

2019-11-07 - M-ITECH - Basic Machine Learning - 201600070

Study: M-ITECH-201600070-1A M-ITECH Basic Machine Learning 201600070

Number of questions: 32
Generated on: Nov 1, 2019

Contents:		Pages:
▪ A. Voorpagina		1
▪ B. Vragen		19

2019-11-07 - M-ITECH - Basic Machine Learning - 201600070

Study: M-ITECH Basic Machine Learning 201600070

Exam Master Course Basic Machine Learning Course code: 201600070

Name: _____

Student number: _____

Introduction

Important message: When you have finished your exam do not leave your table to hand in your exam material but raise your hand and wait until one of the supervisors picks up your exam material. After that you are allowed to leave the room.

This exam is open book and consists of multiple-choice questions. You are allowed to use a simple calculator, but not your mobile phone, tablet or laptop or any other electronic means of computation or communication. Fill in your answers on the multiple choice answer form. Tables for the $\sigma(x)$ function and $-p \log_2(p)$ can be found at the end of the exam. In the questions we will use the standard notation of the course book.

Tips:

- Read each question carefully keeping the possible answers covered.
- Try to answer the question yourself, before you look at the answers you are given to choose from. Make a note of your first thoughts and calculations on a scribbling-paper.
- Beware of double negations (negatives) as these can be confusing.
- Do not stay on any one question too long. If you do not know the answer and have spent more than 10 minutes on the question, move on to the next question and come back to this one later.
- Fill in your answers on the answer form and hand it in with your name and student number on it. Also hand in the exam.
- If there is some time left at the end, check your answers before you hand in exam and the answer form. Did you write your name and student id on it?

Good luck!

Number of questions: 32

You can score a total of 32 points for this exam, you need 21.2 points to pass the exam.

1 Suppose that we are training a perceptron using the perceptron learning rule and that the current discriminant is given by the line $2 + 2x_1 - x_2 = 0$. The next feature point in our training set is given by $x = (2, 4)$. Assume that this feature point is misclassified, what will be the new value for the weight w after one update if one applies a learning rate of 0.4?

1 pt.

- a. (1.6, 1.2, -2.6)
- b. (1.2, 1.2, -0.6)
- c. (2.4, 2.8, 0.6)
- d. (2, 1.2, -2.6)

2 Once again consider the situation of the previous question. In addition to the perceptron learning rule with learning rate 0.4 one also applies L1 regularization of the form $|w_1| + |w_2|$ with parameter 0.1 (in this case this means that the product of the learning rate and regularization term is 0.1). What will now be the new value of w ?

1 pt.

- a. (2.4, 2.7, 0.5)
- b. (1.6, 1.3, -2.7)
- c. (1.6, 1.1, -2.5)
- d. (1.5, 1.1, -2.5)

3 Which if the following binary logical functions cannot be implemented by a single perceptron?

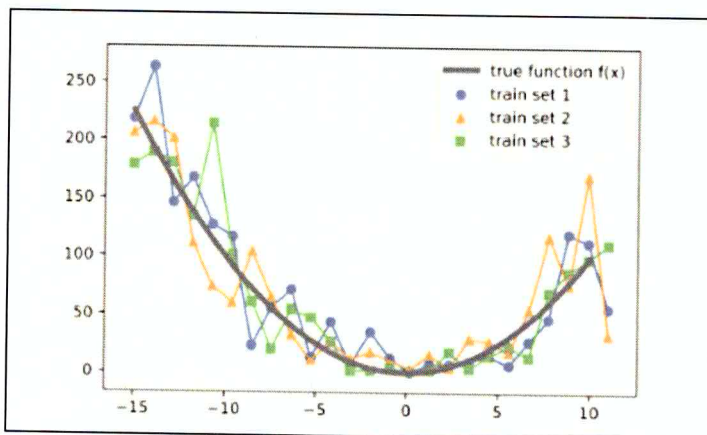
1 pt.

- a. AND
- b. OR
- c. NAND (negation of the AND function)
- d. XOR

4 1 pt. Suppose that we are training a logistic classifier using the standard *cross-entropy* error function and that the current weights of logistic classifier are given by $w = (1, 2, -1)$. The next feature point in our training set is given by $x = (2, 4)$ and the target label for x is 0. What will be the new value for the weight vector w after one stochastic update if one applies a learning rate of 0.7?

