

Student ID Number:*(fill in your NUS Student ID here)*

Institute of Systems Science
National University of Singapore

**MASTER OF TECHNOLOGY IN
INTELLIGENT SYSTEMS**

Graduate Certificate Examination Semester I 2019/20

Subject: Pattern Recognition Systems

Sample Examination Questions

SECTION A

Question	Marks
1	/17
2	/23
TOTAL	/40

SECTION A

Question 1

(Total: 17 Marks)

17

Electrocardiogram (ECG) is widely used by cardiologists to monitor the functionality of the cardiovascular system. Figures 1 and 2 below show an example of ECG signal and a prototype cycle of ECG signal.

The main problem with manual analysis of ECG signals lies in difficulty of detecting and categorizing different waveforms and morphologies. To address the problems raised with the manual analysis, machine learning techniques can be used to accurately detect the possible anomalies in the signals.

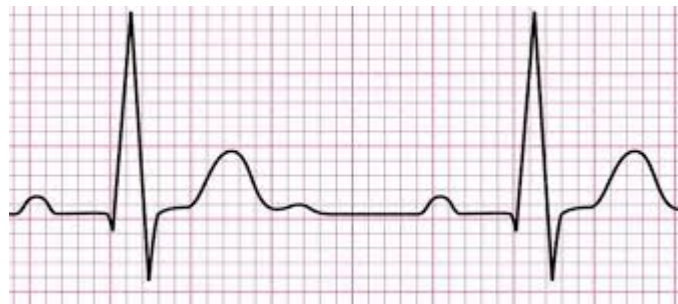


Figure 1. An Example of ECG Signal

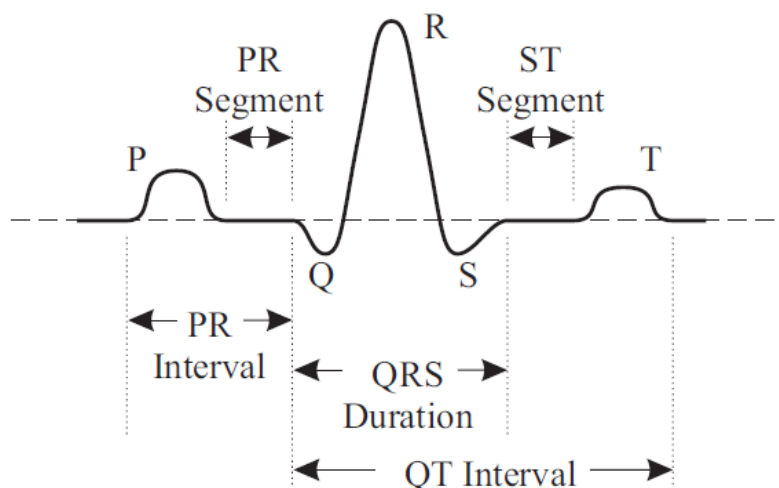


Figure 2. A Prototype Cycle of an ECG Signal

Kent Hospital has accumulated a total of 50,000 ECG records over the last three years. Each record has 200 input variables and one target value. The target value is a diagnosis provided by physicians, 1 for a healthy diagnosis and -1 for a Myocardial Infarction (MI) diagnosis meaning heart attack. 40,000 records represent healthy patients and the remaining 10,000 records have an MI diagnosis.

The 200 input variables and their brief descriptions are listed below:

- V1. patient ID
- V2. patient name
- V3. age in years
- V4. gender, -1=female, 1=male
- V5. Ethnicity, 1 = Chinese, 2 = Malays, 3 = Indians, 4 = Others
- V6. maximum heart rate in beats/min
- V7. minimum heart rate in beats/min
- V8. average time between heart beats in sec
- V9. rms deviation of the mean heart rate in beats/sec
- V10. full width at half maximum for the heart rate distribution
- V11. average QT interval for lead with max T wave
- V12. average QT interval for all leads
- V13. average corrected QT interval for lead with max T wave
- V14. average corrected QT interval for all leads
- V15. average QRS interval for all leads
- V16. average PR interval for lead with maximum P wave
- V17—V200. other variables related to the ECG signal

You are tasked to build a pattern recognition system for automatic diagnosis of ECG signals by using various machine learning techniques.

Answer the following questions:

- a. Assume you are building decision trees for the ECG diagnosis. What variables would you use as input and output variables for your decision tree models? Would you suggest normalizing the data before the modeling? Justify your answer.

(3 Marks)

- b. Physicians have provided you with further diagnosis information. Each MI ECG now has been assigned with one of the diagnoses as listed below:

- Rhythm disturbances
- Heart block and conduction problems
- Electrolytes disturbances and intoxication
- Ischemia and infarction

You are tasked to improve the system by using the newly available data. Propose a possible solution using Support Vector Machines (SVM) to solve this extended problem.

(3 Marks)

- c. Suggest two other machine learning techniques which are also suitable for building the system for part (b) besides the decision tree and SVM used above. Briefly justify your suggestion.

