



University College Dublin
An Coláiste Ollscoile, Baile Átha Cliath

SPRING, 21/22 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

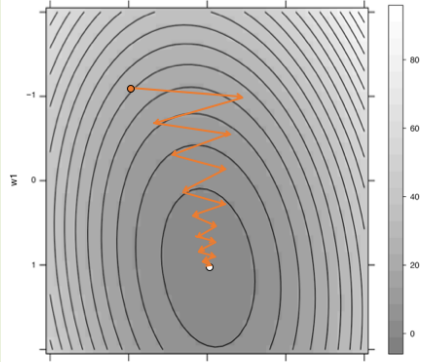
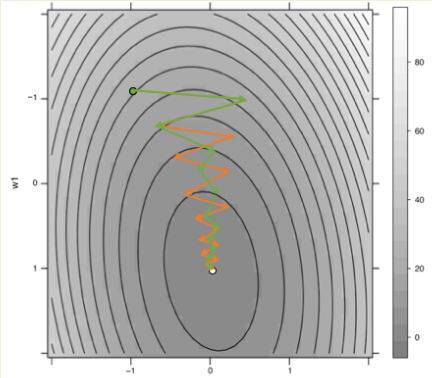
SOLUTIONS

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>softmax</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 $[0.2, -0.9]$. Show your workings.</p>
		[12 marks]
		<p>The network inputs:</p> $\mathbf{d} = \begin{bmatrix} 0.2 \\ -0.9 \end{bmatrix}$ <p>Weights and biases for Layer 1:</p> $\mathbf{W}^{[1]} = \begin{bmatrix} 0.4 & 0.1 \\ 0.4 & -0.5 \\ 0.6 & 0.2 \end{bmatrix}$ $\mathbf{b}^{[1]} = \begin{bmatrix} 1.3 \\ -0.7 \\ 1.2 \end{bmatrix}$ <p>Weights and biases for Layer 2:</p> $\mathbf{W}^{[2]} = \begin{bmatrix} 1.6 & 1.8 & 1.1 \\ 3.2 & -1.4 & -0.1 \end{bmatrix}$ $\mathbf{b}^{[2]} = \begin{bmatrix} -0.2 \\ -1.4 \end{bmatrix}$ <p>Weights and biases for Layer 3:</p> $\mathbf{W}^{[3]} = \begin{bmatrix} -0.1 & 2.4 \\ 2.1 & 1.2 \end{bmatrix}$ $\mathbf{b}^{[3]} = \begin{bmatrix} -0.1 \\ -0.2 \end{bmatrix}$

		<p>To perform a forward propagation for the first layer in the network, first calculate $\text{extbf}\{z\}^{\{1\}}$:</p> $\begin{aligned} z^{[1]} &= \mathbf{W}^{[1]} \mathbf{d} + \mathbf{b}^{[1]} \\ &= \begin{bmatrix} 0.4 & 0.1 \\ 0.1 & -0.5 \\ 0.4 & 0.2 \end{bmatrix} \begin{bmatrix} 0.2 \\ -0.9 \end{bmatrix} + \begin{bmatrix} 1.3 \\ -0.7 \\ 1.2 \end{bmatrix} \\ &= \begin{bmatrix} 1.29 \\ -0.23 \\ 1.1 \end{bmatrix} \end{aligned}$ <p>then apply the activation function, in this case a relu function, to calculate the activation of the nodes at Layer 1:</p> $\begin{aligned} \mathbf{a}^{[1]} &= g(z^{[1]}) \\ &= g\left(\begin{bmatrix} 1.29 \\ -0.23 \\ 1.1 \end{bmatrix}\right) \\ &= \begin{bmatrix} 1.29 \\ 0.0 \\ 1.1 \end{bmatrix} \end{aligned}$ <p>To perform a forward propagation for the second layer in the network, first calculate $\text{extbf}\{z\}^{\{2\}}$:</p> $\begin{aligned} z^{[2]} &= \mathbf{W}^{[2]} \mathbf{a}^{[1]} + \mathbf{b}^{[2]} \\ &= \begin{bmatrix} 1.6 & 1.8 & 1.1 \\ 3.2 & -1.4 & -0.1 \end{bmatrix} \begin{bmatrix} 1.29 \\ 0.0 \\ 1.1 \end{bmatrix} + \begin{bmatrix} -0.2 \\ -1.4 \end{bmatrix} \\ &= \begin{bmatrix} 3.074 \\ 2.618 \end{bmatrix} \end{aligned}$ <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> $\begin{aligned} \mathbf{a}^{[2]} &= g(z^{[2]}) \\ &= g\left(\begin{bmatrix} 3.074 \\ 2.618 \end{bmatrix}\right) \\ &= \begin{bmatrix} 3.074 \\ 2.618 \end{bmatrix} \end{aligned}$ <p>To perform a forward propagation for the third layer in the network, first calculate $\text{extbf}\{z\}^{\{3\}}$:</p> $\begin{aligned} z^{[3]} &= \mathbf{W}^{[3]} \mathbf{a}^{[2]} + \mathbf{b}^{[3]} \\ &= \begin{bmatrix} -0.1 & 2.4 \\ 2.1 & 1.2 \end{bmatrix} \begin{bmatrix} 3.074 \\ 2.618 \end{bmatrix} + \begin{bmatrix} -0.1 \\ -0.2 \end{bmatrix} \\ &= \begin{bmatrix} 5.876 \\ 9.397 \end{bmatrix} \end{aligned}$ <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> $\begin{aligned} \mathbf{a}^{[3]} &= g(z^{[3]}) \\ &= g\left(\begin{bmatrix} 5.876 \\ 9.397 \end{bmatrix}\right) \\ &= \begin{bmatrix} 0.029 \\ 0.971 \end{bmatrix} \end{aligned}$
	(ii)	<p>If the target feature vector for the current input vector is [1.0, 0.0], calculate the loss associated with this training instance using cross entropy loss.</p>
		[2 marks]
		<p>Calculate cross entropy loss</p> <p>Loss = $-(1 \cdot \log(0.0287) + 0 \cdot \log(0.9713))$ = 3.5503</p>
	(b)	<p>Gradient descent with momentum, RMSprop, and adam are three common adaptations to the basic gradient descent algorithm used to train</p>

