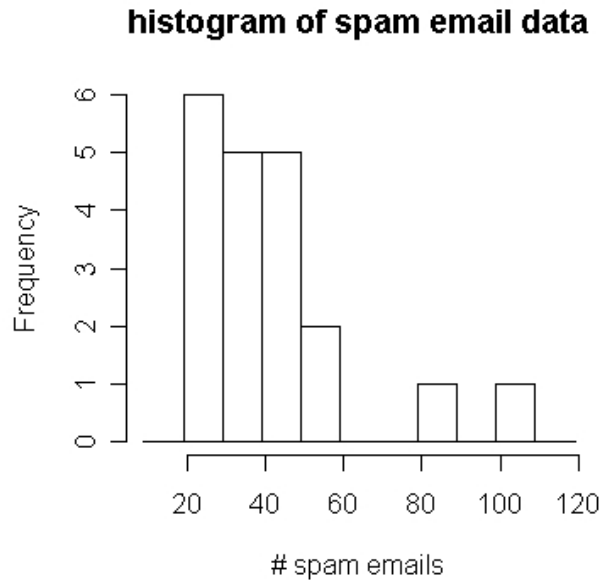


1. Harris recently installed a spam filter software, but he still saw spam emails in his inbox. He made a daily record of the number of spam emails that were delivered to his inbox over the past 20 days. The following is a frequency histogram for his data. The frequency refers to the number of days.



- a) Harris also plotted a stem-and-leaf display for the data. Which of the following is a correct stem-and-leaf display for his data? Check only one answer. [2 marks]

- Stemplot A  
 Stemplot B  
 Stemplot C

- |   |  |  |
|---|--|--|
| <p>A.</p> <p>2   011355</p> <p>3   01467</p> <p>4   12479</p> <p>5   56</p> <p>6  </p> <p>7  </p> <p>8   0</p> <p>9  </p> <p>10   5</p> | <p>B.</p> <p>2   011355</p> <p>3   01467</p> <p>4   12479</p> <p>5   56</p> <p>8   0</p> <p>10   5</p> | <p>C.</p> <p>1   05</p> <p>2   011355</p> <p>3   01467</p> <p>4   12479</p> <p>5   56</p> <p>6  </p> <p>7  </p> <p>8   0</p> |
|---|--|--|

- b) Which of the following is a correct statement about the distribution of the spam email data? Check only one answer. [2 marks]
- The distribution is roughly symmetric, and the mean is about the same as the median.
- The distribution is skewed, and the mean is larger than the median.
- The distribution is skewed, and the mean is smaller than the median.
- c) What is the third quartile of the number of spam emails? Use the stem-and-leaf display you have chosen in part (a) to answer this question. Check only one answer. [2 marks]
- 25
- 41
- 48
- 55
- d) Which of the following pairs of summary statistics best describe the center and the spread of the number of spam emails received daily? Check only one answer and explain briefly. [3 marks]
- mean and standard deviation
- mean and IQR
- median and IQR
- median and variance

Explain:

2. During the boxing week last year, a local bookstore offered discounts on a selection of books. The manager looks at the records of all the 2743 books sold during that week, and constructs the following contingency table:

	Discounted	Not discounted	Total
Paperback	790	389	1179
Hardcover	1276	288	1564
Total	2066	677	2743

- a) What percentage of paperback books sold were discounted? [1 mark]
- b) What percentage of hardcover books sold were discounted? [1 mark]
- c) What can you say about the two events:  $A$  = a book is paperback and  $B$  = the book is discounted? [1 mark]

d) Write down the marginal distribution of the book type variable. [2 marks]

3. The length of trout in a lake is normally distributed with mean  $\mu = 0.95$  feet and an unknown standard deviation  $\sigma$ . If 16% of all trout are shorter than 0.8 feet, what is the value of  $\sigma$ ? [3 marks]

4. Does how long children remain at the lunch table help predict how much they eat? Twenty toddlers at a nursery school were observed. On each toddler, the number of minutes he/she spent at the table when lunch was served and the number of calories that was consumed during lunch were measured. The two variables show a reasonably linear trend with a correlation coefficient  $r = -0.65$ .

The least-squares regression line that predicts the amount of calories consumed from the time stayed at the table during lunch has a slope of  $-3.25$  and an intercept of  $566.5$ .

a) Predict the number of calories consumed for a child who spends 25 minutes at the table during lunch. [1 mark]

b) For the following statements, check all that are correct. [3 marks]

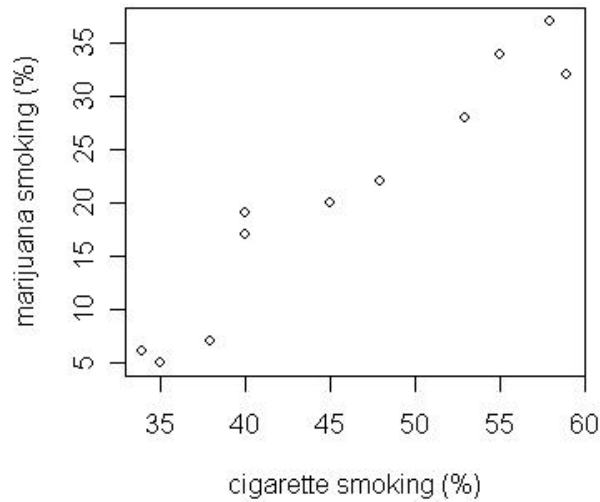
The association between the time spent at the lunch table and the amount of calories consumed is weak because  $r$  is negative.

