

3.6 min / mark

## Section A: Short Answer Questions [13 marks]

Answer each of the questions in this section as briefly as possible. Expect to answer each sub-question in a couple of lines per mark.

### Question 1: General Machine Learning [13 marks]

- (a) Name two problems in *frequentist supervised learning* that can be addressed using *regularisation*, and in each case, explain why. [2 marks]

① solve the ill-posed problem, regulation can vary the complexity and re-condition the problem. For example, ridge regression  
② solve the problem of overfitting, it can penalize weights have high value and control weight in a reasonable range

- (b) Consider the *frequentist supervised learning* setup. Is it always possible to achieve zero training error by increasing *model complexity*? Explain using an example(s), and discuss whether zero training error is a good thing. [2 marks]

Yes, for example, an ANN model with a hidden layer and appropriate activation function can approximate any continuous function very good. No a good thing, overfitting and

- (c) Name an advantage of using *lasso regularisation* method over the *ridge regression*. [1 mark] lack of generalization.

Solutions are sparse and suitable for high-dimension data.

- (d) *Perceptrons* and *hard-margin SVMs* are both *linear classifiers*. What is the difference between the two methods? [1 mark]

For perceptrons, the only goal is to separate the data, but SVM would try to find the best hyperplane by maximizing

- (e) State what mathematical operation is implemented by the *back-propagation* algorithm, as used in *learning artificial neural networks*. [1 mark]

① need a loss function

② set the partial derivatives of loss function to zero and get maximal change of loss function (the derivative)

- (f) In a sentence each, describe what *probabilistic inference* and *statistical inference* are for. [2 marks]

① probabilistic inference compute marginal and conditional distributions from the joint of a PGM using Bayes rules and marginalization (Bayesian)

values)

② statistic inference fits probability table with observation (Frequentist)

③ use those derivative values to update our weight.

- (g) Explain why *Bayesian model selection* tends to prefer simpler models over complex ones, and state a situation when a more complex model would be preferred. [2 marks]

Bayesian penalized those model with higher complexity more, the marginalization process weight.

in it can provide regularization. If we have a large amount of data with high diversity, it would prefer complex model.

- (h) Describe the process of *Bayesian sequential updating*, making reference to the *prior*, *likelihood* and *posterior* distributions. [1 mark]

1. Choose prior initially

2. calculate the posterior with the observe likelihood

3. use it as a prior of next step, repeat, until the model

- (i) *Gaussian mixture modeling* and *k-means algorithm* can both be used for *clustering*. Name another similarity between these methods. [1 mark]

They all need a specific k before training.

## Section B: Method Questions [17 marks]

In this section you are asked to demonstrate your conceptual understanding of a subset of the methods that we have studied in this subject.

## Question 2: Linear Regression [6 marks]

- (a) Outline the steps for setting the value of regularisation parameter  $\lambda$  in ridge regression using cross-validation or held-out validation. [2 marks] Held-out: 1. split the data into training set and test set  
2. use grid search, for different  $\lambda$ , use the training set to fit the model, then use the test set to test the model  
3. choose  $\lambda$  with best results in the evaluation processes.
- (b) Why can't one find the value of  $\lambda$  in the same way weights  $w$  are found? [1 mark]

Because  $\lambda$  and the data is independent.

- (c) Show that the standard method for training a linear regression model (minimising the sum of squared errors) is equivalent to the maximum likelihood estimate for the graphical model,  $y \sim \mathbf{x}'\mathbf{w} + \epsilon$ , where  $y$  are the outputs,  $\mathbf{x}$  the inputs,  $\mathbf{w}$  the model parameters and  $\epsilon$  is zero-mean Gaussian noise. Present your answer mathematically, and show the stages of your working. [3 marks]