		neural networks. Explain how these approaches improve upon basic gradient descent and how they differ from each other.
		[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> <div style="text-align: center; margin: 10px 0;"> $v_{d\mathbf{W}} = \beta v_{d\mathbf{W}} + (1 - \beta) d\mathbf{W}$ $v_{d\mathbf{b}} = \beta v_{d\mathbf{b}} + (1 - \beta) d\mathbf{b}$ </div> <p>Weight and bias terms are then updated with these averaged gradients rather than raw gradient values.</p>

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

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

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.

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

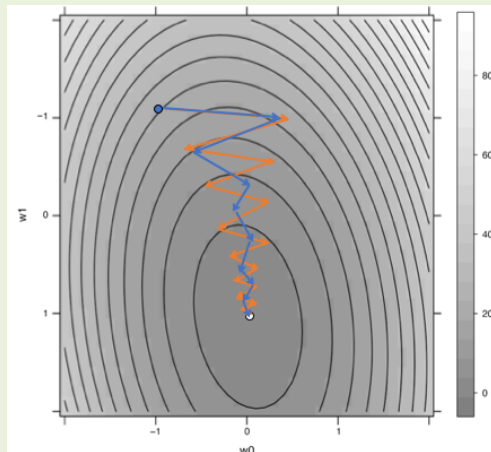
$$v_{d\mathbf{b}} = \beta v_{d\mathbf{b}} + (1 - \beta) d\mathbf{b}^2$$

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:

$$v_{d\mathbf{W}} = \beta_1 v_{d\mathbf{W}} + (1 - \beta_1) d\mathbf{W} \quad \text{Gradient Descent with Momentum}$$

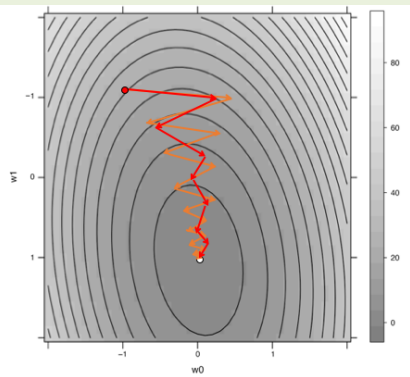
$$v_{d\mathbf{b}} = \beta_1 v_{d\mathbf{b}} + (1 - \beta_1) d\mathbf{b}$$

$$s_{d\mathbf{W}} = \beta_2 s_{d\mathbf{W}} + (1 - \beta_2) d\mathbf{W}^2 \quad \text{RMSprop}$$

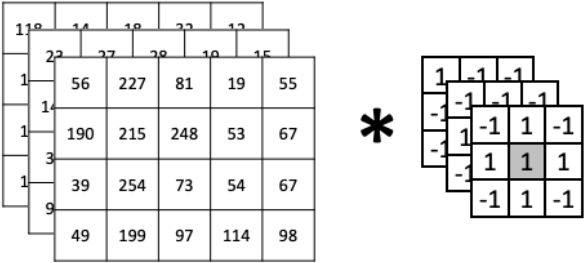
$$s_{d\mathbf{b}} = \beta_2 s_{d\mathbf{b}} + (1 - \beta_2) d\mathbf{b}^2$$

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

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

		<p>This gives a very efficient route across an error surface.</p> 
	(c)	<p>Regularisation can be defined as “<i>any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error</i>” (Goodfellow et al, 2016).</p> <p>Describe two regularisation approaches that can be used when training deep neural networks and explain how they help to reduce generalisation error.</p>
		[8 marks]
		<p>Students can describe any regularisation approach, but these might include:</p> <ul style="list-style-type: none"> • L2 regularisation • Dropout • Data augmentation <p>4 marks awarded for each.</p>

2.	(a)	<p>You have been tasked with training a neural network to diagnose pneumonia from chest x-rays. The model should produce diagnoses on a three-point scale: <i>grade-0</i>= no pneumonia, <i>grade-1</i>= mild pneumonia, <i>grade-2</i>= severe pneumonia. The only input to the model is a 128 pixel by 128 pixel greyscale image of a chest x-ray.</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="507 562 1267 927" data-label="Diagram"> </div> <p>(a) A multi-layer perceptron network for diabetes diagnosis</p> <div data-bbox="379 1021 1394 1400" data-label="Diagram"> </div> <p>(b) A convolutional neural network for diabetes diagnosis</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: $128 * 128 = 16,384$</p> <p>Layer 1: $16,384 * 4,096 + 4,096 = 67,112,960$</p> <p>Layer 2: $4,096 * 1,024 + 1,024 = 4,195,328$</p> <p>Layer 3: $1,024 * 256 + 256 = 262,400$</p>

