



University College Dublin
An Coláiste Ollscoile, Baile Átha Cliath

SPRING, 22/23 TRIMESTER EXAMINATIONS

COMP47590

Advanced Machine Learning

Module Coordinator: Assoc Professor Brian Mac Namee

Student Number

--	--	--	--	--	--	--	--

Seat Number

--	--	--	--

Time Allowed: 120 minutes

Materials Permitted in the Exam Venue:

Non-programmable or scientific calculator

Programmable calculator

Materials to be Supplied to Students:

8 Page Answer Booklets

New Cambridge Statistical Tables

Instructions to Students:

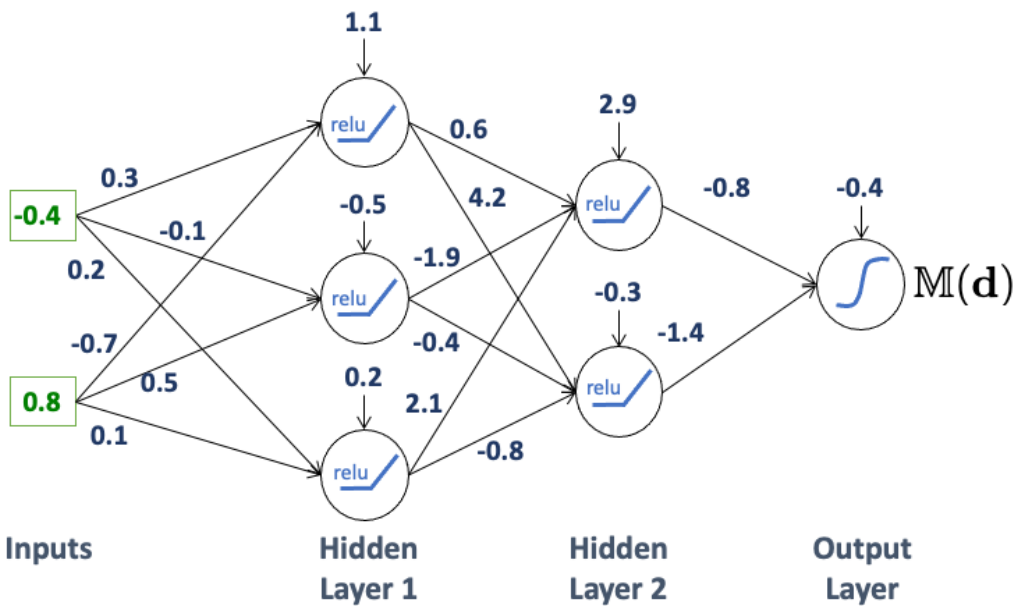
Answer any three out of four questions. All questions carry equal marks. Total marks available 90. The value of each part of each question is shown in brackets next to it.

SOLUTIONS

SOLUTIONS

SOLUTIONS

SOLUTIONS

1.	(a)	<p>The image below shows a <i>feed forward artificial network</i>. The computational units in the two hidden layers use <i>rectified linear (relu)</i> activation functions and the output layer unit uses a <i>sigmoid</i> activation function. The <i>weights</i> and <i>biases</i> are shown along the links in the network.</p>  <p style="text-align: center;"> Inputs Hidden Layer 1 Hidden Layer 2 Output Layer </p>
	(i)	<p>Perform a forward propagation through the network using an input feature vector of $[-0.4, 0.8]$. Show your workings.</p>
		[12 marks]
		<h2 style="text-align: center;">2 Setup Input, Weight and Bias Matrices</h2> <p>The network inputs:</p> $\mathbf{d} = \begin{bmatrix} -0.4 \\ 0.8 \end{bmatrix}$ <p>Weights and biases for Layer 1:</p> $\mathbf{W}^{[1]} = \begin{bmatrix} 0.3 & -0.7 \\ -0.1 & 0.5 \\ 0.2 & 0.1 \end{bmatrix}$ $\mathbf{b}^{[1]} = \begin{bmatrix} 1.1 \\ -0.5 \\ 0.2 \end{bmatrix}$ <p>Weights and biases for Layer 2:</p> $\mathbf{W}^{[2]} = \begin{bmatrix} 0.6 & -1.9 & 2.1 \\ 4.2 & -0.4 & -0.8 \end{bmatrix}$ $\mathbf{b}^{[2]} = \begin{bmatrix} 2.9 \\ -0.3 \end{bmatrix}$

Weights and biases for Layer 3:

$$\mathbf{W}^{[3]} = \begin{bmatrix} -0.8 & -1.4 \end{bmatrix}$$

$$\mathbf{b}^{[3]} = \begin{bmatrix} -0.4 \end{bmatrix}$$

3 Forward Propagate

To perform a forward propagation for the first layer in the network, first calculate $\mathbf{z}^{[1]}$:

$$\begin{aligned}\mathbf{z}^{[1]} &= \mathbf{W}^{[1]} \mathbf{d} + \mathbf{b}^{[1]} \\ &= \begin{bmatrix} 0.3 & -0.7 \\ -0.1 & 0.5 \\ 0.2 & 0.1 \end{bmatrix} \begin{bmatrix} -0.4 \\ 0.8 \end{bmatrix} + \begin{bmatrix} 1.1 \\ -0.5 \\ 0.2 \end{bmatrix} \\ &= \begin{bmatrix} 0.42 \\ -0.06 \\ 0.2 \end{bmatrix}\end{aligned}$$