In your answer to the above, may find the formulation of the Gaussian distribution handy:

$$N(v|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(v-\mu)^2\right\}. \quad \begin{aligned} y &\sim \mathbf{x}'\mathbf{w} + \epsilon \\ &\sim N(\mathbf{x}'\mathbf{w} + \epsilon | 0, \sigma^2) \\ p &= P(y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_n) = \prod_i P(y_i | x_i) \end{aligned}$$

$$\log(p) = \sum_i \log(P(y_i | x_i)) = \sum_i -\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(y_i - \mathbf{x}'\mathbf{w})^2$$

## Question 3: Support Vector Machines [4 marks]

- (a) A method is called non-parametric if it does not rely on a fixed number of parameters. Explain why a SVM can be considered non-parametric. [2 marks]

The model of SVM depends on support vector, it make prediction via sign of  $b + \sum_{i=1}^n A_i y_i \mathbf{x}' \mathbf{x}$

- (b) Consider a definition of a kernel as the dot product in some feature space. Show that a sum of two kernels is a kernel. [2 marks] proof, suppose:  $K_1(u, v), K_2(u, v)$  is valid kernel

$$\therefore \forall A \in \mathbb{R} \quad \therefore A^T K_1 A \geq 0, A^T K_2 A \geq 0$$

not on a fixed number of parameters,

$$A^T K_1 A + A^T K_2 A$$

$$= A^T (K_1 + K_2) A$$

$$\geq 0$$

## Question 4: Ensemble Methods [3 marks]

Outline how the bagging approach to ensemble learning works. What is the key assumption underlying the method? Explain why this leads to bagging often making better predictions than just using a single classifier.

1. use bagging to generate different datasets for different model

train those model with the dataset they get from bagging,

make prediction based on the majority (or average) result of the model

对了?

## Question 5: Probabilistic Models [4 marks]

- (a) Bayesian models often use a conjugate prior. With the aid of an example, state what it means to be conjugate, and explain why it is of practical importance. [2 marks]

A conjugate prior exist when product of likelihood  $\times$  prior results

in the same distribution as the prior, for example, Beta  $\sim$  Beta  $\times$  Binomial

in this way, we can help us sequentially update our Bayesian model and solve discrete problem, which is very practical.

3. In bias-variance decomposition, bagging can reduce variance

- (b) The Maximum a Posteriori (MAP) method is sometimes referred to as "poor man's Bayes." Explain how MAP is different to full Bayesian inference, and describe a situation in which the two methods will produce different results. [2 marks]

MAP make point estimation on a given condition, we only get the parameter with the highest probability.

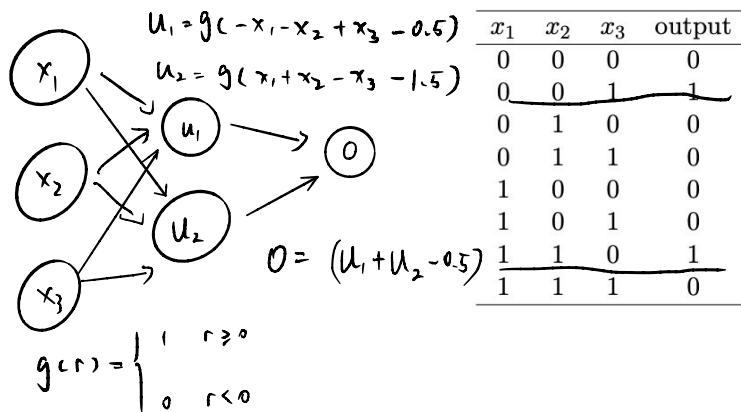
while full Bayesian inference keep the complete distribution of the posterior. when the output of the full Bayesian inference is a probability distribution but the MAP can only provide the top 1% different

## Section C: Numeric Questions [10 marks]

In this section you are asked to demonstrate your understanding of a subset of the methods that we have studied in this subject, in being able to perform numeric calculations.

### Question 6: Artificial Neural Networks [2 marks]

*Artificial neural networks (ANN)* are capable of representing arbitrary *Boolean functions*. Draw an *ANN* that represents a *Boolean function* defined by the table below. Specify the value of each *weight*, and indicate which *activation function* is used for each node. Explicitly draw *bias* terms if these are used. Note  $x_1, x_2, x_3$  are the inputs to the network.