		<p>Layer 4: $256 * 3 + 3 = 771$</p> <p>Total parameters: 71,571,459</p> <p>Convolutional neural network:</p> <p>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: $128 \times 128 \times 1$</p> <p>Layer 1: $5 * 5 * 1 * 16 + 16 = \mathbf{416}$</p> <p>Layer 1 Output Dim: $124 \times 124 \times 16$</p> <p>Layer 1 Output Dim After Pooling: $62 \times 62 \times 16$</p> <p>Layer 2: $3 * 3 * 16 * 32 + 32 = \mathbf{4,640}$</p> <p>Layer 2 Output Dim: $60 \times 60 \times 32$</p> <p>Layer 2 Output Dim After Pooling: $30 \times 30 \times 32$</p> <p>Layer 2 Dimensionality after flattening: $30 \times 30 \times 32 = \mathbf{28,800}$</p> <p>Layer 3: $28,800 * 716 + 716 = \mathbf{20,621,516}$</p> <p>Layer 4: $716 * 3 + 3 = \mathbf{2,151}$</p> <p>Total parameters: 20,628,723</p>
	(b)	<p>The image below shows a 3 channel input that is being convolved (cross correlated) with a 3 channel kernel.</p> <div style="text-align: center;">  </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>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><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>Channel 1</div><div><table><tr><td>?</td><td>-295</td><td>-208</td></tr><tr><td>-323</td><td>286</td><td>?</td></tr></table></div></div> <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><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>Channel 2</div><div><table><tr><td>185</td><td>143</td><td>137</td></tr><tr><td>-563</td><td>?</td><td>-514</td></tr></table></div></div> <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>Channel 3</div><div><table><tr><td>?</td><td>116</td><td>?</td></tr><tr><td>196</td><td>145</td><td>-149</td></tr></table></div></div> <div><div>Final Output</div><table><tr><td>1392</td><td>?</td><td>?</td></tr><tr><td>?</td><td>?</td><td>-933</td></tr></table></div>	118	14	18	32	12	11	145	18	14	11	13	12	145	39	17	17	31	23	133	16	1	-1	-1	-1	1	-1	-1	-1	1	?	-295	-208	-323	286	?	23	27	28	19	15	140	145	148	153	167	33	54	83	92	94	94	99	107	110	98	-1	-1	-1	1	1	1	-1	-1	-1	185	143	137	-563	?	-514	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	?	116	?	196	145	-149	1392	?	?	?	?	-933
118	14	18	32	12																																																																																																													
11	145	18	14	11																																																																																																													
13	12	145	39	17																																																																																																													
17	31	23	133	16																																																																																																													
1	-1	-1																																																																																																															
-1	1	-1																																																																																																															
-1	-1	1																																																																																																															
?	-295	-208																																																																																																															
-323	286	?																																																																																																															
23	27	28	19	15																																																																																																													
140	145	148	153	167																																																																																																													
33	54	83	92	94																																																																																																													
94	99	107	110	98																																																																																																													
-1	-1	-1																																																																																																															
1	1	1																																																																																																															
-1	-1	-1																																																																																																															
185	143	137																																																																																																															
-563	?	-514																																																																																																															
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																																																																																																															
?	116	?																																																																																																															
196	145	-149																																																																																																															
1392	?	?																																																																																																															
?	?	-933																																																																																																															
		[7 marks]																																																																																																															
		<div>Simple convolutions are calculated for each channel and then summed for the final output.</div> <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><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>Channel 1</div><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><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><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>Channel 2</div><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><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>Channel 3</div><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><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>	118	14	18	32	12	11	145	18	14	11	13	12	145	39	17	17	31	23	133	16	1	-1	-1	-1	1	-1	-1	-1	1	322	-295	-208	-323	286	-270	23	27	28	19	15	140	145	148	153	167	33	54	83	92	94	94	99	107	110	98	-1	-1	-1	1	1	1	-1	-1	-1	185	143	137	-563	-533	-514	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	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																																																																																																													
1	-1	-1																																																																																																															
-1	1	-1																																																																																																															
-1	-1	1																																																																																																															
322	-295	-208																																																																																																															
-323	286	-270																																																																																																															
23	27	28	19	15																																																																																																													
140	145	148	153	167																																																																																																													
33	54	83	92	94																																																																																																													
94	99	107	110	98																																																																																																													
-1	-1	-1																																																																																																															
1	1	1																																																																																																															
-1	-1	-1																																																																																																															
185	143	137																																																																																																															
-563	-533	-514																																																																																																															
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																																																																																																															
885	116	165																																																																																																															
196	145	-149																																																																																																															
1392	-36	94																																																																																																															
-690	-102	-933																																																																																																															
	(c)	<div>The ability of convolutional networks to learn accurate models with many fewer weights than multi-layer perceptrons of similar depth is often attributed to shared weights and sparse connections. Explain the meaning of these two terms.</div>																																																																																																															
		[5 marks]																																																																																																															
		<div>Sample Answer</div> <div>The connections between layers in a convolutional neural network are much less dense than the connections within simpler feed-forward networks (e.g. multi-layer perceptrons). The image below illustrates this.</div>																																																																																																															