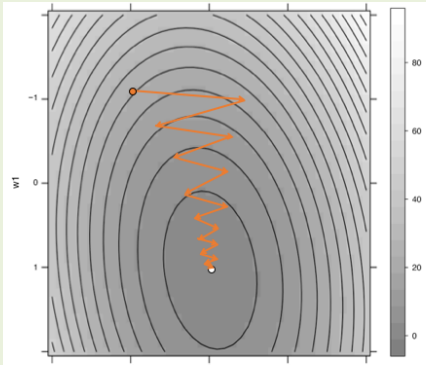
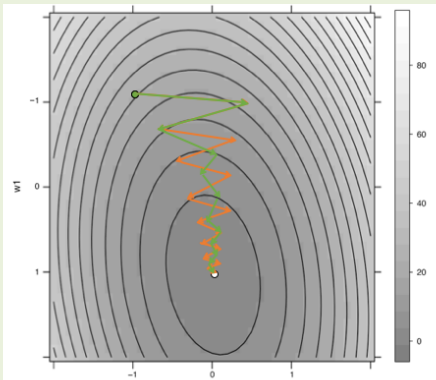
then apply the activation function, in this case a relu function, to calculate the activation of the nodes at Layer 1:

$$\begin{aligned}\mathbf{a}^{[1]} &= g(\mathbf{z}^{[1]}) \\ &= g\left(\begin{bmatrix} 0.42 \\ -0.06 \\ 0.2 \end{bmatrix}\right) \\ &= \begin{bmatrix} 0.42 \\ 0.0 \\ 0.2 \end{bmatrix}\end{aligned}$$

To perform a forward propagation for the second layer in the network, first calculate $\mathbf{z}^{[2]}$:

$$\begin{aligned}\mathbf{z}^{[2]} &= \mathbf{W}^{[2]} \mathbf{a}^{[1]} + \mathbf{b}^{[2]} \\ &= \begin{bmatrix} 0.6 & -1.9 & 2.1 \\ 4.2 & -0.4 & -0.8 \end{bmatrix} \begin{bmatrix} 0.42 \\ 0.0 \\ 0.2 \end{bmatrix} + \begin{bmatrix} 2.9 \\ -0.3 \end{bmatrix} \\ &= \begin{bmatrix} 3.572 \\ 1.304 \end{bmatrix}\end{aligned}$$

			<p>then apply the activation function, in this case a $\text{extbf}\{\text{relu}\}$, to calculate the activation of the output nodes of the network:</p> $\begin{aligned} \mathbf{a}^{[2]} &= g(\mathbf{z}^{[2]}) \\ &= g\left(\begin{bmatrix} 3.572 \\ 1.304 \end{bmatrix}\right) \\ &= \begin{bmatrix} 3.572 \\ 1.304 \end{bmatrix} \end{aligned}$ <p>To perform a forward propagation for the third layer in the network, first calculate $\mathbf{z}^{[3]}$:</p> $\begin{aligned} \mathbf{z}^{[3]} &= \mathbf{W}^{[3]} \mathbf{a}^{[2]} + \mathbf{b}^{[3]} \\ &= \begin{bmatrix} -0.8 & -1.4 \end{bmatrix} \begin{bmatrix} 3.572 \\ 1.304 \end{bmatrix} + \begin{bmatrix} -0.4 \end{bmatrix} \\ &= \begin{bmatrix} -5.083 \end{bmatrix} \end{aligned}$ <p>then apply the activation function, in this case a $\text{extbf}\{\text{softmax function}\}$, to calculate the activation of the output nodes of the network:</p> $\begin{aligned} \mathbf{a}^{[3]} &= g(\mathbf{z}^{[3]}) \\ &= g\left(\begin{bmatrix} -5.083 \end{bmatrix}\right) \\ &= \begin{bmatrix} 0.006 \end{bmatrix} \end{aligned}$ <p>D: Half an attempt D+: First layer matrices set up correctly C: First layer activation calculated B-: Matrices set up wrong, but calculations made B: Set up correctly but failed to apply relu B: Correct workings but did not work out final sigmoid or some other significant calculation error (e.g. put matrices in wrong order in calculation) B+: Small calculation error but correct workings A: No, or minimal, workings A+: Correct answer and well formatted workings.</p>
		(ii)	<p>If the target feature value for the current input vector is 1.0, calculate the loss associated with this training instance using cross entropy loss.</p>
			[2 marks]
			<p>Calculate cross entropy loss</p> $\begin{aligned} \text{Loss} &= -(1 \cdot \log(0.00616) + 0 \cdot \log(0.99384)) \\ &= 5.0894 \end{aligned}$ <p>E: stated log loss equation D+: Incorrectly applied log loss C+: Correctly applied log loss but calculation error B+: log10 instead natural log</p>

