

# Class test II (1 December 2016)

Please use the exam account that has already been logged in. Do not log into Windows using your own credentials. During the test you are not allowed to talk to or otherwise communicate with other students (email, instant messaging, etc.) or access the internet / course material.

At the end of the test, please log on to Moodle and upload your script by clicking on the link *Upload class test 2* on Moodle.

Please make sure you regularly save your R script, just in case RStudio crashes.

**You have to attempt all questions.**

1. Your Desktop contains two data files which you will use in this question.

The file `stations.csv` contains the list of bike stations of the Bay Area Bike Share system in the San Francisco Bay Area. It has the following columns.

<code>id</code>	Numeric identifier of the station
<code>name</code>	Name of the station
<code>lat</code>	Latitude of the station
<code>long</code>	Longitude of the station
<code>docks</code>	Number of docks at the station
<code>city</code>	City in which the station is located

The file `trips.csv` contains all trips made at the end of August 2013. It has the following columns.

<code>start_id</code>	Numeric identifier of the station where the trip started
<code>end_id</code>	Numeric identifier of the station where the trip ended
<code>start_date</code>	Date and time the trip started
<code>end_date</code>	Date and time the trip ended
<code>subscription_type</code>	User type ("Subscriber" or "Customer")

Give the R code required to answer the questions below.

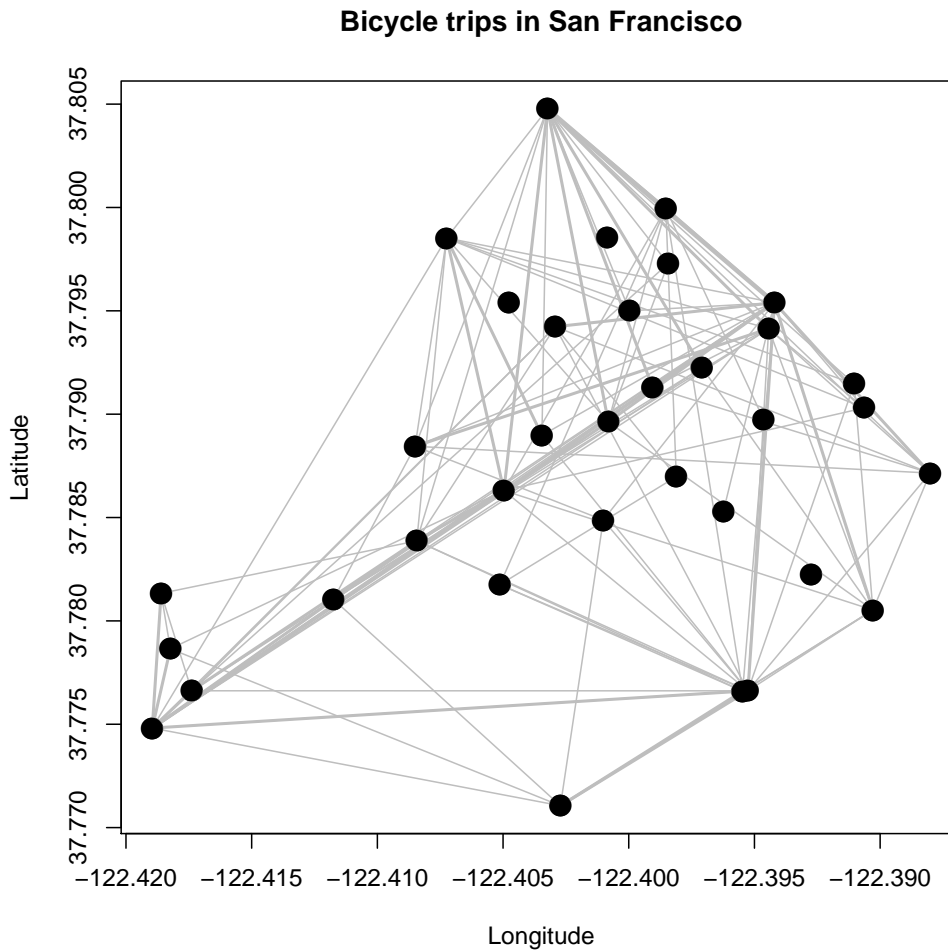
- Read both data files into R and store them in a data frame each. [4 marks]
- How many docks are there in the entire system? How many docks are there on the city of Mountain View? [2 marks]
- What proportion of trips were made by subscribers and what proportion of trips were made by customers? [2 marks]
- Find the station name of the station with the fewest docks. [2 marks]
- Create a plot of the GPS coordinates of the stations in San Francisco. The label of the x-axis should be "Longitude" and the label of the y-axis should be "Latitude". The title of the plot should be "Bicycle trips in San Francisco". [2 marks]
- Create a so-called origin-destination matrix. The matrix should be  $70 \times 70$  and the  $(i, j)$ -th entry should contain the number of trips made from station  $i$  to station  $j$ . [4 marks]
- Add lines to your plot representing the number of trips between the stations in San Francisco. Do not show lines corresponding to trips involving stations outside San Francisco.  
 If there are less than 5 trips between  $i$  and  $j$  (either way), do not show a line on your plot.  
 If there are between 5 and 9 trips between  $i$  and  $j$  (either way), draw a line of width 1.  
 If there are between 10 and 19 trips between  $i$  and  $j$  (either way), draw a line of width 2.  
 If there are more than 20 trips between  $i$  and  $j$  (either way), draw a line of width 3.

