



Coláiste na Tríonóide, Baile Átha Cliath
Trinity College Dublin

Ollscoil Átha Cliath | The University of Dublin

Faculty of Engineering, Mathematics and Science

School of Computer Science & Statistics

Integrated Computer Science Programme
Year 3

Hilary Term 2019

ST3009: Statistical Methods for Computer Science

DD MMM YYYY

Venue

00.00 – 00.00

Doug Leith

Instructions to Candidates:

Attempt **all** questions.

You may not start this examination until you are instructed to do so by the invigilator.

Materials Permitted for this examination:

Non-programmable calculators are permitted for this examination – please indicate the make and model of your calculator on each answer book used.

1. (i) Define the terms “sample space” and “random event”. Give an example of each.

[5 marks]

(ii) You perform the following experiment: you take a six-sided die, and roll it. If the number that comes up is six, you stop; otherwise you repeat.

- (a) What is the probability that you roll the die once and then the experiment stops?
What is the probability that you roll the die 5 times before the experiment stops?

[5 marks]

(b) Suppose now that you roll three dice and stop when one or more of the dice come up six. What is the probability that you roll the dice 5 times before the experiment stops?

[5 marks]

(iii) Define the conditional probability of an event and state Bayes Rule. [5 marks]

(iv) Suppose two websites A and B take hotel bookings. Site A takes 60% of all bookings and site B takes 40%. However, only 75% of the bookings made on site A result in positive reviews after the hotel stay, while on site B it is 90%. Given that a booking received a positive review, what is the probability that booking was made on site B? Hint: use Bayes Rule. [5 marks]

2. (i) Define the term “random variable” and give an example. What is the probability mass function of a discrete random variable? Give an example. [5 marks]

(ii) Define what it means for two random variables to be independent. [5 marks]

Let X and Y be independent random variables that take values in the set $\{1, 2, 3\}$. Assume that X and Y are uniformly distributed on $\{1, 2, 3\}$ i.e. the probability of each value occurring is the same. Let $V = 2X + 2Y$.

(a) Are V and X independent? Explain using the definition of independence given above. [10 marks]

(b) Compute the PMF of V . [5 marks]

3. A web server is connected to the internet via two links. The server is required to serve an object consisting of n packets to a client and has to choose which link to send the n packets across. The time taken to send the n packets across link i is $T_i = \sum_{j=1}^n X_j^{(i)}$ where $X_j^{(i)}$ is a random variable whose value is the time that would be taken to send packet j across link i .

(i) Suppose the transmit times $X_j^{(1)}, j = 1, 2, \dots$ on link 1 are independent and identically distributed with mean 0.1 and variance 0.05.

- (a) What is the expected value of T_1 ? Hint: use the linearity of the expectation i.e. that $E[X+Y]=E[X]+E[Y]$. [5 marks]
- (b) What is the variance of T_1 ? Hint: use the linearity of the variance when variables are independent i.e. that $\text{var}(X+Y)=\text{var}(X)+\text{var}(Y)$. [5 marks]
- (ii) Using Chebyshev's inequality give an upper bound on the probability that T_1 is greater than or equal to $0.5n$. Recall that for random variable X Chebyshev's inequality is: $P(|X - \mu| \geq k) \leq E[(X - \mu)^2]/k^2$ where $\mu = E[X]$. [5 marks]
- (iii)
- (a) Calculate the expected value and variance of T_1/n . [5 marks]
- (b) As the number of packets n goes to infinity, explain what Chebyshev's inequality tells us about the probability that T_1/n is greater than its expected value? Hint: recall the weak law of large numbers. [5 marks]
4. Suppose that we observe data (x_i, y_i) , $i=1,2,\dots,n$ for n people sampled from the population, where x_i is the distance of the person's home from Dublin city centre, and $y_i=1$ if the i 'th person has visited TCD in the last year and 0 otherwise. We model y_i as an observation of random variable Y_i and x_i as an observation of random variable X_i . Recall that a logistic regression classifier uses the model
- $$P(Y_i=y | X_i=x_i) = 1/(1+\exp(-yz_i))$$
- with $z_i=\theta x_i$, θ a parameter and y either 0 or 1.
- (i) Assume the Y_i 's are independent and that parameter θ is known. Using the above model write an expression for the likelihood $P(\{Y_i=y_i, i=1,2,\dots,n\} | \{X_i=x_i, i=1,2,\dots,n\}, \theta)$ i.e. for the probability that data $y_i, i=1,2,\dots,n$ is observed under this model. [5 marks]
- (ii) Now suppose parameter θ is not known. What is meant by the maximum likelihood estimate of parameter θ ? [5 marks]
- (iii) Suppose that in addition to the distance of their home from Dublin we also know for person i the number w_i of their children who currently study at TCD. Describe how we can modify the above logistic regression classifier to include this information. Hint: we now have two inputs/features x_i and w_i . [5 marks]
- (iv) This type of logistic model is often said to be most effective for linearly separable data. Explain what linearly separable means, giving an example. What can happen if the data is not linearly separable? [5 marks]
- (v) Suppose the data is not linearly separable. How might we change the inputs/features used in the model to accommodate this? [5 marks]