### Question 7: Gaussian Mixture Models [3 marks]

Assume a *Gaussian mixture model* with parameters  $\theta$ . In your answers to both parts of this question below you can leave the Gaussian equation as  $N(\mathbf{x}| \dots)$ , replacing “ $\dots$ ” with the parameters of a Gaussian.

- (a) Write down the probability  $P(\mathbf{x}_i|\theta)$  of a point  $\mathbf{x}_i \in \mathbf{R}^d$ . Identify what parameters are included in vector  $\theta$  and state any constraints that apply to the parameters. [1 mark]

$$P(\mathbf{x}_i|\theta) = \sum_{j=1}^k w_j \cdot N(\mathbf{x}_i | \mu_j, \Sigma_j) = \sum_{j=1}^k P(c_j) \cdot P(\mathbf{x}_i | c_j)$$

here  $\sum_{j=1}^k P(c_j) = 1$  and  $P(c_j) \geq 0$

- (b) Write down the conditional probability of point  $i$  originating from cluster  $c$  given data  $P(z_i = c | \mathbf{x}_i, \theta)$ . [2 marks]

$$P(z_i = c | \mathbf{x}_i, \theta) = \frac{P(\mathbf{x}_i | z_i = c, \theta) \cdot P(z_i = c)}{P(\mathbf{x}_i)} = \frac{w_c \cdot N(\mathbf{x}_i | \mu_c, \Sigma_c)}{\sum_{j=1}^k w_j \cdot N(\mathbf{x}_i | \mu_j, \Sigma_j)}$$

### Question 8: Dimensionality Reduction [2 marks]

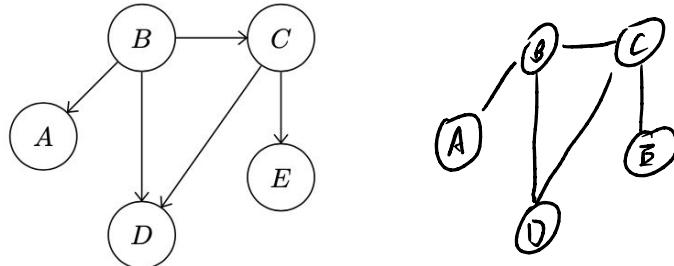
Let  $\mathbf{x}$  be a vector, and  $A$  be a matrix so that  $A\mathbf{x}$  is defined. Let  $\alpha$  be a real number. Recall that if  $A\mathbf{x} = \alpha\mathbf{x}$ , then  $\mathbf{x}$  is called an *eigenvector* of  $A$  with the corresponding *eigenvalue*  $\alpha$ .

Consider a set of  $m$  *high-dimensional points mapped onto a line* (i.e., to 1D). Let  $\mathbf{x}$  be an  $m$ -dimensional vector representing 1D coordinates of the mapped points. According to *Laplacian eigenmaps* method, the best mapping is the one that minimizes  $\mathbf{x}'L\mathbf{x}$ , where  $L$  is the *Laplacian* of the *similarity graph* constructed from *high-dimensional points*, and  $\mathbf{x}$  is restricted to have a fixed norm, e.g.,  $\mathbf{x}'\mathbf{x} = 1$ . Show that the optimal  $\mathbf{x}^*$  is an *eigenvector* of  $L$ .

$$\begin{aligned} A &= \mathbf{x}'L\mathbf{x} - \lambda(\mathbf{x}'\mathbf{x} - 1) \\ \frac{dA}{dx} &= 2L\mathbf{x} - 2x\lambda = 0 \\ L\mathbf{x} &= x\lambda \\ \mathbf{x} &\text{ is an eigenvector of } L \end{aligned}$$

### Question 9: Probabilistic Inference [3 marks]

Consider the following directed *probabilistic graphical model* (PGM) over five binary-valued random variables, denoted  $A, B, C, D$  and  $E$ .