- a. $(0.49, 0.98, -3.05)$
- b. $(-0.27, 0.54, -3.92)$
- c. $(1.51, 3.02, 1.05)$
- d. $(0.90, 1.80, -1.40)$

5 1 pt. Consider the following plots of a regression model trained on different training sets drawn from the same overall dataset.



How would you characterize this model in terms of bias and variance?

- a. Model has high variance and bias.
- b. Model has low variance and bias.
- c. Model has high variance and low bias.
- d. Model has low variance and high bias.

6 Suppose that we train a one hidden layer Neural Network on a regression problem. During the model selection phase we increase the amount of neurons in the hidden layers step by step. What will happen with the variance and the bias of the model?
1 pt.

- a. Bias will increase and variance will increase.
- b. Bias will increase and variance will decrease.
- c. Bias will decrease and variance will decrease.
- d. Bias will decrease and variance will increase.

7 Consider the following statements:
1 pt.

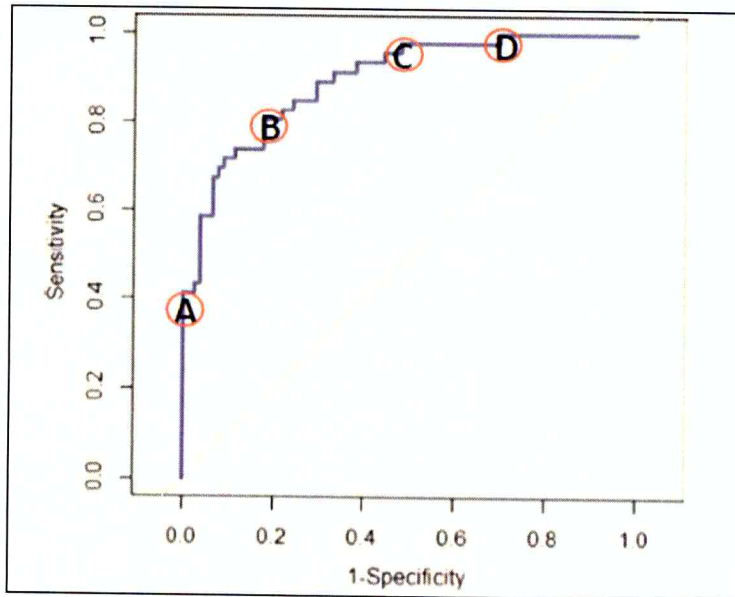
- (I) Overfitting can be prevented by pre-training.
- (II) Overfitting can be prevented by constructing additional features.
- (III) Overfitting can be prevented by weight sharing.

Which of the above statements are true?

- a. Only (I) and (III) are true.
- b. Only (II)
- c. Only (I)
- d. All three statements are true.

8 Consider the following ROC curve of a probabilistic classifier for a serious disease.

1 pt.



A doctor using this classifier does not want to make a wrong diagnose for any patient having the disease. The threshold corresponding to which point on the ROC curve should be used?

- a. A
- b. B
- c. C
- d. D

9 Consider a two class classification in which $t = 1$ denotes the positive class and $t = 0$ denotes the negative class. On a training set we trained a model M and on a small validation set the model M generates the following outputs:

1 pt.

index n	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
target t	1	0	0	1	1	0	0	0	1	1	0	0	0	1	1	0
output M	0.6	0.4	0.6	0.5	0.7	0.2	0.6	0.9	0.4	0.9	0.1	0.2	0.8	0.3	0.8	0.2

What is the TP-rate and FP-rate of M on this validation set for a threshold value $\theta = 0.75$?

- a. TP-rate=0.22 and FP-rate=0.29
- b. TP-rate=0.71 and FP-rate=0.78
- c. TP-rate= 0.42 and FP-rate= 0.33
- d. TP-rate=0.29 and FP-rate=0.22

10 Consider the following confusion matrix

1 pt.

		Predicted class		
		C_1	C_2	C_3
Actual Class	C_1	110	8	7
	C_2	16	130	10
	C_3	26	5	120

What is the accuracy of this classifier?