(2 Marks)

- d. There are too many variables in the data set which may affect the model performance. You aim to reduce the number of features and achieve better system performance. Propose a hybrid intelligent system approach which uses two machine learning techniques together to realize these. Explain briefly how it works.

(3 Marks)

- e. A principal component analysis was conducted on the augmented dataset described in part (b) above. The following shows the python code and output after running a principal component analysis:

```
In [32]: len(eigval)
Out[32]: 188

In [33]: sum(eigval)
Out[33]: 8.892215806505435

In [34]: eigval[0:10]
Out[34]:
array([3.09539125, 2.53826793, 0.83451697, 0.48165584, 0.27284785,
       0.24349596, 0.21309404, 0.15343467, 0.14830454, 0.1090943 ])
```

PC	Eigenvalue	Variance	Cumulative
1	3.09539125	0.348101229	0.348101229
2	2.53826793	0.285448305	0.633549534
3	0.83451697	0.093848034	0.727397568
4	0.48165584	0.054166009	0.781563577
5	0.27284785	0.030683899	0.812247475
6	0.24349596	0.027383047	0.839630522
7	0.21309404	0.02396411	0.863594632
8	0.15343467	0.017254942	0.880849575
9	0.14830454	0.016678019	0.897527593
10	0.1090943	0.012268517	0.90979611

- (1) Provide a hand sketch of the Scree Plot for the first 7 principal components.

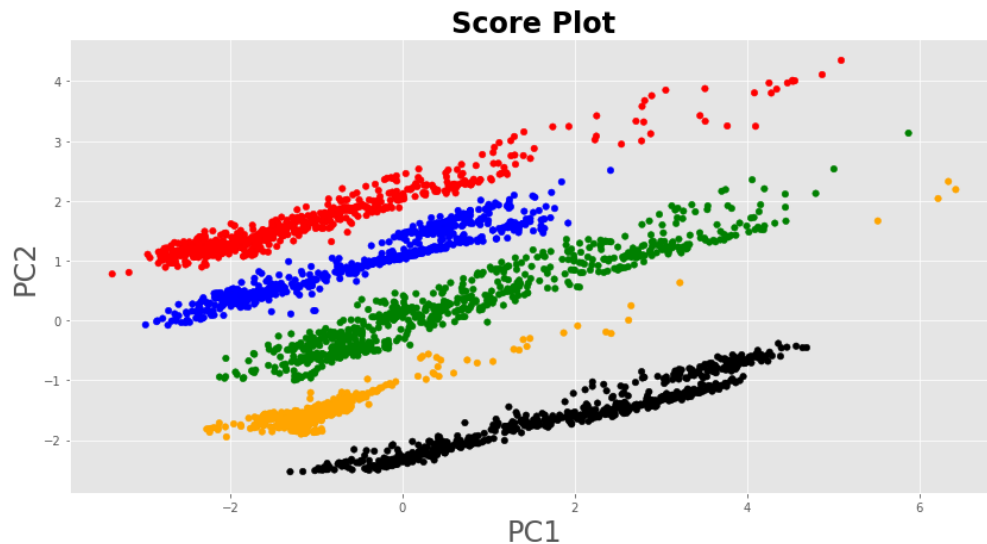
(2 Marks)

- (2) Explain how many principal components you would choose to extract.

(2 marks)

(3) The plot below shows the Score Plot for PC1 (x) and PC2 (y) with the following colored labels.

- Red: Rhythm disturbances
- Blue: Heart block and conduction problems
- Green: Electrolytes disturbances and intoxication
- Orange: Ischemia and infarction
- Black: Normal



Explain the use of a Score Plot. What can you observe about the data structure of the ECG dataset?

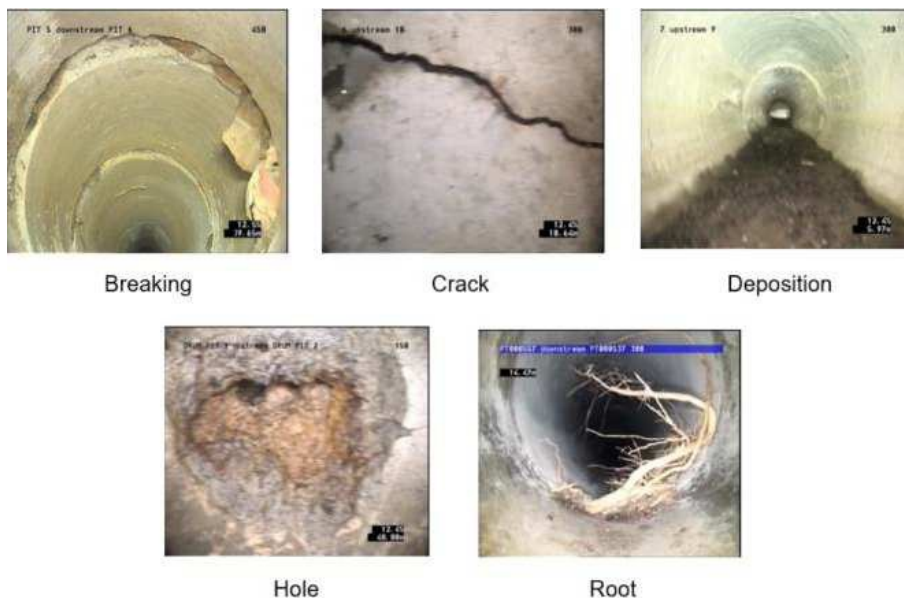
(2 Marks)

Question 2*(Total: 23 Marks)*

In heavy rain, impervious surfaces such as roads and driveways can contribute large amount of polluted stormwater. To prevent flooding, streets are usually lined with stormwater drains to quickly move stormwater off the road. These drains are connected to stormwater pipes that channel stormwater directly to nearby water body.

