

SEMESTER 2 EXAMINATION 2020 - 2021

ADVANCED MACHINE LEARNING

DURATION 120 MINS (2 Hours)

---

This paper contains 4 questions

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

An outline marking scheme is shown in brackets to the right of each question.

This examination is worth 80%. The coursework was worth 20%.

This is a take home examination. Detailed instructions are given on line.
---

University approved calculators MAY be used.

A foreign language dictionary is permitted ONLY IF it is a paper version of a direct 'Word to Word' translation dictionary AND it contains no notes, additions or annotations.

9 page examination paper.

## Section A

### Question A1.

- (a) Explain what it means to normalise the input features and why this is often a good thing to do.

Indicative Solution for Question A1(a).

---

We normalise the features so they all have the same mean and variance. It is often useful as many learning machines will treat larger features as more important. If we don't know which features are important normalising them will remove an unjustified bias.

---

[5 marks]

- (b) What is the purpose of cross-validation and what is its cost?

Indicative Solution for Question A1(b).

---

Cross-validation allows us to get a more accurate approximation of the generalisation performance while using a large proportion of examples to train the weights. The cost is that we have to train our learning machine multiple times.

---

[5 marks]

- (c) To find the maximum or minimum of a function explain the major ways (i) gradient descent (ii) Newton's method and (iii) quasi-Newton methods differ in terms of the information they use.

Indicative Solution for Question A1(c).

- 
- (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 marks]

- (d) Explain why CNNs capture the structure of typical image datasets.

Indicative Solution for Question A1(d).

---

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). Objects tend to be localised.

---

[5 marks]

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

Indicative Solution for Question A1(e).

---

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 boundary.

---

[5 marks]

- (f) Show that an empirical covariance matrix,  $\mathbf{C}$ , can be written as  $\mathbf{X}\mathbf{X}^\top$  and hence prove that its eigenvalues are non-negative.

Indicative Solution for Question A1(f).

---

The empirical covariance matrix is given by

$$\mathbf{C} = \frac{1}{m-1} \sum_{i=1}^m (\mathbf{X}_i - \hat{\boldsymbol{\mu}})(\mathbf{X}_i - \hat{\boldsymbol{\mu}})^\top$$

where  $\hat{\boldsymbol{\mu}}$  is the empirical mean vector. Defining  $\mathbf{X}$  as a matrix whose  $i^{th}$  column is equal to  $(\mathbf{X}_i - \hat{\boldsymbol{\mu}})/\sqrt{m-1}$  then  $\mathbf{C} = \mathbf{X}\mathbf{X}^\top$ . If  $\mathbf{v}$  is an eigenvector of  $\mathbf{C}$  then  $\mathbf{C}\mathbf{v} = \lambda\mathbf{v}$ . Multiplying on the left by  $\mathbf{v}^\top$

$$\mathbf{v}^\top \mathbf{C} \mathbf{v} = \mathbf{v}^\top \mathbf{X} \mathbf{X}^\top \mathbf{v} = \lambda \mathbf{v}^\top \mathbf{v}.$$

Defining  $\mathbf{u} = \mathbf{X}^\top \mathbf{v}$  then

$$\lambda = \frac{\mathbf{u}^\top \mathbf{u}}{\mathbf{v}^\top \mathbf{v}} = \frac{\|\mathbf{u}\|^2}{\|\mathbf{v}\|^2} \geq 0$$

[5 marks]

- (g) Prove that  $f(x) = \exp(cx)$  is a convex-up function.

Indicative Solution for Question A1(g).

---

To prove this we need to show  $f''(x) \geq 0$  for all  $x$ . But  $f''(x) = c^2 \exp(cx)$ . Since  $c^2 \geq 0$  and  $\exp(cx) > 0$  it follows that  $f(x)$  is convex-up.

---

[5 marks]

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

Indicative Solution for Question A1(h).

---

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

---

**TURN OVER**

Answers

## Section B

### Question B1.

(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

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

Indicative Solution for Question B1(a).

---

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


---

[10 marks]

(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.

Indicative Solution for Question B1(b).

---

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  $\mathbb{E}[X_i] = \mu$ , that is, the true value we are trying to predict. The constant  $\rho$  is 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.

---

**TURN OVER**

Answers

**Question B2.**

(a) Consider gradient descent

$$x^{(t+1)} = x^{(t)} - r f'(x^{(t)})$$

acting on a function  $f(x) = x^2$ . Explain what happens starting from some point  $x^{(0)}$  if (i)  $0 < r < \frac{1}{2}$ , (ii)  $r = \frac{1}{2}$ , (iii)  $\frac{1}{2} < r < 1$ , (iv)  $r = 1$  and (v)  $r > 1$ .