- a. $110/125+130/156+120/151$
- b. $110/152+130/143+120/137$
- c. $360/432$
- d. $110/125$

11 Consider the following confusion matrix

1 pt.

		Predicted class		
		C_1	C_2	C_3
Actual Class	C_1	110	8	7
	C_2	16	130	10
	C_3	26	5	120

Given the above confusion matrix. What is the recall for class C_1 ?

- a. $110/125$
- b. $110/152$
- c. $110/172$
- d. $110/432$

12 Assume we have probabilistic classifier $Cl(x)$ which classifies a datapoint x as class C_1 if $Cl(x) > \theta$

1 pt.

. What happens with the precision and recall for class C_1 if θ increases?

- a. Both precision and recall will increase.
- b. Both precision and recall will decrease.
- c. Precision will increase and recall will decrease.
- d. Precision will decrease and recall will increase.

13

1 pt.

Consider a two class problem with classes A and B and a data point x . Moreover assume that $p(x|A)=0.80$, $P(A)=0.60$ and $P(x|B)=0.30$. What is the probability that x belongs to class A , i.e. $P(A|x)$?

- a. 0.60
- b. 0.12
- c. 0.48
- d. 0.80

14

1 pt.

Consider the neural network (NN) for which the input is 3 dimensional and that there are 2 neurons in the hidden layer and that there are 2 output neurons. The activation for all neurons is the sigmoid function σ . The weights of the NN are as follows.

Hidden layer:

$$w_{1,0}^{(1)} = -4, w_{1,1}^{(1)} = 1, w_{1,2}^{(1)} = 1, w_{1,3}^{(1)} = 0$$

$$w_{2,0}^{(1)} = -3, w_{2,1}^{(1)} = 1, w_{2,2}^{(1)} = 2, w_{2,3}^{(1)} = 1$$

Output layer:

$$w_{1,0}^{(2)} = -3, w_{1,1}^{(2)} = 1, w_{1,2}^{(2)} = 2$$

$$w_{2,0}^{(2)} = 1, w_{2,1}^{(2)} = -1, w_{2,2}^{(2)} = 1$$

What will be the output for the NN on the input $x = (x_1, x_2, x_3) = (2, 2, -1)$? A table with values for $\sigma(x)$ can be found at the end of the exam. Select the alternative which is closest to your answer.

- a. (0.68, 0.20)
- b. (0.74, 1.38)
- c. (0.22, 0.54)
- d. (0.32, 0.80)

- 15** Consider the neural network (NN) for which the input is 3 dimensional and that there are 2 neurons in the hidden layer and that there are 2 output neurons. The activation for all neurons is the sigmoid function σ . The weights of the NN are as follows.

Hidden layer:

$$w_{1,0}^{(1)} = -4, w_{1,1}^{(1)} = 1, w_{1,2}^{(1)} = 1, w_{1,3}^{(1)} = 0$$

$$w_{2,0}^{(1)} = -3, w_{2,1}^{(1)} = 1, w_{2,2}^{(1)} = 2, w_{2,3}^{(1)} = 1$$

Output layer:

$$w_{1,0}^{(2)} = -3, w_{1,1}^{(2)} = 1, w_{1,2}^{(2)} = 2$$

$$w_{2,0}^{(2)} = 1, w_{2,1}^{(2)} = -1, w_{2,2}^{(2)} = 1$$

Consider the NN . Assume that for a given input the output is (0.6,0.4) and the target output is (0,1). Moreover assume that one applies stochastic gradient descent and the error function is given by:

$$\frac{1}{3} [|y_1 - t_1|^3 + |y_2 - t_2|^3]$$

What will be the vector $(\delta_1^{(2)}, \delta_2^{(2)})$ for the output neurons

- a. (-0.09, 0.09)
- b. (-0.36, 0.36)
- c. (0.36, -0.36)
- d. (0.09, -0.09)

- 16** Consider the neural network (NN) for which the input is 3 dimensional and that there are 2 neurons in the hidden layer and that there are 2 output neurons. The activation for all neurons is the sigmoid function σ . The weights of the NN are as follows.

1 pt.