Your plot should look similar to the plot shown overleaf.

If your code from part (f) is not working you can create a fake origin-destination matrix using the code below.

```
R> od <- matrix(rpois(4900, 5), ncol=70)
```

[4 marks]



2. (a) The number of ways in which we can draw  $n$  objects from a population of  $N$  objects is given by  $\frac{N!}{(N-n)!}$ , assuming that we are retaining the order. For large  $N$  this formula suffers from numerical overflow.

Write a function `approx.permut`, which takes  $N$  and  $n$  as arguments and which returns the approximation

$$\sqrt{\frac{N}{N-n}} \cdot \left(\frac{N}{N-n}\right)^N \cdot (N-n)^n \cdot \exp(-n).$$

[5 marks]

- (b) Alternatively we can calculate  $\frac{N!}{(N-n)!}$  by rewriting it as

$$N \cdot (N-1) \cdots (N-2) \cdot (N-n+1) = \prod_{i=0}^{n-1} (N-i) \quad (1)$$

This formula is exact and numerically stable. Write a function `exact.permut` which takes  $N$  and  $n$  as arguments and which returns the value of  $\frac{N!}{(N-n)!}$  calculated as set out in equation (1).

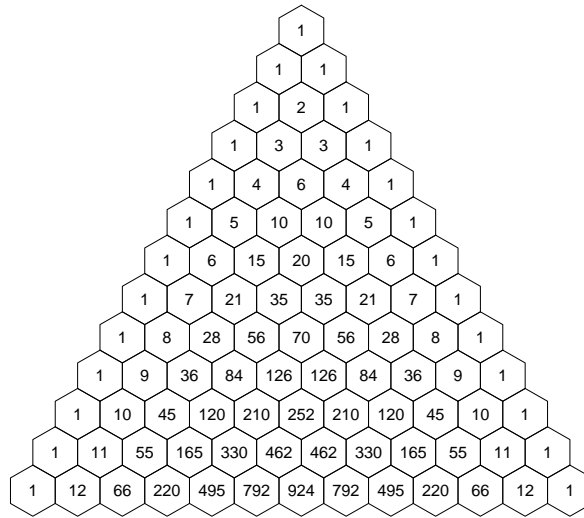
*Do not calculate the result as  $N!$  divided by  $(N-n)!$ , as this would be numerically unstable.* [5 marks]

- (c) Calculate  $\frac{200!}{(200-10)!}$  using your functions `approx.permut` from part (a) and `exact.permut` from part (b).  
You can compare your results to  $\exp(\text{lgamma}(200+1) - \text{lgamma}(200-10+1))$ . [2 marks]

3. Pascal's triangle is an illustration of Pascal's rule for binomial coefficients:

$$\binom{k}{i} = \binom{k-1}{i-1} + \binom{k-1}{i}.$$

Consider a triangle such that the top row consists of  $\binom{0}{0} = 1$  and the second row consists of  $\binom{1}{0} = 1$  and  $\binom{1}{1} = 1$ . The  $k$ -th row consists of  $\binom{k}{0}, \binom{k}{1}, \dots, \binom{k}{k}$ . The recursive formula from above just implies that each entry in the triangle is the sum of the two entries immediately above it (see below).