Indicative Solution for Question B2(a).

---

We note that

$$x^{(t+1)} = (1 - 2r) x^{(t)}$$

- (i)  $0 < r < \frac{1}{2}$ :  $x^{(t)}$  **converges monotonically towards the optimum at  $x = 0$ .**
  - (ii)  $r = \frac{1}{2}$ : **we jump to the optimum in one step**
  - (iii)  $\frac{1}{2} < r < 1$ :  $x^{(t)}$  **jumps to the other side of the optimum but converges towards it**
  - (iv)  $r = 1$ :  $x^{(t)} = (-1)^t x^{(0)}$
  - (v)  $r > 1$ :  $x^{(t)}$  **diverges from the optimum.**
- 

[10 marks]

(b) Describe the problems that can arise in finding the optimum of a high-dimensional loss function and solution to them.

Indicative Solution for Question B2(b).

---

Solutions might discuss the following

- **Generally we will not be in a convex landscape so we may have to travel a long way to the optimum. We can follow the gradient.**
  - **There will typically be a large range of curvatures meaning there is no clear step size.**
    - **We could use second order methods but in very high-dimensional spaces this is often not practical.**
    - **We can seek to adapt the step size on different variables although this doesn't necessarily help if variables have correlated effects on the loss.**
    - **We can use momentum.**
- 

[10 marks]

**TURN OVER**

**Question B3.**

- (a) Explain how Markov chain Monte-Carlo (MCMC) techniques are used to solve Bayesian inference problems and what problem do they solve.

Indicative Solution for Question B3(a).

---

**MCMC is used to obtain random independent samples from the posterior  $f(\theta|\mathcal{D})$ . In Bayesian inference**

$$f(\theta|\mathcal{D}) = \frac{f(\mathcal{D}|\theta) f(\theta)}{f(\mathcal{D})}.$$

Although we can compute the likelihood,  $f(\mathcal{D}|\theta)$  and prior  $f(\theta)$  the posterior marginal likelihood or evidence,  $f(\mathcal{D})$ , is often intractable. In MCMC we only care about ratios of the joint probability  $f(\mathcal{D}, \theta) = f(\mathcal{D}|\theta) f(\theta)$  which we can compute. Usually this joint probability is too complicated to sample using transformation methods and the rejection rate would be far too high to use rejection methods as most parameters have very low likelihoods. MCMC allow us to transition between likely states, although they require many transitions for the samples to be independent.

---

[10 marks]

- (b) Describe the Metropolis algorithm and show that it satisfies detailed balance. Explain why this is important.

Indicative Solution for Question B3(b).

---

**The Metropolis algorithm is an iterative algorithm where we change the parameters  $\theta$  of our model. We choose a random parameters  $\theta' \sim p(\theta'|\theta)$  where  $p(\theta'|\theta) = p(\theta|\theta')$  (we are using Metropolis not Metropolis-Hastings). We update our parameter if**

$$\frac{f(\theta', \mathcal{D})}{f(\theta, \mathcal{D})} > U$$

where  $U$  is a uniform random variate  $U \sim U(0, 1)$ . That is, the transition probability from  $\theta$  to  $\theta'$  is

$$W(\theta', \theta) = p(\theta'|\theta) \min\left(\frac{f(\theta', \mathcal{D})}{f(\theta, \mathcal{D})}, 1\right).$$

We want to converge to the posterior  $f(\theta|\mathcal{D})$ . This will happen (assuming ergodicity) if we satisfy the detailed balance equation

$$W(\theta', \theta) f(\theta|\mathcal{D}) = W(\theta, \theta') f(\theta'|\mathcal{D})$$

or

$$\frac{W(\theta', \theta)}{W(\theta, \theta')} = \frac{f(\theta'|\mathcal{D})}{f(\theta|\mathcal{D})} = \frac{f(\theta', \mathcal{D})}{f(\theta, \mathcal{D})}.$$

But by the definition of  $W(\theta', \theta)$  we see that if  $f(\theta', \mathcal{D}) \leq f(\theta, \mathcal{D})$

$$\frac{W(\theta', \theta)}{W(\theta, \theta')} = \frac{\frac{p(\theta'|\theta) f(\theta', \mathcal{D})}{f(\theta, \mathcal{D})}}{p(\theta|\theta') \times 1} = \frac{f(\theta', \mathcal{D})}{f(\theta, \mathcal{D})}$$



since by construction  $p(\theta'|\theta) = p(\theta|\theta')$ . We get an identical result if  $f(\theta', \mathcal{D}) > f(\theta, \mathcal{D})$ .  
Satisfying detailed balance ensure that eventually we converge to the posterior distribution.

---

[10 marks]

Answers

**END OF PAPER**