Hidden layer:

$$w_{1,0}^{(1)} = -4, w_{1,1}^{(1)} = 1, w_{1,2}^{(1)} = 1, w_{1,3}^{(1)} = 0$$

$$w_{2,0}^{(1)} = -3, w_{2,1}^{(1)} = 1, w_{2,2}^{(1)} = 2, w_{2,3}^{(1)} = 1$$

Output layer:

$$w_{1,0}^{(2)} = -3, w_{1,1}^{(2)} = 1, w_{1,2}^{(2)} = 2$$

$$w_{2,0}^{(2)} = 1, w_{2,1}^{(2)} = -1, w_{2,2}^{(2)} = 1$$

Consider the above NN structure and weights. Assume that the δ vector $(\delta_1^{(2)}, \delta_2^{(2)})$ of the output layer is $(-0.4, -0.6)$ and that the output of the hidden neuron 1 is 0.5 and the output of the hidden neuron 2 is 0.3. What will be the delta $(\delta_1^{(1)})$ of the hidden neuron 1?

- a. 0.05
- b. -0.05
- c. 0.20
- d. -0.29

- 17 Consider the neural network (NN) for which the input is 3 dimensional and that there are 2 neurons in the hidden layer and that there are 2 output neurons. The activation for all neurons is the sigmoid function σ . The weights of the NN are as follows.

1 pt.

Hidden layer:

$$w_{1,0}^{(1)} = -4, w_{1,1}^{(1)} = 1, w_{1,2}^{(1)} = 1, w_{1,3}^{(1)} = 0$$

$$w_{2,0}^{(1)} = -3, w_{2,1}^{(1)} = 1, w_{2,2}^{(1)} = 2, w_{2,3}^{(1)} = 1$$

Output layer:

$$w_{1,0}^{(2)} = -3, w_{1,1}^{(2)} = 1, w_{1,2}^{(2)} = 2$$

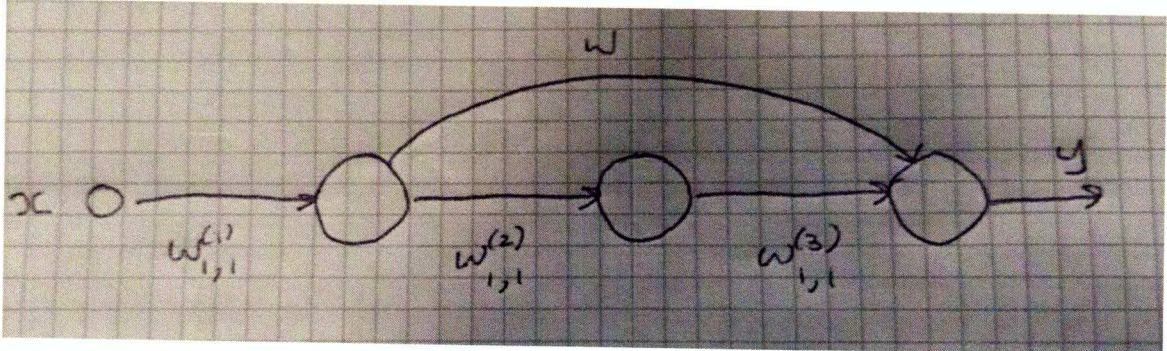
$$w_{2,0}^{(2)} = 1, w_{2,1}^{(2)} = -1, w_{2,2}^{(2)} = 1$$

Consider the same NN structure and weights as in the question above. Moreover the delta vector $\delta^{(2)} = (\delta_1^{(2)}, \delta_2^{(2)})$ of the output layer is $(-0.4, -0.6)$ and that the output of the hidden neuron 1 is 0.5 and the output of the hidden neuron 2 is 0.3. But now we assume that the NN shares the following weights: $w_{1,1}^{(1)} = w_{2,1}^{(1)}$, meaning these two variables are identical. What will be the adaptation dw to the weight $w_{1,1}^{(1)}$ if we apply a learning rate of 1? Use for the input x the values $x = (2, 2, -1)$

- a. 0.12
- b. -0.24
- c. 0.24
- d. -0.12

18
1 pt.

Assume we have a very simple neural network consisting of two hidden layers, each consisting of one neuron, one input neuron and one output neuron. The activation of each neuron is denoted by f . See figure below



This neural network as a so-called skip connection from the neuron in the first hidden layer to the output neuron. i.e the output of this neuron also goes directly to the output neuron. This connection as a weight denoted by w . Rest of the notation is standard and as usual $\delta_j^i = \frac{\partial E}{\partial a_j^{(i)}}$ with E the error or loss function.