Channel 2

23	27	28	19	15
140	145	148	153	167
33	54	83	92	94
94	99	107	110	98

-1

-1

-1

1

1

1

-1

-1

-1

Channel 2

185	143	137
-563	?	-514

Channel 3

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

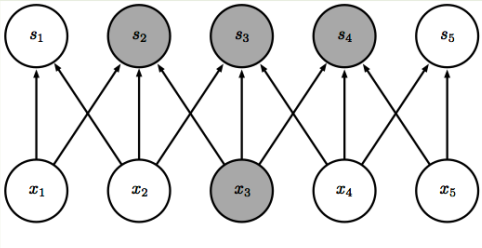
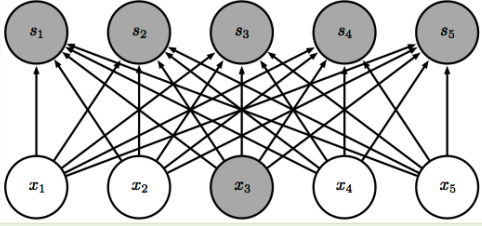
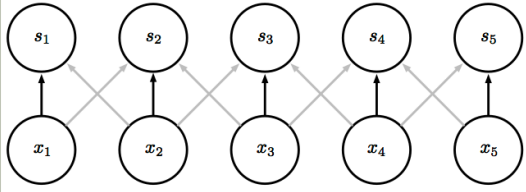
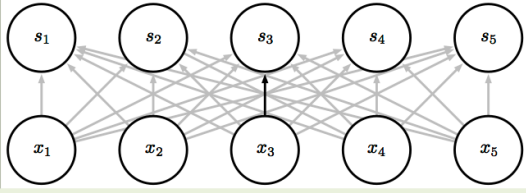
-1

Channel 3

?	116	?
196	145	-149

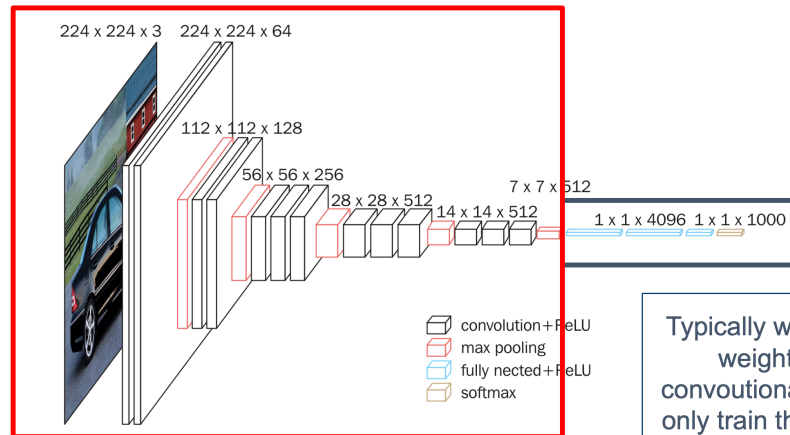
Final Output

1392	?	?
?	?	-933

		<p>This arises from the fact that small convolutional kernels are used and so only a small number of the inputs to a network are actually connected to any hidden layer unit in the network.</p> <div style="display: flex; align-items: center;"> <div style="flex: 1;"> <p>Sparse connections due to small convolution kernel</p> </div> <div style="flex: 2;">  </div> </div> <hr/> <div style="display: flex; align-items: center;"> <div style="flex: 1;"> <p>Dense connections</p> </div> <div style="flex: 2;">  </div> </div> <p>Convolutions are similarly responsible for the fact that weights in a CNN are shared. The bottom part of the image below shows the connections between two dense layers. In this scenario the weight on each link only affects one pair of computational units. In the CNN scenario above, however, the weights in the convolutional filter are applied at multiple positions within a network and so are shared across hidden units.</p> <div style="display: flex; align-items: center;"> <div style="flex: 1;"> <p>Convolution shares the same parameters across all spatial locations</p> </div> <div style="flex: 2;">  </div> </div> <hr/> <div style="display: flex; align-items: center;"> <div style="flex: 1;"> <p>Traditional matrix multiplication does not share any parameters</p> </div> <div style="flex: 2;">  </div> </div>
	(d)	<p>Transfer learning using pre-trained weights has become a standard approach to training large convolutional neural networks for image classification. Explain what this means, and how it is helpful for training convolutional neural networks with small datasets.</p>
		[6 marks]
		<p>Students should provide an explanation similar to the following. Large network models (e.g. VGG-16), pre-trained on broad coverage datasets (e.g. ImageNet), can be loaded into most machine learning packages now. By doing this we can take advantage of what has been learned before. Typically all but the layers after the flattening layer in a network are frozen so that weights don't change during training and only the final layers are</p>

trained. Something like the image below (albeit with much less detail) would be useful.