			<p>B: Wrong sign</p> <p>A+: Correctly applied log loss function to result from previous section.</p>
	(b)	<p>Gradient descent with momentum, RMSprop, and adam are three common adaptations to the basic gradient descent algorithm used to train neural networks. Explain how these approaches improve upon basic gradient descent and how they differ from each other.</p>	
		[8 marks]	
		<p><u>Sample Answer</u></p> <p>(Diagrams such as those included in this answer would be useful.)</p> <p>The gradient descent algorithm optimizes network weight values during a journey across an error (or loss) surface.</p>  <p>All of the adaptations listed improve this process by taking a more direct route across the error surface. This is achieved by adding different types of momentum terms.</p> <p>Gradient descent with momentum adds a momentum term to the gradient descent process to avoid sudden changes in exploration of the error surface. A diagram such as the following would help.</p>  <p>This is achieved by using an exponentially weighted moving average applied to the gradient terms over subsequent iterations.</p>	

$$v_{d\mathbf{W}} = \beta v_{d\mathbf{W}} + (1 - \beta) d\mathbf{W}$$

$$v_{d\mathbf{b}} = \beta v_{d\mathbf{b}} + (1 - \beta) d\mathbf{b}$$

Weight and bias terms are then updated with these averaged gradients rather than raw gradient values.

$$\mathbf{W} = \mathbf{W} - \alpha v_{d\mathbf{W}}$$

$$\mathbf{b} = \mathbf{b} - \alpha v_{d\mathbf{b}}$$

RMSprop builds on top of the same idea as gradient descent with momentum but uses the square of the gradients in the exponentially weighted moving average.

$$v_{d\mathbf{W}} = \beta v_{d\mathbf{W}} + (1 - \beta) d\mathbf{W}^2$$

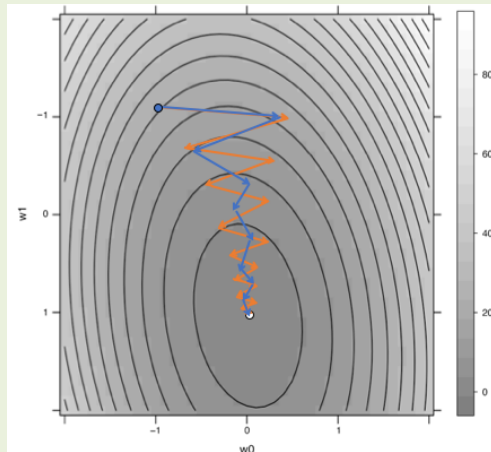
$$v_{d\mathbf{b}} = \beta v_{d\mathbf{b}} + (1 - \beta) d\mathbf{b}^2$$

These values are then used in the weight and bias update rules.

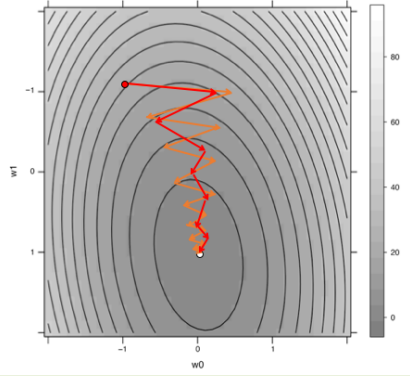
$$\mathbf{W} = \mathbf{W} - \alpha \frac{d\mathbf{W}}{\sqrt{v_{d\mathbf{W}}} + \epsilon}$$

$$\mathbf{b} = \mathbf{b} - \alpha \frac{d\mathbf{b}}{\sqrt{v_{d\mathbf{b}}} + \epsilon}$$

This gives a more direct route across the error surface.



Adam mixes the gradient descent with momentum and RMSprop approaches in a weighted sum:

		$v_{d\mathbf{W}} = \beta_1 v_{d\mathbf{W}} + (1 - \beta_1) d\mathbf{W}$ $v_{d\mathbf{b}} = \beta_1 v_{d\mathbf{b}} + (1 - \beta_1) d\mathbf{b}$ $s_{d\mathbf{W}} = \beta_2 s_{d\mathbf{W}} + (1 - \beta_2) d\mathbf{W}^{*2}$ $s_{d\mathbf{b}} = \beta_2 s_{d\mathbf{b}} + (1 - \beta_2) d\mathbf{b}^{*2}$ <p style="text-align: right; color: red;">Gradient Descent with Momentum</p> <p style="text-align: right; color: red;">RMSprop</p> $\mathbf{W} = \mathbf{W} - \alpha \frac{v_{d\mathbf{W}}}{\sqrt{s_{d\mathbf{W}}} + \epsilon}$ $\mathbf{b} = \mathbf{b} - \alpha \frac{v_{d\mathbf{b}}}{\sqrt{s_{d\mathbf{b}}} + \epsilon}$ <p>This gives a very efficient route across an error surface.</p>  <p>E: Explaaintion given, but doesn't really match what these approaches do.</p> <p>D+: very minimal explanation of three approaches</p> <p>C+: solid exp[lanation of three approaches (no equations)</p> <p>B: solid explantion with equations</p> <p>A: Correct explanation of three approaches and relevant advantages/disadvantages/insight for each.</p>	
	(c)	<p>We can describe what artificial neural networks do as learning representations (or embeddings) of input data that make downstream tasks, for example classification, straight forward. Good embeddings are often said to be:</p> <ul style="list-style-type: none"> • Distributed • Abstract • Invariant • Disentangled 	