One standard deviation (SD) increase in the number of minutes spent at the lunch table is associated with one SD decrease in the number of calories consumed.

If two variables are positively correlated, a given increase in one variable will lead to an increase in the other.

None of the above.

5. A survey was conducted in 11 countries to determine the percentage of teenagers who had smoked cigarettes and used marijuana. The scatterplot for the two variables is shown below:



One more country participated in the survey, and the percentages of teenagers who have smoked cigarettes and used marijuana were found to be 68% and 15%, respectively. The correlation coefficient  $r$  is then recalculated. How do the values of  $r$  before and after the inclusion of the new observation compare? Check only one answer and explain briefly. [3 marks]

- $r(\text{before}) < r(\text{after}) < 0$   
  $0 < r(\text{after}) < r(\text{before})$   
  $r(\text{before}) < r(\text{after}) < 1$

Explain:

6. You need to drive past two traffic lights on the way from your house to the nearest grocery store. The probability that you hit a red light is 0.5 at the first intersection and 0.4 at the second intersection. The probability that you run into a red light at both intersections is 0.25. On a random day you drive from home to that grocery store.

Define the following events:

$E_1$  = you run into a red light at the first intersection

$E_2$  = you run into a red light at the second intersection

Based on the information given, are  $E_1$  and  $E_2$  independent events? Explain briefly why or why not. [3 marks]

7. Consider the following two studies:

*Study 1:* A study compared 120 patients with brain cancer to 246 healthy patients without brain cancer. The patients' cell phone use was measured using a questionnaire. The brain cancer patients used cell phones more often, on the average.

*Study 2:* A study exposed rats to two common types of cell phone radiation for four hours a day, five days a week, for two years. One third of the rats were randomized to be exposed to analog cell phone frequency, one third to digital cell phone frequency, and one third served as controls and received no radiation. At the end of two years, their brains were examined for cancerous tumors. No statistically significant difference in the percentage of brain cancer was found among the groups.

a) Check whether each of the following statements is true or false. [1.5 marks each]

i. Study 1 shows that cell phone use causes brain cancer.

\_\_\_ True \_\_\_ False

ii. In both studies, the response variable is the presence of brain cancer.

\_\_\_ True \_\_\_ False

b) Identify the following elements for Study 2:

i. the experimental unit [1.5 marks]: \_\_\_\_\_

ii. the treatments [1.5 marks]:

\_\_\_\_\_

8. A city council was planning to turn a major street in the city from a primary traffic artery to a secondary traffic artery. It sent out a questionnaire to all the 36,589 households living in the city requesting for their input concerning the plan. Thirty four percent of the 10,375 households who returned the questionnaires opposed the plan.

For the following statements, check all that are correct. [3 marks]

\_\_\_ This survey conducted by the city council is likely to suffer from nonresponse bias.

\_\_\_ The 10,375 households that returned the questionnaires formed a random sample of the population.

\_\_\_ The percentage of the 10,375 households that opposed the plan, 34%, is a parameter.

9. In a university parking database with 5600 registered vehicles, records show that 43% of the registered vehicles are Asian makes, 23% are European makes and the remaining are American makes. Among all the 5600 cars, 20% once received a parking ticket.

a) You randomly pick three vehicles with replacement from the database (“with replacement” means any drawn vehicle will be put back to the database before the next vehicle is drawn). What is the probability that at most two of the three are American makes? [3 marks]

b) Consider a random sample of 100 vehicles selected from the database.

The sample proportion of the 100 selected vehicles that had never received a parking ticket has an approximate 95% chance of falling between \_\_\_\_\_ and \_\_\_\_\_.

Fill in the blanks and show your calculation below. [4 marks]

10. Each day the value of a particular stock goes up one unit with probability 0.3, stays the same with probability 0.5 or else goes down one unit with probability 0.2. [7 marks]

a) Consider the random variable: change in value of the stock in a day. Find the mean of this random variable.

b) The standard deviation of the change in the value of the stock in a day is given to be 0.70. A stockbroker reported that the average change in the stock value over 500 independent days exceeds 0.01 unit. Do you think what the stockbroker reported is unusual? Justify your answer using z-score or probability.