (a) explain why ensembling many machines can reduce the variance in the bias-variance dilemma if the machine predictions are not heavily correlated. Use this to explain the success of random forest. (10 marks)

We can think of X_i as the prediction made by a learning machine. For a finite training set this prediction will vary. If our machines are unbiased $\langle X \rangle_i = \mu$, that is the true output. ρ is a measure of the Pearson correlation between the predictions. By ensembling many machines (high n) we reduce the second term. However, because the machines are correlated our prediction will still vary from the true prediction by $\rho \sigma^2$. In random forest we try to make the predictions uncorrelated by basing them on a different set of features and using bootstrapping so they learn on a slightly different training set.

End of question B2

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

$\frac{10}{10}$

Question B 3

(a) We can write the loss for ridge regression as

$$L(w) = \|Xw - y\|^2 + \eta \|w\|^2$$

where X is the design matrix and y is a vector of target values. Calculate the weights, w^* , that minimise the loss. By writing $X = USV^T$ show that we can write $w^* = V\hat{S}^+U^T y$, where the elements of \hat{S}^+ are zero everywhere except for the diagonal where $\hat{S}_{ii}^+ = s_i/(s_i^2 + \eta)$ with s_i being the singular values of X (i.e. $s_i = S_{ii}$). (15 marks)

The gradient of the loss is given by

$$\nabla L(w) = 2X^T X w - 2X^T y + 2\eta w.$$

Setting this to zero and solving for w

$$w^* = (X^T X + \eta I)^{-1} X^T y.$$

Writing $X = USV^T$ then

$$X^T X + \eta I = V (S^T S + \eta I) V^T.$$

So that

$$w^* = V (S^T S + \eta I)^{-1} S U^T y$$

(note that we used the fact that V and U are orthogonal matrices so that $VV^T = I$ and $U^T U = I$).

Thus, $w^* = V\hat{S}^+U^T y$ where

$$\hat{S}^+ = (S^T S + \eta I)^{-1} S.$$

Since the elements of S are zero everywhere except on the diagonal then $S^T S$ will also be zero everywhere except for the diagonal which are s_i^2 and consequently \hat{S}^+ will be zero everywhere except on the diagonal where $\hat{S}_{ii}^+ = s_i/(s_i^2 + \eta)$.

- (b) Use the result that you derive to explain how adding the L_2 regulariser $\eta \|w\|^2$ improves the conditioning of the solution and is likely to improve generalisation. (5 marks)

The effect of adding a regulariser is to change the usual pseudo-inverse ($S_{ii}^+ = s_i^{-1}$) to $\hat{S}_{ii}^+ = s_i / (s_i^2 + \eta)$. If for any singular value $s_i = 0$ the problem would be ill-posed as there would be an infinity of solutions—the pseudo-inverse is not defined. The problem now has a unique solution. That is, $\hat{S}_{ii}^+ = 0$ if $s_i = 0$. Similarly if s_i is small $\hat{S}_{ii}^+ = s_i^{-1}$ would be very large and the pseudo-inverse poorly conditioned (very sensitive to the training data—i.e. with a large variance). However, if $s_i < \eta$ then $\hat{S}_{ii}^+ < 1$ and the regularised solution will be much better conditioned.

End of question B3

(a) $\frac{\quad}{15}$	(b) $\frac{\quad}{5}$	Total $\frac{\quad}{20}$
------------------------	-----------------------	--------------------------

5

Question B 4

- (a) Write a Lagrangian for the linear programming problem of choosing x to minimise $c^T x$, subject to the constraints $Mx = b$. Show that the Lagrangian can be rewritten to obtain the dual problem where the roles of the variables x and the Lagrange multipliers are exchanged. Write down the dual problem as a maximisation problem plus a new set of constraints. (5 marks)

The Lagrangian will be

$$\mathcal{L}(x, \alpha) = c^T x - \alpha^T (Mx - b)$$

where α are a set of Lagrange multipliers. Note that the optimum is given by minimising with respect to x , but maximising with respect to α . We can rewrite this as

$$\mathcal{L}(x, \alpha) = b^T \alpha - x^T (M^T \alpha - c)$$

which we can interpret as a linear programming problem where we choose α to maximise $b^T \alpha$ subject to the constraints $M^T \alpha = c$.

- (b) Describe the Wasserstein distance as a linear programming problem and describe the dual problem. Describe how this is used in the Wasserstein GAN. (15 marks)

(Not all details are necessary for full marks.)

We can pose the Wasserstein distance between two probability densities $p(x)$ and $q(x)$ as a minimisation problem were

$$W(p, q) = \min_{\gamma(x, y) \in \Lambda} \int \gamma(x, y) d(x, y) dx dy$$

where $d(x, y)$ is the Euclidean distance between x and y and Λ is the set of joint probability distributions with marginal given by

$$\int \gamma(x, y) dy = p(x), \quad \int \gamma(x, y) dx = q(y).$$

This can be seen as a linear programme of finding $\gamma(x, y)$ that minimises a linear objective function (the integral in the definition of $W(p, q)$), subject to constraints (given above). Note that $\gamma(x, y)$ can be consider a transportation policy transferring probability mass from p at a point x to q at a point y .

The Lagrangian can be written as

$$\begin{aligned} \mathcal{L} = & \int \gamma(x, y) d(x, y) dx dy - \int \alpha(x) \left(\int \gamma(x, y) dy - p(x) \right) dx \\ & - \int \beta(y) \left(\int \gamma(x, y) dx - q(y) \right) dy \end{aligned}$$

15

where $\alpha(x)$ and $\beta(y)$ are Lagrange multipliers and $\gamma(x, y) \geq 0$. This can also be written

$$\mathcal{L} = \int \alpha(x) p(x) dx + \int \beta(y) q(y) dy + \int \gamma(x, y) (d(x, y) - \alpha(x) - \beta(y)) dx dy$$

Leading to a dual problem of finding $\alpha(x)$ and $\beta(y)$ that maximises $\int \alpha(x) p(x) dx + \int \beta(y) q(y) dy$ subject to the constraint $\alpha(x) + \beta(y) \leq d(x, y)$. Since $d(x, x) = 0$ we have that $\beta(x) = -\alpha(x)$. Thus we have to find a function $\alpha(x)$ that maximises $\int \alpha(x) (p(x) - q(x)) dx$ subject to the constraint $\alpha(x) - \alpha(y) \leq d(x, y)$, that is, $\alpha(x)$ is Lipschitz-1. In a Wasserstein GAN we learn a function $\alpha(x)$ that discriminate between samples, x , drawn from the true distribution $p(x)$ and samples, y , drawn from a generator with distribution $g(y)$. It does this by maximising $\int \alpha(x) (p(x) - g(x)) dx$. This is our critic or discriminator network, but it should be Lipschitz-1. We also learn a generator that tries to minimise this objective function (decreasing the Wasserstein distance).

End of question B4

(a) $\frac{5}{5}$ (b) $\frac{15}{15}$ Total $\frac{20}{20}$

END OF PAPER