VGG-16

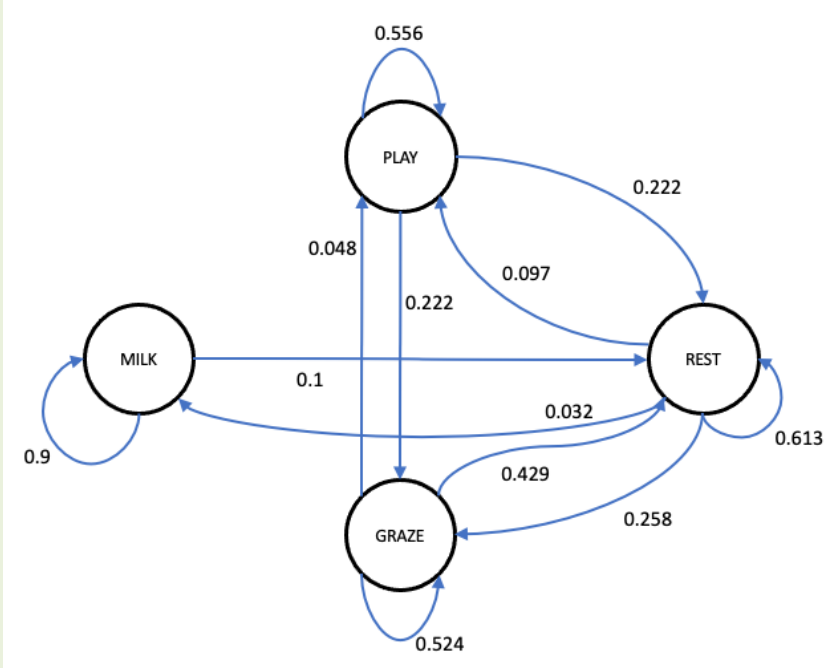
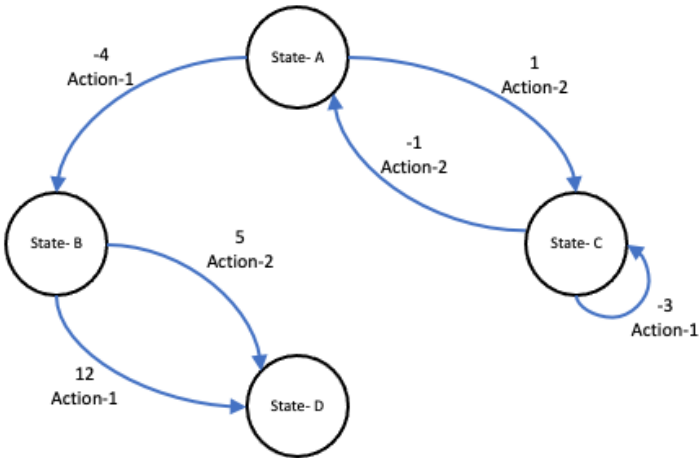


Typically we freeze the weights in the convolutional layers and only train the weights in the dense layers for our specific problem

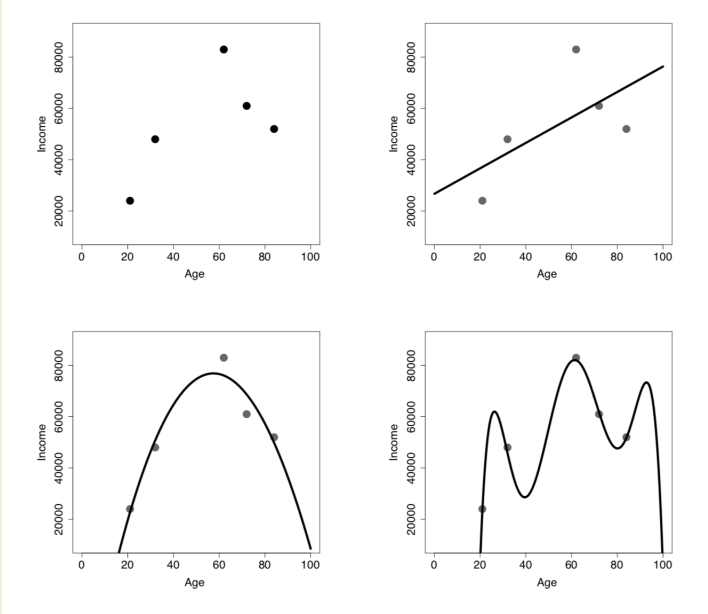
Students should then explain that by using pre-trained weights early layers of the network will be able to identify broadly useful features within images. These ought to be useful for any image classification problem. This is particularly useful when we have a small dataset for an image classification problem as we get the benefit of a large neural network, without needing the large dataset that is typically required to train such a network.

3.	(a)	Describe the concept of discounted return that is frequently used in reinforcement learning.
		[4 marks]
		<p><u>Sample Answer</u></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>
	(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> $\{r_0 = -33, r_1 = -11, r_2 = -12, r_3 = 27, r_4 = 87, r_5 = 156\}$ <p>Compare the discounted returns calculated at time $t = 0$ based on this reward sequence when discounting factors of 0.72 and 0.22 are used.</p>
		[6 marks]
		<p><u>Sample Answer</u></p> <p>Students should begin by calculating the discounted returns using:</p> $G_\gamma = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots + \gamma^{e-t} r_e$ <p>For discounting factor of 0.72:</p> $ \begin{aligned} G &= (0.72^0 \times -33) + (0.72^1 \times -11) + (0.72^2 \times -12) \\ &\quad + (0.72^3 \times 27) + (0.72^4 \times 87) + (0.72^5 \times 156) \\ &= (-33) + (0.72 \times -11) + (0.5184 \times -12) \\ &\quad + (0.3732 \times 27) + (0.2687 \times 87) + (0.1935 \times 156) \\ &= 16.4985 \end{aligned} $ <p>For discounting factor of 0.22:</p> $ \begin{aligned} G &= (0.22^0 \times -33) + (0.22^1 \times -11) + (0.22^2 \times -12) \\ &\quad + (0.22^3 \times 27) + (0.22^4 \times 87) + (0.22^5 \times 156) \\ &= (-33) + (0.22 \times -11) + (0.0484 \times -12) \\ &\quad + (0.0106 \times 27) + (0.0023 \times 87) + (0.0005 \times 156) \\ &= -35.4365 \end{aligned} $