		Describe what each of these terms means in relation to embeddings generated by artificial neural networks.
		[8 marks]
		<p>A short description of each term along the following lines is required.</p> <ul style="list-style-type: none"> - Distributed: Expressed across multiple features. - Abstract: Can capture abstract concepts. - Invariant: Are invariant to small and local changes in input data. - Disentangled: Are not overly connected to a single task. <p style="text-align: right;">[2 marks each]</p> <p>E: Explanation given but doesn't really match what these approaches do.</p> <p>D+: very minimal explanation of three approaches</p> <p>C+: solid explanation of four ideas</p> <p>B+: Really good explanation of three ideas and why they are important</p> <p>A+: Correct explanation of four ideas and relevant advantages/disadvantages/insight for each.</p>

2.	(a)	<p>You have been tasked with training a neural network to control a self-driving car from image input. The model should output one of four control signals - <i>left</i>, <i>right</i>, <i>brake</i>, or <i>accelerate</i> - from each input image frame. The only input to the model is a 224 pixel by 224 pixel greyscale image from the front of the car.</p> <p>Image (a) shows the architecture of a multi-layer perceptron neural network designed for this problem. Image (b) shows the architecture of a convolutional neural network designed for this problem. Both architectures are composed of four layers.</p> <div data-bbox="494 571 1260 940"> </div> <p>(a) Multi-layer perceptron network architecture</p> <div data-bbox="359 1030 1420 1422"> </div> <p>(b) Convolutional neural network architecture</p> <p>Calculate the number of parameters (weights and biases) that need to be learned for each network architecture.</p>
		[12 marks]
		<p><u>Sample Answer</u></p> <p>Multi-layer perceptron:</p> <p>This is pretty straight-forward as it involves multiplying through the sizes of the network layers.</p> <p>Input: $224 * 224 = 50,176$</p> <p>Layer 1: $50,176 * 8,000 + 8,000 = \mathbf{401,416,000}$</p>

	<p> Layer 2: $8,000 * 3,000 + 3,000 = \mathbf{24,003,000}$ Layer 3: $3,000 * 800 + 800 = \mathbf{2,400,800}$ Layer 4: $800 * 4 + 4 = \mathbf{3,204}$ Total parameters: $\mathbf{427,823,004}$ </p> <p> C: Wrong input size. B: Missed a layer B+-: right workings but calculation mistake A+: Right answer </p> <p> Convolutional neural network: Students need determine the number of weights based on the size of each filter and the size of each layer. To calculate the number of activations at the flattening layer they also need to keep track of the number of activations flowing through the network. </p> <p> Layer 1 Dim: $224 * 224 * 1$ Layer 1: $5 * 5 * 1 * 12 + 12 = \mathbf{312}$ Layer 1 Output Dim: $220 * 220 * 12$ Layer 1 Output Dim After Pooling: $110 * 110 * 12$ </p> <p> Layer 2: $3 * 3 * 12 * 24 + 24 = \mathbf{2,616}$ Layer 2 Output Dim: $108 * 108 * 24$ Layer 2 Output Dim After Pooling: $54 * 54 * 24$ Layer 2 Dimensionality after flattening: $54 * 54 * 24 = \mathbf{69,984}$ </p> <p> Layer 3: $69,984 * 800 + 800 = \mathbf{55,988,000}$ </p> <p> Layer 4: $800 * 4 + 4 = \mathbf{3,204}$ </p> <p> Total parameters: $\mathbf{55,994,132}$ </p>
--	---

		<p>D: used image size not convolutions for size</p> <p>D+: multiple issues</p> <p>C+: Incorrect flattening</p> <p>B: Missed a layer</p> <p>B: FORGOT BIAS TERMS</p> <p>B+: right workings but calculation mistake</p> <p>A+: Right answer</p>
(b)	<p>The image below shows a 3-channel input that is being convolved (cross correlated) with a 3-channel 3 x 3 kernel.</p> <div><div><div><div>11814183212</div><div>2327281915</div><div>11145181411</div><div>13121453917</div><div>17312313316</div></div><div><div>56227811955</div><div>1902152485367</div><div>39254735467</div><div>491999711498</div></div><div><div>1-1-1</div><div>-11-1</div><div>-1-11</div></div><div>*</div><div><div>1-1-1</div><div>-11-1</div><div>-1-11</div></div></div><p>The image below expands the three-channel input and three-channel kernel so that all values can be seen and shows the intermediate convolution result for each channel as well as the final output. Calculate the values marked with a ? in the intermediate convolution results and the final output.</p><div><div><div><div>Channel 1</div><div><div>11814183212</div><div>11145181411</div><div>13121453917</div><div>17312313316</div></div><div><div>1-1-1</div><div>-11-1</div><div>-1-11</div></div><div><div>Channel 1</div><div><div>322? -208</div><div>-323286 -270</div></div></div></div><div><div><div>Channel 2</div><div><div>2327281915</div><div>140145148153167</div><div>3354839294</div><div>949910711098</div></div><div><div>-1-1-1</div><div>111</div><div>-1-1-1</div></div><div><div>Channel 2</div><div><div>185143137</div><div>? -533 -514</div></div></div></div><div><div><div>Channel 3</div><div><div>56227811955</div><div>1902152485367</div><div>39254735467</div><div>491999711498</div></div><div><div>-11-1</div><div>111</div><div>-11-1</div></div><div><div>Channel 3</div><div><div>885116165</div><div>196? -149</div></div></div></div><div><div><div>Final Output</div><div><div>1392-36?</div><div>-690? -933</div></div></div></div></div></div></div></div></div>	
		[8 marks]
		Simple convolutions are calculated for each channel and then summed for the final output.