These stormwater pipes are often monitored in order to understand the status of the stormwater system, so that local authority can plan for maintenance and future facility planning. Currently these pipes are inspected on-site: a certified technician guides a robot with camera mounted and visually look for defects. Once a defect is found, the technician rotates and zooms in the camera manually to get a better understanding of the defect. This process is time consuming and expensive.

The local authority decided this process should be automated. They came to the start-up you are working in and sought for solution. They are looking at a system to identify five types of common defects in a pipe (See the below figure) automatically.



Currently they have captured a few thousands of images (grabbed from videos) and the size of each image given to you will be of 128 x 128 pixels. You are tasked to build a deep learning model for this problem.

- a. To start with, you are asked to come out a simple convolutional neural network to serve as a baseline model for further studies. Propose a model that uses only convolutional layers, max-pooling layers, flatten layer and dense layers. To restrict complexity, only 3 convolutional layers (without padding) are allowed. Furthermore, the total number of parameters in the model must be below 100,000, as the model is likely to be run at the edge.

Make a table to clearly illustrate the architecture. Explain the design of your architecture. For each layer, you need to specify the kernel size, the output shape, and the number of parameters. You are also required to inform the choice of activation function and initialization method for the layers. Explain the choices. Specify the loss function to be used for the training of the model.

(8 marks)

- b. The images given to you are captured from two different brands of cameras. You are told in future videos will be captured by at least 5 different brands of cameras, and thus your model should be robust against the variation of images in terms of lighting, colour and quality. Do you think this is an issue for your baseline model? What is the single measure that can make your model more robust against the above variation? How would you implement the measure for your baseline model?

(2 marks)

- c. Your colleague suggests putting in a set of residual layers. Remake the table of your baseline model to illustrate how the residual layers will be implemented (No need to indicate the number of parameters in each layer). For the added layers, you are required specify the input(s) to each added layer. Do you think it is a good idea to have more than 1 set of residual layers in the baseline model? Explain why.

(3 marks)

- d. To further enhance the automated robot vision-based inspection solution, your company has decided to deploy a set of IoT sensors to provide continuous pipe monitoring, such as water flow sensor. Table 1 provides an example of water flow sensor dataset. These sensor data are updated at an interval of one hour. You are asked to predict the risk of choked drainage using the water flow sensor data. The change of the water flow sensor reading over the time might give you some insight. For example, if the pipe is choked, the water flow speed might suddenly decrease to smaller values. Evaluate whether the Haar wavelet transformation can indicate the pipe choke (identify major change of the water flow using sensor reading over the time) from the water level data provided in Table 1. You need to show your calculation to justify your answers.

(3 marks)

Table 1. An example of water flow sensor dataset used for pip monitoring. These sensor data are updated at an interval of one hour.

Signal index	1	2	3	4	5	6	7	8
Water flow sensor reading (meter/second)	0.6	0.7	0.8	0.4	0.3	0.3	0.6	0.6

- e. To integrate both water level sensor data and robot vision data, you are asked to develop a multimodal feature-based sense making solution to make decision on whether you need to take action to do cleaning or not. In your proposed solution, you need to (i) choose appropriate deep learning techniques to extract features from water flow sensor data and robot vision data, respectively; (ii) describe how to fuse these two multimodal data to make decision, which is a two-category classification problem (Action/No action).

(3 marks)

- f. You are asked to build a decision-based fusion system. First, the water flow sensor data is used to make a decision on Low or High risk of choke, second, the robot vision data is used to make decision on Normal/Root defect/Crack defect. Based on these two decisions from individual type of data, you are asked to build a probability-based classification system to make decision (Action/No action). Given the new test data [High risk of choke, Crack defect], apply your classification method to make decision (i.e., Action/No action). Show your calculation steps to justify your answer.

Table 2. An example of training dataset collected for decision-based fusion system.

Training data index	Decision based on water flow sensor data (Low / High risk of choke)	Decision based on robot vision inspection (Normal / Root defect / Crack defect)	Type of action needed (Action / No action)
1	Low risk of choke	Normal	No action
2	Low risk of choke	Root defect	No action
3	Low risk of choke	Crack defect	No action
4	High risk of choke	Normal	Action
5	High risk of choke	Root defect	Action
6	Low risk of choke	Root defect	Action
7	High risk of choke	Crack defect	Action
8	High risk of choke	Normal	No action

(4 marks)