What is the backpropagation formula for $\delta_1^{(1)}$ in terms of $\delta_1^{(3)}$, $\delta_1^{(2)}$, the activation function f , the derivative f' and the weights of the network?

- $\delta_1^{(1)} = [w\delta_1^{(3)} + w_{1,1}^{(2)}\delta_1^{(2)}]f'(a_1^{(1)})$
- $\delta_1^{(1)} = [w\delta_1^{(3)} + w_{1,1}^{(2)}\delta_1^{(2)} + w_{1,1}^{(3)}\delta_1^{(3)}]f'(a_1^{(1)})$
- $\delta_1^{(1)} = [w\delta_1^{(3)}f'(a_1^{(3)}) + w_{1,1}^{(2)}\delta_1^{(2)}]f'(a_1^{(1)})$
- None of the above.

19
1 pt.

Assume that:

- We have a four-class classification problem in a 2-dimensional space.
- We apply Bayes law to estimate $P(C_k|x)$, $k = 1, 2, 3, 4$
- We assume that the likelihoods are modelled by normal (Gaussian) probability distributions.

How many parameters does one need to estimate?

- 24
- 28
- 27
- 23

- 20** Consider a 1-dimensional dataset $A = \{5.0, 6.0, 7.0, 8.0, 9.0, 9.5, 9.6, 9.8\}$ and one is using for the conditional class likelihood $p(x|A)$ a Kernel density estimator of the form

1 pt.

$$p(x|A) = \frac{1}{N} \sum_{n=1}^N k(x - a_n) \text{ with } k(y) = 1 \text{ if } |y| \leq 1/2 \text{ and } k(y) = 0 \text{ otherwise}$$

In the above expression N is the number of elements in A and a_n is the n -th element in A . Given the above what is the value for $p(9.4|A)$?

- a. 3/8
 - b. 2/8
 - c. 4/8
 - d. 1/8
- 21** Consider a two class classification problem, classes C_1 and C_2 , for which we apply a probabilistic approach. The loss matrix for this classification problem is given by:

1 pt.

$$\begin{vmatrix} 0 & 2 \\ 3 & 0 \end{vmatrix}$$

Rows correspond to true classes, C_1 first row and C_2 second row, columns correspond to classification by the classifier. First (second) column corresponds to C_1 (C_2).

Assume that we apply a classification rule of the form: if $P(C_1|x) > \theta$ then $P(C_1|x) > \theta$ is classified as C_1 . What is the optimal value for θ given the loss matrix above?

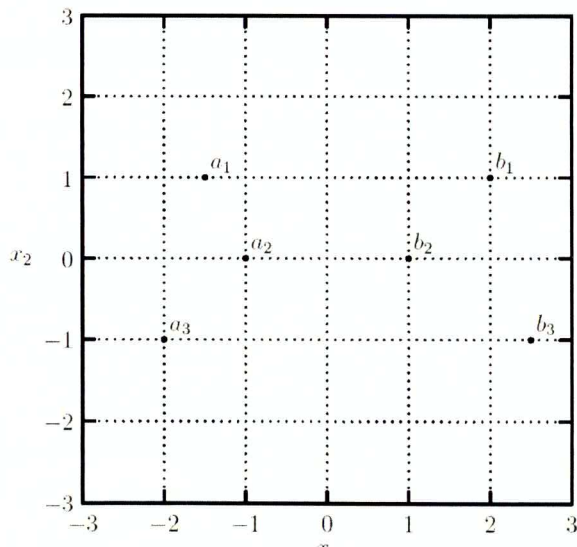
- a. 0.4
- b. 0.3
- c. 0.5
- d. 0.6

- 22** Naive Bayes makes the simplifying assumption that all of the input dimensions are conditionally independent given the class. Which of the following statements is correct:
1 pt.
- a. The naive Bayes model is identical to a model with Gaussian class-conditional distributions with diagonal covariance matrices.
 - b. The independence assumption makes the model faster to train, but does not reduce its performance
 - c. When the independence assumption is incorrect, Naive Bayes will incorrectly classify the datapoint.
 - d. When the independence assumption is incorrect, Naive Bayes will tend to overestimate the probability of the most likely class

- 23** Which statement is true for k-Nearest-Neighbours?
1 pt.
- a. There is no good way of dealing with multi-class classifications
 - b. Increasing the value of k reduces the risk of overfitting, but may lead to oversmoothing
 - c. As long as k is odd it is impossible to have region in the space where the only strategy is to take a random guess
 - d. Having a low value of k is only useful for speed reasons

24 Consider the dataset depicted below:

1 pt.