		<div><div><div>Channel 1</div><table><tr><td>118</td><td>14</td><td>18</td><td>32</td><td>12</td></tr><tr><td>11</td><td>145</td><td>18</td><td>14</td><td>11</td></tr><tr><td>13</td><td>12</td><td>145</td><td>39</td><td>17</td></tr><tr><td>17</td><td>31</td><td>23</td><td>133</td><td>16</td></tr></table></div><div><div>Channel 2</div><table><tr><td>23</td><td>27</td><td>28</td><td>19</td><td>15</td></tr><tr><td>140</td><td>145</td><td>148</td><td>153</td><td>167</td></tr><tr><td>33</td><td>54</td><td>83</td><td>92</td><td>94</td></tr><tr><td>94</td><td>99</td><td>107</td><td>110</td><td>98</td></tr></table></div><div><div>Channel 3</div><table><tr><td>56</td><td>227</td><td>81</td><td>19</td><td>55</td></tr><tr><td>190</td><td>215</td><td>248</td><td>53</td><td>67</td></tr><tr><td>39</td><td>254</td><td>73</td><td>54</td><td>67</td></tr><tr><td>49</td><td>199</td><td>97</td><td>114</td><td>98</td></tr></table></div><div><div>*</div><table><tr><td>1</td><td>-1</td><td>-1</td></tr><tr><td>-1</td><td>1</td><td>-1</td></tr><tr><td>-1</td><td>-1</td><td>1</td></tr></table></div><div><div>*</div><table><tr><td>-1</td><td>-1</td><td>-1</td></tr><tr><td>1</td><td>1</td><td>1</td></tr><tr><td>-1</td><td>-1</td><td>-1</td></tr></table></div><div><div>*</div><table><tr><td>-1</td><td>1</td><td>-1</td></tr><tr><td>1</td><td>1</td><td>1</td></tr><tr><td>-1</td><td>1</td><td>-1</td></tr></table></div><div><div>Channel 1</div><table><tr><td>322</td><td>-295</td><td>208</td></tr><tr><td>-323</td><td>286</td><td>-270</td></tr></table></div><div><div>Channel 2</div><table><tr><td>185</td><td>143</td><td>137</td></tr><tr><td>-563</td><td>-533</td><td>-514</td></tr></table></div><div><div>Channel 3</div><table><tr><td>885</td><td>116</td><td>165</td></tr><tr><td>196</td><td>145</td><td>-149</td></tr></table></div><div><div>Final Output</div><table><tr><td>1392</td><td>-36</td><td>94</td></tr><tr><td>-690</td><td>-102</td><td>-933</td></tr></table></div></div>	118	14	18	32	12	11	145	18	14	11	13	12	145	39	17	17	31	23	133	16	23	27	28	19	15	140	145	148	153	167	33	54	83	92	94	94	99	107	110	98	56	227	81	19	55	190	215	248	53	67	39	254	73	54	67	49	199	97	114	98	1	-1	-1	-1	1	-1	-1	-1	1	-1	-1	-1	1	1	1	-1	-1	-1	-1	1	-1	1	1	1	-1	1	-1	322	-295	208	-323	286	-270	185	143	137	-563	-533	-514	885	116	165	196	145	-149	1392	-36	94	-690	-102	-933
118	14	18	32	12																																																																																																													
11	145	18	14	11																																																																																																													
13	12	145	39	17																																																																																																													
17	31	23	133	16																																																																																																													
23	27	28	19	15																																																																																																													
140	145	148	153	167																																																																																																													
33	54	83	92	94																																																																																																													
94	99	107	110	98																																																																																																													
56	227	81	19	55																																																																																																													
190	215	248	53	67																																																																																																													
39	254	73	54	67																																																																																																													
49	199	97	114	98																																																																																																													
1	-1	-1																																																																																																															
-1	1	-1																																																																																																															
-1	-1	1																																																																																																															
-1	-1	-1																																																																																																															
1	1	1																																																																																																															
-1	-1	-1																																																																																																															
-1	1	-1																																																																																																															
1	1	1																																																																																																															
-1	1	-1																																																																																																															
322	-295	208																																																																																																															
-323	286	-270																																																																																																															
185	143	137																																																																																																															
-563	-533	-514																																																																																																															
885	116	165																																																																																																															
196	145	-149																																																																																																															
1392	-36	94																																																																																																															
-690	-102	-933																																																																																																															
		<p>E: POORLY APPLIED CONVOLUTION.</p> <p>C+: Missing final computation.</p> <p>B+: Simple calculation error.</p> <p>A+: All numbers correct.</p>																																																																																																															
(c)	<p>The 2017 paper “<i>Attention is all you need</i>” by Vaswani et al is now one of the most cited papers in machine learning research. Explain what attention is and how the transformer architecture utilises it.</p>																																																																																																																
		[10 marks]																																																																																																															
		<p><u>Sample Answer</u></p> <p>Attention is a mechanism used in deep learning that allows a model to selectively focus on certain parts of the input during processing. It helps the model to identify the most important features and context of the input data while ignoring irrelevant information.</p> <p>The Transformer architecture is a type of neural network architecture that is specifically designed to handle sequence-to-sequence tasks such as language translation, text summarization, and speech recognition. It replaces the recurrent and convolutional layers commonly used in these tasks with a self-attention mechanism.</p>																																																																																																															