- (a) State the form of the joint probability density function,  $P(A, B, C, D, E)$ . [1 mark]

$$P(A, B, C, D, E) = P(B) \cdot P(A|B) \cdot P(C|B) \cdot P(D|B, C) \cdot P(E|C)$$

- (b) State all the *independence relations* that hold between  $B$  and  $E$ , considering both *marginal* and *conditional* independence. [2 marks]

$$B \perp E | C ; \quad A \perp E | B ; \quad E \perp D | B$$

### Section D: Design and Application Questions [10 marks]

In this section you are asked to demonstrate that you have gained a high-level understanding of the methods and algorithms covered in this subject, and can apply that understanding. Expect your answer to each question to be from one third of a page to one full page in length. These questions may require significantly more thought than those in Sections A–C and should be attempted only after having completed the earlier sections.

### Question 10: Movie Recommendation [5 marks]

Your task is to design an automatic *movie recommendation* system for a company that provides an online movie streaming service. This company maintains a large *collection* of movies, and each movie can be accessed for a small fee. The service is only available for *registered users*. The *users preview the movies* and can pay for access to the full movie. For each *user*, the system keeps track of what movies were *previewed* and *what movies were purchased*. The aim of your system is to provide personalised *movie recommendations* to each user, as to maximise the number of purchases. Note that your system will only be applied for users that already have a history of a large number of *previews*. With all these considerations in mind, outline the design of your *movie recommendation* system.

- (a) Formulate the *movie recommendation task* as a *supervised learning* problem. Explain, what your training instances are, and what is the target variable. What features you will use, and how do you construct them from the data? [2 marks]

*my training instances are the user's previewed and purchased*

- (b) You have decided to use a kernel classifier for this purpose. Outline why using a kernel might convey an advantage for this problem over other methods, such as a linear classifiers. [1 mark]

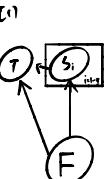
*not-linear separable → use kernel to transform  
the features to another dimension, in which they  
are linear separable.*

- (c) It is also possible to use *unsupervised learning methods* to assist with the task. What *unsupervised learning methods* might you apply for this task, and how? [2 marks]

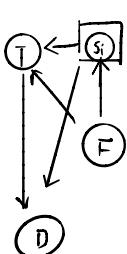
### Question 11: Modelling diseases and symptoms [5 marks]

Influenza is a common disease, commonly known as the flu. It manifests with symptoms including ~~muscle pain, coughing, fever and sneezing~~. A doctor will make a judgement about whether or not a patient has the flu, on the grounds of ~~an individual's symptoms~~. The doctor may also make use of a ~~clinical diagnostic test~~. These tests provide more reliable information than using symptoms alone, but are not perfect, in that they can return false positive or false negative results.

- (a) Draw a directed *probabilistic graphical model* (PGM) to best model the above scenario, based on the following random variables:  $F$ , whether the individual has the flu;  $S_i$ , the presence of symptom type  $i$ ,  $i = 1 \dots I$ ; and  $T$  the result of the test. All random variables are binary valued. You may choose to include other random variables, as needed. [2 marks]



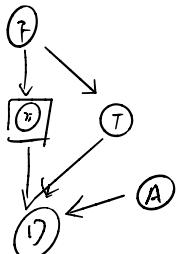
- (b) Now incorporate  $D$ , the doctor's diagnosis. How might the PGM differ if the doctor is known to be perfectly accurate, or inaccurate, when making their diagnosis based on the evidence? You may choose to include other random variables, as needed. [1 mark]



- (c) In making a diagnosis, the PGM can be used to compute a distribution over  $F$ . State the desired distribution, based on an individual attending the doctor with specific symptoms. [1 mark]

$$P(F = 1 \mid D, S_1 = s_1, S_2 = s_2, \dots, S_n = s_n)$$

- (d) The parameters of the PGM are in the form of conditional probability tables. Explain how these parameters can be learned from data, and what kind of data would be required. [1 mark]



We could use statistic inference to learn the data,  
ideally we would have a record of previous patients  
symptoms and test results and can use ML to  
find the parameters that better explain the data.