		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.																																																																																																																																																
	(c)	<p>As part of a project to develop a farm simulator, the behaviour of a cow has been observed over a day. During the day the cow can REST, GRAZE grass, PLAY with other cows, or enter the milking parlour to MILK. The behaviour of the cow over the course of a day is captured in the table below (with time flowing down through the columns and left to right).</p> <table><tr><td>1</td><td>REST</td><td>13</td><td>PLAY</td><td>25</td><td>PLAY</td><td>37</td><td>MILK</td><td>49</td><td>REST</td><td>61</td><td>REST</td></tr><tr><td>2</td><td>GRAZE</td><td>14</td><td>REST</td><td>26</td><td>REST</td><td>38</td><td>MILK</td><td>50</td><td>GRAZE</td><td>62</td><td>REST</td></tr><tr><td>3</td><td>REST</td><td>15</td><td>REST</td><td>27</td><td>REST</td><td>39</td><td>MILK</td><td>51</td><td>GRAZE</td><td>63</td><td>PLAY</td></tr><tr><td>4</td><td>REST</td><td>16</td><td>GRAZE</td><td>28</td><td>REST</td><td>40</td><td>REST</td><td>52</td><td>GRAZE</td><td>64</td><td>PLAY</td></tr><tr><td>5</td><td>REST</td><td>17</td><td>GRAZE</td><td>29</td><td>REST</td><td>41</td><td>REST</td><td>53</td><td>REST</td><td>65</td><td>PLAY</td></tr><tr><td>6</td><td>GRAZE</td><td>18</td><td>REST</td><td>30</td><td>MILK</td><td>42</td><td>REST</td><td>54</td><td>REST</td><td>66</td><td>GRAZE</td></tr><tr><td>7</td><td>GRAZE</td><td>19</td><td>PLAY</td><td>31</td><td>MILK</td><td>43</td><td>REST</td><td>55</td><td>GRAZE</td><td>67</td><td>GRAZE</td></tr><tr><td>8</td><td>GRAZE</td><td>20</td><td>GRAZE</td><td>32</td><td>MILK</td><td>44</td><td>REST</td><td>56</td><td>GRAZE</td><td>68</td><td>REST</td></tr><tr><td>9</td><td>REST</td><td>21</td><td>REST</td><td>33</td><td>MILK</td><td>45</td><td>REST</td><td>57</td><td>GRAZE</td><td>69</td><td>REST</td></tr><tr><td>10</td><td>REST</td><td>22</td><td>GRAZE</td><td>34</td><td>MILK</td><td>46</td><td>GRAZE</td><td>58</td><td>REST</td><td>70</td><td>REST</td></tr><tr><td>11</td><td>PLAY</td><td>23</td><td>GRAZE</td><td>35</td><td>MILK</td><td>47</td><td>GRAZE</td><td>59</td><td>GRAZE</td><td>71</td><td>REST</td></tr><tr><td>12</td><td>PLAY</td><td>24</td><td>PLAY</td><td>36</td><td>MILK</td><td>48</td><td>GRAZE</td><td>60</td><td>REST</td><td>72</td><td>REST</td></tr></table>	1	REST	13	PLAY	25	PLAY	37	MILK	49	REST	61	REST	2	GRAZE	14	REST	26	REST	38	MILK	50	GRAZE	62	REST	3	REST	15	REST	27	REST	39	MILK	51	GRAZE	63	PLAY	4	REST	16	GRAZE	28	REST	40	REST	52	GRAZE	64	PLAY	5	REST	17	GRAZE	29	REST	41	REST	53	REST	65	PLAY	6	GRAZE	18	REST	30	MILK	42	REST	54	REST	66	GRAZE	7	GRAZE	19	PLAY	31	MILK	43	REST	55	GRAZE	67	GRAZE	8	GRAZE	20	GRAZE	32	MILK	44	REST	56	GRAZE	68	REST	9	REST	21	REST	33	MILK	45	REST	57	GRAZE	69	REST	10	REST	22	GRAZE	34	MILK	46	GRAZE	58	REST	70	REST	11	PLAY	23	GRAZE	35	MILK	47	GRAZE	59	GRAZE	71	REST	12	PLAY	24	PLAY	36	MILK	48	GRAZE	60	REST	72	REST
1	REST	13	PLAY	25	PLAY	37	MILK	49	REST	61	REST																																																																																																																																							
2	GRAZE	14	REST	26	REST	38	MILK	50	GRAZE	62	REST																																																																																																																																							
3	REST	15	REST	27	REST	39	MILK	51	GRAZE	63	PLAY																																																																																																																																							
4	REST	16	GRAZE	28	REST	40	REST	52	GRAZE	64	PLAY																																																																																																																																							
5	REST	17	GRAZE	29	REST	41	REST	53	REST	65	PLAY																																																																																																																																							
6	GRAZE	18	REST	30	MILK	42	REST	54	REST	66	GRAZE																																																																																																																																							
7	GRAZE	19	PLAY	31	MILK	43	REST	55	GRAZE	67	GRAZE																																																																																																																																							
8	GRAZE	20	GRAZE	32	MILK	44	REST	56	GRAZE	68	REST																																																																																																																																							
9	REST	21	REST	33	MILK	45	REST	57	GRAZE	69	REST																																																																																																																																							
10	REST	22	GRAZE	34	MILK	46	GRAZE	58	REST	70	REST																																																																																																																																							
11	PLAY	23	GRAZE	35	MILK	47	GRAZE	59	GRAZE	71	REST																																																																																																																																							
12	PLAY	24	PLAY	36	MILK	48	GRAZE	60	REST	72	REST																																																																																																																																							
	(i)	Based on this behaviour sequence, calculate a transition matrix that gives the probability of moving between the four states.																																																																																																																																																
		[8 marks]																																																																																																																																																
		<p><u>Sample Answer</u></p> <p>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:</p> <table><tr><td></td><td>PLAY</td><td>MILK</td><td>REST</td><td>GRAZE</td></tr><tr><td>PLAY</td><td>5.0</td><td>0.0</td><td>2.0</td><td>2.0</td></tr><tr><td>MILK</td><td>0.0</td><td>9.0</td><td>1.0</td><td>0.0</td></tr><tr><td>REST</td><td>3.0</td><td>1.0</td><td>19.0</td><td>8.0</td></tr><tr><td>GRAZE</td><td>1.0</td><td>0.0</td><td>9.0</td><td>11.0</td></tr></table> <p>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:</p> <table><tr><td></td><td>PLAY</td><td>MILK</td><td>REST</td><td>GRAZE</td></tr><tr><td>PLAY</td><td>0.556</td><td>0.000</td><td>0.222</td><td>0.222</td></tr><tr><td>MILK</td><td>0.000</td><td>0.900</td><td>0.100</td><td>0.000</td></tr><tr><td>REST</td><td>0.097</td><td>0.032</td><td>0.613</td><td>0.258</td></tr><tr><td>GRAZE</td><td>0.048</td><td>0.000</td><td>0.429</td><td>0.524</td></tr></table>		PLAY	MILK	REST	GRAZE	PLAY	5.0	0.0	2.0	2.0	MILK	0.0	9.0	1.0	0.0	REST	3.0	1.0	19.0	8.0	GRAZE	1.0	0.0	9.0	11.0		PLAY	MILK	REST	GRAZE	PLAY	0.556	0.000	0.222	0.222	MILK	0.000	0.900	0.100	0.000	REST	0.097	0.032	0.613	0.258	GRAZE	0.048	0.000	0.429	0.524																																																																																														
	PLAY	MILK	REST	GRAZE																																																																																																																																														
PLAY	5.0	0.0	2.0	2.0																																																																																																																																														
MILK	0.0	9.0	1.0	0.0																																																																																																																																														
REST	3.0	1.0	19.0	8.0																																																																																																																																														
GRAZE	1.0	0.0	9.0	11.0																																																																																																																																														
	PLAY	MILK	REST	GRAZE																																																																																																																																														
PLAY	0.556	0.000	0.222	0.222																																																																																																																																														
MILK	0.000	0.900	0.100	0.000																																																																																																																																														
REST	0.097	0.032	0.613	0.258																																																																																																																																														
GRAZE	0.048	0.000	0.429	0.524																																																																																																																																														