		<p>The self-attention mechanism in the Transformer architecture allows the model to weigh the importance of different parts of the input sequence when making predictions. The self-attention layer computes a weighted sum of the input sequence, where the weights are based on the similarity between each input element and every other element in the sequence.</p> <p>The Transformer architecture utilizes multi-head attention, where the self-attention mechanism is applied multiple times in parallel, each with its own set of weights. This allows the model to capture different relationships between different parts of the input sequence and is particularly effective for capturing long-range dependencies.</p> <p>D+: Basic idea of context-sensitive representation only C: REASONABLE EXPLANATION OF ATTENTION MECHANISM WITH KEY IDEAS. B+: Excellent explanation of what attention is but no motivation for why. A+: Excellent explanation with details and motivation.</p>
--	--	---

3.	(a)	Describe the concept of discounted return that is frequently used in reinforcement learning.
		[4 marks]
		<p>Sample Answer</p> <p>Calculating the expected return from a sequence of actions is key in developing reinforcement learning systems. In a basic formulation of expected return that simply sums rewards, e.g.</p> $G = r_t + r_{t+1} + r_{t+2} + r_{t+3} + \dots + r_e$ <p>expected future rewards are considered to be as valuable as the immediate reward that the agent will receive from taking the next immediate action. Just like we might be more excited about receiving a gift of \$100 today than a promise to receive a gift of \$100 in a year's time, it is reasonable when calculating expected return to pay more attention to the immediate reward we expect to receive from taking the next action, than to the rewards that we expect to receive in 10 or even 100 action's time. This is known as discounted return. We can define discounted return as:</p> $G_\gamma = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots + \gamma^{e-t} r_e$ <p>where γ is a discount factor that can take a value between 0 and 1.</p> <p>F: Explanation of something else. D: Incorrect explanation, but something right. C: Good explanation but no equation. B: Good explanation plus equation A: Excellent explanation plus equation with some insight</p>
	(b)	<p>An intelligent agent trained to play a video game completes an episode and receives the following sequence of rewards over six timesteps:</p> $\{r_0 = 84, r_1 = 98, r_2 = -57, r_3 = -104, r_4 = -96\}$ <p>Compare the discounted returns calculated at time $t = 0$ based on this reward sequence when discounting factors of 0.9 and 0.1 are used.</p>
		[6 marks]
		<p>Sample Answer</p> <p>Students should begin by calculating the discounted returns using:</p> $G_\gamma = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots + \gamma^{e-t} r_e$

