

## Lab 2: Part 2

The aim of this part of the lab exercise is to give you practical experience in working with distributions and confidence intervals using R Studio and an R Notebook.

### 1 Distributions

1. The evidence shows that 25% of people exposed to a specific virus will show symptoms. If we have a group of 5 people all exposed to this virus.
  - (a) Use the `dbinom` function to compute the probabilities that 0, 1, 2, 3, 4, 5 of the people who are exposed will show symptoms.
  - (b) What is the probability that at least one will show symptoms? (hint: use `pbinom` to help you find the probability of no one showing symptoms)
2. Assume that among diabetics the fasting blood level of glucose is approximately normally distributed with a mean of 105 mg per 100 ml and a standard deviation of 9mg per 100ml.
  - (a) Plot the density function for this distribution.
  - (b) What proportion of diabetics have levels between 90 and 125 mg per 100ml? (This quantity is represented by the area under the normal curve between  $x=90$  and  $x=125$ )
  - (c) You should see from the plot that most of the area under the normal curve is contained between the two vertical lines. Therefore we should expect that the proportion of diabetics with levels between 90 and 125 mg per 100 ml to be high. Use the `pnorm` function to calculate this proportion. (you may need to use it twice)
  - (d) What level cuts off the lower 10% of diabetics? (hint: use `qnorm`)

### 2 Confidence Intervals

1. Open a new R Studio notebook and load any libraries you plan to use
2. Read in the `auto-mpg.csv` data set into a new R notebook
3. Visualise the data and check which columns from the continuous variables look normally distributed
4. Compute a 95% confidence interval for the mean for one of the variable that appears to be normally distributed
5. The model years range from the 70s and the 80s. What proportion of cars are from the 80s?
6. Compute a 90 % confidence interval for the proportion of cars from the 80s?

OPTIONAL Compute the confidence intervals for the mean using the "bootstrap method".

## 2.1 About the data

- mpg: continuous
- cylinders: multi-valued discrete
- displacement: continuous
- horsepower: continuous
- weight: continuous
- acceleration: continuous
- model year: multi-valued discrete
- origin: multi-valued discrete
- car name: string (unique for each instance)

Dataset source: UCI Machine Learning Repository

## 3 R code

You may find this code useful to complete this worksheet:

```
1
2 #to read in a csv file
3 x<-read.csv("filename.csv")
4
5 # to explore the contents of the data
6 summary(x)
7
8 # for the binomial distribution
9 # the probability mass function for a given number of successes $v$, number of
   trials $n$ and probability of success $p$
10 dbinom(v,n,p)
11
12 #the cumulative distribution function for the same as above
13 pbinom(v,n,p)
14
15 #For the Normal distribution (the PDF) in order to find the y-values
   corresponding to a range of x values $(a,b)$
16 #with a mean of $m$ and a standard deviation of $sd$
17 dnorm(x, mean=m, sd=sd)
18
19 #the CDF for a given value of x
20 pnorm(x, mean, sd)
21
22 #the x value that corresponds to a given CDF
23 qnorm(p, mean, sd)
24
25 #for Z value the standard normal distribution x values:
26 qnorm(1- significance level)
```