We want to classify this dataset with a support vector machine without kernel function. What are the support vectors in this case?

- a. a_1, a_2, b_2
- b. a_2, b_2, b_3 and a_2, a_3, b_1 give rise to the same margin and are both correct solutions
- c. All points in this simple case are support vectors
- d. a_2, b_2

25 Consider the following two statements about non-linear transformations of input features:

1 pt.

- (i) Classes which are separable remain linear separable after a non-linear transformation of the input features.
- (ii) Classes which are not linear separable can become linear separable after a nonlinear transformation of the input features.

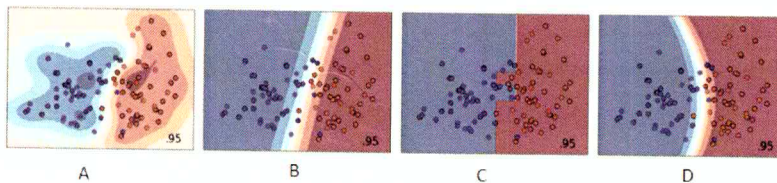
Which of the above statements are true?

- a. Both are true.
- b. Only (ii).
- c. Both are false.
- d. Only (i).

26 Which of the following statements is true? The "dual representation" of a Support Vector Machine (SVM)
1 pt.

- a. Results in a model that uses the training data twice, and therefore requires less training data to obtain the same performance
- b. Allows us to find a sparse solution for the problem
- c. Expresses the weights of the SVM in terms of the training datapoints
- d. Allows us to use a Lagrangian to indicate which training points are "support vectors", and is therefore much faster to compute than the "primal representation"

27 Which of the following plots shows a classification based on a decision tree?
1 pt.



- a. Plot A
- b. Plot C
- c. Plot B
- d. Plot D

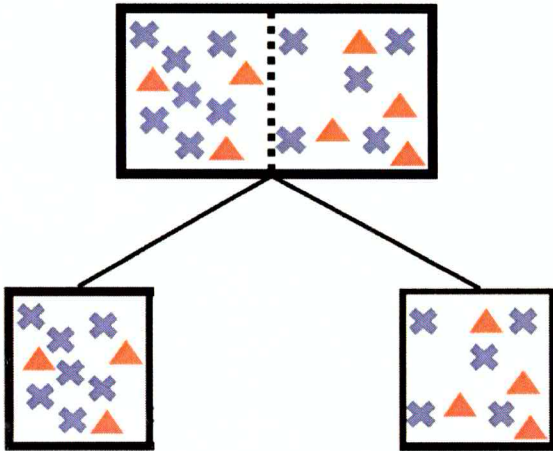
28 Consider a dataset with initial entropy of 0.95. For attribute (feature) B we get the following partition of the dataset with respect to the target labels Y and O:
1 pt.

B	Y	O
b1	40	10
b2	2	18

What is the information gain for attribute B? Select the alternative which is closest to your answer.

- a. 0.65
- b. 0.35
- c. 0.30
- d. 0.60

- 29** Having a dataset consisting of data points of either class A (blue crosses) or class B (red triangles), and split by a decision tree according to the figure below, 1 pt.



Which of the below statements about Kappa, denoted by K , and the usefulness of the split for a threshold of 0.05 are correct?

- a. $K = 0.0224$, the split is useful.
 - b. $K = 0.4247$, the split is useless.
 - c. $K = 0.0224$, the split is useless.
 - d. $K = 0.4247$, the split is useful.
- 30** Consider a K -class classification problem with $K > 2$. For this classification we can use several ML models: 1 pt.

- (1) A Decision Tree.
- (2) A Neural Network with K outputs.
- (3) $K(K - 1)/2$ binary classifiers one-versus-one with majority vote.