		The following tables shows the calculations for discounting factor of 0.9 and 0.1:																																																																																																							
		Time		0	1	2	3	4																																																																																																	
		Reward		84	98	-57	-104	-96																																																																																																	
				1	0.9	0.81	0.729	0.6561																																																																																																	
		Discounted Return	0.9	84	88.2	-46.17	-75.816	-62.9856	-12.7716																																																																																																
				1	0.1	0.01	0.001	0.0001																																																																																																	
		Discounted Return	0.1	84	9.8	-0.57	-0.104	-0.0096	93.1164																																																																																																
		Students should then discuss that with the lower discount factor much less attention is paid to later rewards and so the overall return is much lower.																																																																																																							
		D:																																																																																																							
		D: small mistake																																																																																																							
		C+: Correct answers but little or no insight																																																																																																							
		A: Correct answers plus good insight																																																																																																							
	(c)	To try to better understand the slightly baffling behaviour of her new baby girl, Maria - a scientifically minded new mother - monitored her baby over a period of time, recording her activity at 20 minute intervals. The activity stream looked like this (with time flowing down through the columns):																																																																																																							
		<table><tr><td>0</td><td>SLEEPING</td><td>12</td><td>SLEEPING</td><td>24</td><td>SLEEPING</td><td>36</td><td>CRYING</td></tr><tr><td>1</td><td>CRYING</td><td>13</td><td>SLEEPING</td><td>25</td><td>HAPPY</td><td>37</td><td>HAPPY</td></tr><tr><td>2</td><td>SLEEPING</td><td>14</td><td>SLEEPING</td><td>26</td><td>CRYING</td><td>38</td><td>HAPPY</td></tr><tr><td>3</td><td>SLEEPING</td><td>15</td><td>CRYING</td><td>27</td><td>SLEEPING</td><td>39</td><td>HAPPY</td></tr><tr><td>4</td><td>SLEEPING</td><td>16</td><td>CRYING</td><td>28</td><td>SLEEPING</td><td>40</td><td>HAPPY</td></tr><tr><td>5</td><td>HAPPY</td><td>17</td><td>SLEEPING</td><td>29</td><td>HAPPY</td><td>41</td><td>HAPPY</td></tr><tr><td>6</td><td>HAPPY</td><td>18</td><td>SLEEPING</td><td>30</td><td>HAPPY</td><td>42</td><td>SLEEPING</td></tr><tr><td>7</td><td>HAPPY</td><td>19</td><td>HAPPY</td><td>31</td><td>HAPPY</td><td>43</td><td>SLEEPING</td></tr><tr><td>8</td><td>SLEEPING</td><td>20</td><td>SLEEPING</td><td>32</td><td>HAPPY</td><td>44</td><td>SLEEPING</td></tr><tr><td>9</td><td>SLEEPING</td><td>21</td><td>HAPPY</td><td>33</td><td>HAPPY</td><td>45</td><td>SLEEPING</td></tr><tr><td>10</td><td>SLEEPING</td><td>22</td><td>HAPPY</td><td>34</td><td>HAPPY</td><td>46</td><td>SLEEPING</td></tr><tr><td>11</td><td>SLEEPING</td><td>23</td><td>CRYING</td><td>35</td><td>CRYING</td><td>47</td><td>SLEEPING</td></tr></table>								0	SLEEPING	12	SLEEPING	24	SLEEPING	36	CRYING	1	CRYING	13	SLEEPING	25	HAPPY	37	HAPPY	2	SLEEPING	14	SLEEPING	26	CRYING	38	HAPPY	3	SLEEPING	15	CRYING	27	SLEEPING	39	HAPPY	4	SLEEPING	16	CRYING	28	SLEEPING	40	HAPPY	5	HAPPY	17	SLEEPING	29	HAPPY	41	HAPPY	6	HAPPY	18	SLEEPING	30	HAPPY	42	SLEEPING	7	HAPPY	19	HAPPY	31	HAPPY	43	SLEEPING	8	SLEEPING	20	SLEEPING	32	HAPPY	44	SLEEPING	9	SLEEPING	21	HAPPY	33	HAPPY	45	SLEEPING	10	SLEEPING	22	HAPPY	34	HAPPY	46	SLEEPING	11	SLEEPING	23	CRYING	35	CRYING	47	SLEEPING
0	SLEEPING	12	SLEEPING	24	SLEEPING	36	CRYING																																																																																																		
1	CRYING	13	SLEEPING	25	HAPPY	37	HAPPY																																																																																																		
2	SLEEPING	14	SLEEPING	26	CRYING	38	HAPPY																																																																																																		
3	SLEEPING	15	CRYING	27	SLEEPING	39	HAPPY																																																																																																		
4	SLEEPING	16	CRYING	28	SLEEPING	40	HAPPY																																																																																																		
5	HAPPY	17	SLEEPING	29	HAPPY	41	HAPPY																																																																																																		
6	HAPPY	18	SLEEPING	30	HAPPY	42	SLEEPING																																																																																																		
7	HAPPY	19	HAPPY	31	HAPPY	43	SLEEPING																																																																																																		
8	SLEEPING	20	SLEEPING	32	HAPPY	44	SLEEPING																																																																																																		
9	SLEEPING	21	HAPPY	33	HAPPY	45	SLEEPING																																																																																																		
10	SLEEPING	22	HAPPY	34	HAPPY	46	SLEEPING																																																																																																		
11	SLEEPING	23	CRYING	35	CRYING	47	SLEEPING																																																																																																		
		Maria noticed that her baby could occupy one of three states - HAPPY, CRYING, or SLEEPING - and moved quite freely between them.																																																																																																							
	(i)	Based on the sequence of states given above calculate a transition matrix that gives the probability of moving between each of the three states.																																																																																																							
		[10 marks]																																																																																																							

Sample Answer

The first step to building the transition matrix is to count the frequency of each possible state transition. Working down through the list of states we can count the number of times we move from one state to the next. This gives a transition frequency table:

	SLEEPING	CRYING	HAPPY
SLEEPING	15	2	5
CRYING	4	2	1
HAPPY	3	3	12

By normalising each row in the table (dividing each value by the sum of values in the row) we can calculate the final transition matrix:

	SLEEPING	CRYING	HAPPY
SLEEPING	0.682	0.091	0.227
CRYING	0.571	0.286	0.143
HAPPY	0.167	0.167	0.667

D:

C:

B: Correct counts but not normalised (or incorrect normalisation)

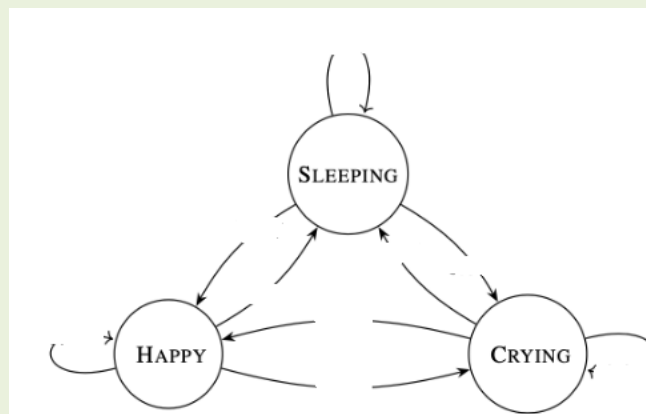
A+: Perfect

- (ii) Draw a Markov process diagram to capture the behaviour of a small baby as described above.

[5 marks]

Sample Answer

A diagram like this one is appropriate for this answer.