		(ii)	Draw a Markov process diagram to capture the behaviour of the cow as described by the transition matrix.
			[3 marks]
			<p>Sample Answer</p> <p>A diagram like this one is appropriate for this answer.</p> 
	(d)	<p>The following image shows a simple state transition diagram for a domain in which an agent can occupy one of four states and has two actions available to it.</p>  <p>Actions that connect states are shown above the arrows in the diagram. This environment is fully deterministic (actions always lead to the states shown in the diagram above). The rewards associated with each state</p>	

		<p>transition are shown in the diagram above the name of the action taken to complete the state transition.</p> <p>The action-value-function-table for this environment is shown in the table below.</p> <table><tr><th>State</th><th>Action</th><th>Value</th></tr><tr><td>State-A</td><td>Action-1</td><td>3.24</td></tr><tr><td>State-A</td><td>Action-2</td><td>-1.24</td></tr><tr><td>State-B</td><td>Action-1</td><td>6.78</td></tr><tr><td>State-B</td><td>Action-2</td><td>4.56</td></tr><tr><td>State-C</td><td>Action-1</td><td>-12.56</td></tr><tr><td>State-C</td><td>Action-2</td><td>-2.87</td></tr><tr><td>State-D</td><td>Action-1</td><td>0</td></tr><tr><td>State-D</td><td>Action-2</td><td>0</td></tr></table>	State	Action	Value	State-A	Action-1	3.24	State-A	Action-2	-1.24	State-B	Action-1	6.78	State-B	Action-2	4.56	State-C	Action-1	-12.56	State-C	Action-2	-2.87	State-D	Action-1	0	State-D	Action-2	0
State	Action	Value																											
State-A	Action-1	3.24																											
State-A	Action-2	-1.24																											
State-B	Action-1	6.78																											
State-B	Action-2	4.56																											
State-C	Action-1	-12.56																											
State-C	Action-2	-2.87																											
State-D	Action-1	0																											
State-D	Action-2	0																											
	(i)	If the agent begins in State-C which action will it select following a greedy action selection strategy ?																											
		[3 marks]																											
		Action-2 returns the highest reward so would be selected.																											
	(ii)	What state would the agent occupy after taking the action selected in Part (a) of this question and what reward would the agent receive after taking the action.																											
		[3 marks]																											
		State-A would be the next state visited.																											
	(iii)	Assuming that Q-learning is being used, update the entry in the action value table above for State-C and the action selected in Part (a) of this question. In your calculations assume that $\alpha = 0.1$ and that $\gamma = 0.9$.																											
		[3 marks]																											
		The Q value can be calculated as follows.																											
		<table><tr><td>Revised Q(State-C, Action-2)</td><td>=</td><td>-2.87 +</td><td>0.1 * (</td><td>-1 +</td><td>0.9 * 3.24 -</td><td>-2.87</td></tr><tr><td></td><td>=</td><td>-2.87 +</td><td>0.1 * (</td><td></td><td>4.786</td><td></td></tr><tr><td></td><td>=</td><td>-2.39</td><td></td><td></td><td></td><td></td></tr></table>	Revised Q(State-C, Action-2)	=	-2.87 +	0.1 * (-1 +	0.9 * 3.24 -	-2.87		=	-2.87 +	0.1 * (4.786			=	-2.39										
Revised Q(State-C, Action-2)	=	-2.87 +	0.1 * (-1 +	0.9 * 3.24 -	-2.87																							
	=	-2.87 +	0.1 * (4.786																								
	=	-2.39																											