Which of these classifiers (options) suffer from ambiguous regions?

- a. Only option (1).
- b. Only option (2) and (3).
- c. Only option (3)
- d. Only option (1) and (3).

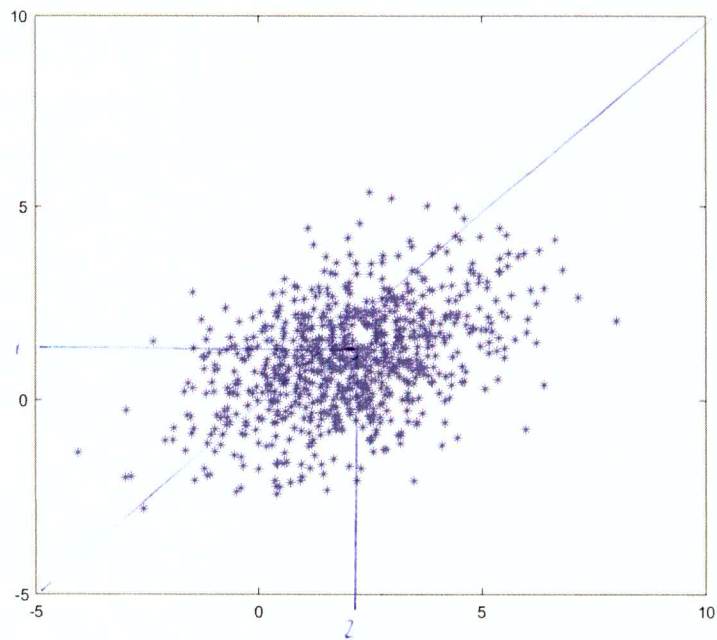
31 Which of the following statements is true?

1 pt.

- a. Auto-encoders are a dimensionality reduction method that explicitly filter out noise.
- b. Auto-encoders are a neural networks where the predicted output is a reconstruction of the (potentially corrupted) input. These allow us to create useful non-linear transformations of the data.
- c. Deep auto-encoders require heavy regularisation to work, and the only regularisation that is useful in practice is the L_1 regularisation.
- d. Auto-encoders just perform PCA, but they are still useful because they are much faster than the matrix decomposition that PCA requires.

32 Considering the following dataset.

1 pt.



Which covariance matrix describes the covariance in this dataset best?

a. $\begin{vmatrix} 3 & 0 \\ 0 & 2 \end{vmatrix}$

b. $\begin{vmatrix} 3 & -1 \\ -1 & 2 \end{vmatrix}$

c. $\begin{vmatrix} 0 & 2 \\ 3 & 0 \end{vmatrix}$

d. $\begin{vmatrix} 3 & 1 \\ 1 & 2 \end{vmatrix}$

Table for $-p \log_2(p)$

p	$-p \log_2(p)$	p	$-p \log_2(p)$	p	$-p \log_2(p)$
0	0	1/8	0.38	1/10	0.33
1	0	2/8	0.50	2/10	0.46
1/2	0.50	3/8	0.53	3/10	0.52
1/3	0.53	4/8	0.50	4/10	0.53
2/3	0.39	5/8	0.42	5/10	0.50
1/4	0.50	6/8	0.31	6/10	0.44
2/4	0.50	7/8	0.17	7/10	0.36
3/4	0.31	1/9	0.35	8/10	0.26
1/5	0.46	2/9	0.48	9/10	0.14
2/5	0.53	3/9	0.53	1/11	0.31
3/5	0.44	4/9	0.52	2/11	0.45
4/5	0.26	5/9	0.47	3/11	0.51
1/6	0.43	6/9	0.39	4/11	0.53
2/6	0.53	7/9	0.28	5/11	0.52
3/6	0.50	8/9	0.15	6/11	0.48
4/6	0.39			7/11	0.42
5/6	0.22			8/11	0.33
1/7	0.40			9/11	0.24
2/7	0.51			10/11	0.13
3/7	0.52				
4/7	0.46				
5/7	0.35				
6/7	0.19				

Table for $\sigma(x)$

x	$\sigma(x)$
-4	0.02
-3	0.05
-2	0.12
-1	0.27
0	0.50
1	0.73
2	0.88
3	0.95
4	0.98