		<p>D: Diagram but no numbers</p> <p>C:</p> <p>B:</p> <p>A: Perfect</p>
	(d)	<p>The image below shows an illustration of reinforcement learning using actor-critic method.</p> <p>Describe the role of the actor and critic models in this approach.</p>
		[5 marks]
		<p>Attempt to merge best of policy gradient methods with value function methods. Composed of two models</p> <ul style="list-style-type: none"> - Critic is an action value function model $Q_{M_{WC}}(s_t, a_t)$ - Actor is a policy model $\pi_{M_{WA}}(s_t)$ <p>During learning when the actor model is being updated it uses outputs from the critic action value function to estimate expected return</p>

4.	(a)	<p>Some benchmark experiments have found ensemble models based on gradient boosting can be prone to overfitting to incorrectly labelled instances in the training dataset (for example a positive instance mislabelled as a negative instance). Explain why this is the case and describe how a learning rate can be introduced to the gradient boosting algorithm to mitigate this.</p>
		[10 marks]

		<p><u>Sample Answer</u></p> <p>Gradient boosting builds an ensemble by iteratively training models and building models that correct the predictions of their predecessors. In this way the next model trained will focus on the parts of the training set that the previous model struggled with.</p> <p>The reason that boosting is sensitive to noise in the target features than bagging is that it is prone to over fitting to these noisy instances as models will typically struggle to correctly predict them which means subsequent models will focus too much on them.</p> <p>A learning rate in gradient boosting means that in the aggregation process later models are assigned less weight than earlier models.</p> <p><i>A+ amazing explanations of both.</i></p> <p><i>B+ decent explanation of both problem – model focuses on the error cases. - and how learning rate helps – focuses less on models trained later..</i></p> <p><i>C+ good problem explanation – model focuses on the error cases. No or bad explanation of how learning rate helps – focuses less on models trained later. Or vice versa</i></p> <p><i>D+ hovering around right ideas.</i></p> <p><i>E</i></p> <p><i>NG</i></p>
	(b)	<p>When developing a machine learning model that will be deployed to perform a task for a user, we can describe three different goals of evaluation:</p> <ol style="list-style-type: none"> 1.to determine which model is the most suitable for a task 2.to estimate how the model will perform after deployment 3.to convince users that the model will meet their needs <p>Describe the differences between these goals, and how the evaluation methods used to achieve each of them can be different.</p>
		[10 marks]
		<p><u>Sample Answer</u></p> <p><i>Students should outline that when performing an evaluation to prepare a model for deployment the first goal is to determine which approach might work best. Decisions to be made based on this evaluation include which modelling approach will be used, which data pre-processing techniques might be used, and the optimal algorithm hyper-parameters to be used. This evaluation is very much driven by the machine learning practitioner. Typically, k-fold cross validation can be used as the main evaluation method for this kind of evaluation.</i></p>

		<p><i>Once all of the decisions about what modelling approach to take have been made the next goal attempts to evaluate the likely performance of a model after deployment. This is largely concerned with estimating generalisation error of the model. The only sensible way to evaluate this is to use a hold out test set that has been involved at all in training the model. This is the only reasonable way to evaluate generalisation error.</i></p> <p><i>Once practitioners are convinced that a modelling approach will achieve acceptable performance after deployment the last step is to convince the people for whom the model is being built that it will meet their needs. This can be done with the same kinds of experiments used to evaluate likely generalisation error, but often different performance measures need to be used so as to speak to a different audience (not machine learning practitioners). It is also important to consider a deployment experiment at this stage too.</i></p> <p><i>A+</i></p> <p><i>B+</i></p> <p><i>C+</i></p> <p><i>C: Basic explanation of three approaches.</i></p> <p><i>D+</i></p> <p><i>D Little more than restating the question!</i></p> <p><i>E</i></p> <p><i>NG</i></p>
	(c)	<p>The following is a definition of machine learning:</p> <p><i>The field of machine learning is concerned with the question of how to construct computer programs that automatically improve with experience.</i> - Tom Mitchell</p> <p>Do you believe that this definition accurately defines machine learning? In your answer discuss the appropriateness of the definition, the scope of the definition, and any recommendations for improvements you would suggest.</p>
		[10 marks]
		<u>Sample Answer</u>

	<p>The definition given by Tom Mitchell accurately defines machine learning as a field of study that deals with the development of computer programs that can improve their performance over time with experience. The definition is appropriate as it highlights the central aspect of machine learning, which is learning from data.</p> <p>However, the definition's scope is limited to the concept of "improving with experience" and doesn't account for other aspects of machine learning such as pattern recognition, prediction, and decision-making. These aspects are also critical to machine learning and should be included in the definition.</p> <p>In terms of improvements to the definition, it could be expanded to include the concept of generalization. In machine learning, the ultimate goal is to develop models that can make accurate predictions on new, unseen data.</p> <p>Therefore, a better definition could be: "The field of machine learning is concerned with the development of computer programs that automatically learn patterns from data, generalize these patterns to make predictions on new, unseen data, and improve their performance over time through experience."</p> <p>Overall, the original definition accurately defines machine learning but could benefit from an expansion to include other critical aspects of the field.</p> <p><i>A+ Exceptional discussion of appropriateness, scope , and improvements. Probably bringing in other definitions,.</i></p> <p><i>B+ decent discussion of appropriateness, scope , and improvements.</i></p> <p><i>C +Missing one of appropriateness, scope , and improvements.</i></p> <p><i>D+</i></p> <p><i>E</i></p> <p><i>NG</i></p>
--	--