4.	(a)	Machine learning algorithms face a constant struggle between over-fitting and under-fitting . Explain what this means.
		[10 marks]
		<p><u>Sample Answer</u></p> <p>1) What is meant by an ill-posed problem and what are the implications of this for machine learning.</p> <ul style="list-style-type: none"> Inductive machine learning algorithms essentially search through a hypothesis space to find the best hypothesis that is consistent with the training data used. It is possible to find multiple hypotheses that are consistent with a given training set (i.e. agrees with all training examples). It is for this reason that inductive machine learning is referred to as an ill-posed problem as there is typically not enough information in the training data used to build a model to choose a single best hypothesis. Inductive machine learning algorithms must somehow choose one of the available hypotheses as the best. An example like that shown in the figure below would be useful at this point  <p>The figure consists of four scatter plots arranged in a 2x2 grid, all showing 'Income' on the y-axis (ranging from 20,000 to 80,000) and 'Age' on the x-axis (ranging from 0 to 100). Each plot contains six data points. The top-left plot shows the raw data points without a fitted line. The top-right plot shows a linear fit line, which is a straight line passing through the points, representing underfitting. The bottom-left plot shows a quadratic fit curve, which is a smooth parabola that passes through all six data points, representing a good fit. The bottom-right plot shows a high-degree polynomial fit curve, which is a wiggly line that passes through all six data points but oscillates wildly between them, representing overfitting.</p> <p>2) How do machine learning algorithms deal with the fact that machine learning is ill posed.</p> <ul style="list-style-type: none"> Because inductive learning is ill-posed, we have to make some extra assumptions to have a unique solution with the data we have. The set of assumptions we make to have learning possible is called the inductive bias of the learning algorithm - this is the main implication of inductive machine learning being ill-posed. <p>3) Define what is meant by inductive bias:</p> <ul style="list-style-type: none"> The inductive bias of a learning algorithm:

		<ol style="list-style-type: none"> 1. is a set of assumption about what the true function we are trying to model looks like. 2. defines the set of hypotheses that a learning algorithm considers when it is learning. 3. guides the learning algorithm to prefer one hypothesis (i.e. the hypothesis that best fits with the assumptions) over the others. 4. is a necessary prerequisite for learning to happen because inductive learning is an ill posed problem. <ul style="list-style-type: none"> • An example of the specific inductive bias introduced by particular machine learning algorithms would be good here. E.g.: <ol style="list-style-type: none"> ○ Maximum margin: when drawing a boundary between two classes, attempt to maximize the width of the boundary. This is the bias used in Support Vector Machines. The assumption is that distinct classes tend to be separated by wide boundaries. ○ Minimum cross-validation error: when trying to choose among hypotheses, select the hypothesis with the lowest cross-validation error. <p>4) The importance and difficulty of selecting the right inductive bias:</p> <ul style="list-style-type: none"> • As learning is not possible without inductive bias the question becomes how to choose the right bias? This, however, is important and difficult because: <ol style="list-style-type: none"> ○ If the inductive bias of the learning algorithm constrains the search to only consider simple hypotheses, we may have excluded the real function from the hypothesis space. In other words, the true function is unrealizable in the chosen hypothesis space, (i.e., we are underfitting). ○ If the inductive bias of the learning algorithm allows the search to consider complex hypotheses, the model may hone in on irrelevant factors in the training set. In other words, the model with overfit the training data.
	(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. Decision s to be made based on this evaluation include which modelling approach will be used, which data pre-processing techniques</i></p>

		<p><i>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>
	(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]
		<p><u>Sample Answer</u></p> <p>This is an open ended question. To score well students should discuss the following points:</p> <ul style="list-style-type: none"> • The GDPR refers to automated processing which will include machine learning • The GDPR refers to "automated processing, including profiling, which produces legal effects concerning him or her or similarly significantly effects him or her". This does not cover all applications of machine learning but only those that have a legal or significant effect on a data subject. For example, image classification for improving search (a large application area for deep learning) is unlikely to fall under this definition. • Many refer to a "right to explanation" within the GDPR but there is some debate over how much that is present.

		<ul style="list-style-type: none">• If a right to explanation is required sophisticated machine learning models like deep neural networks will certainly be affected.
--	--	---