2. For  $t = 1, 2, \dots, T$

- i. Draw  $\Delta S_t$  from the binomial distribution given in (2).
- ii. Draw  $\Delta I_t$  from the binomial distribution given in (3).
- iii. Update  $S_t$ ,  $I_t$  and  $R_t$  as set out in (4).

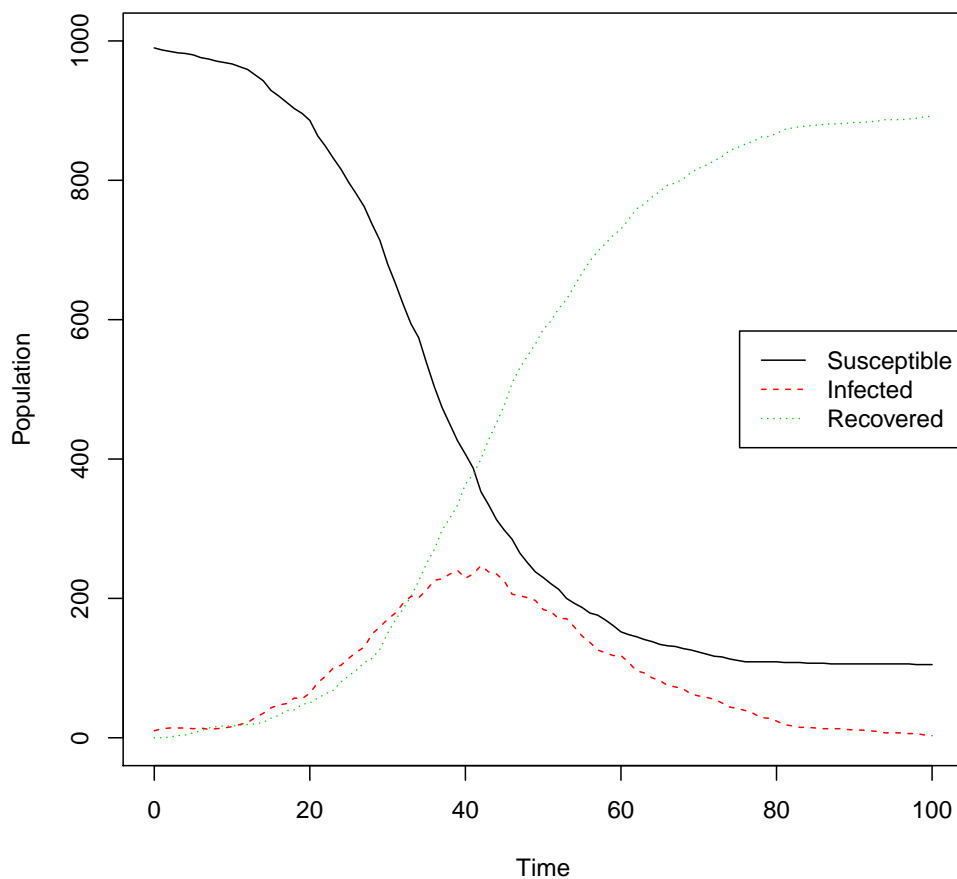
(a) Write a function `simulate.sir` which simulates from the SIR model as set out above. The function should take the parameters `N`, `initial.infected` ( $I_0$ ), `alpha` ( $\alpha$ ), `beta` ( $\beta$ ) and `T` as arguments. It should return a data frame with  $T + 1$  rows, in which the  $(t + 1)$ -st row contains the values of  $S_t$ ,  $I_t$  and  $R_t$ .

Your function should check whether  $0 \leq I_0 \leq N$ ,  $0 < \alpha < 1$  and  $\beta > 0$  and produce an error message if this is not the case.

The default values should be  $N = 1000$ ,  $I_0 = 0.01N$  (rounded to the nearest integer),  $\alpha = \frac{1}{10}$ ,  $\beta = \frac{1}{4}$  and  $T = 100$ .

*Hint: You can draw one realisation from the  $Bi(n, \theta)$  distribution using the R function `rbinom(1, n, theta)`. [10 marks]*

(b) Use your function `simulate.sir` to create one realisation from the SIR model (using the above default values for the parameters) and plot the three columns as three curves in one plot as shown below. The curves should use a different colour and/or line type. Add a legend to your plot.



*Please note that the SIR model is a stochastic model, so every realisation will look slightly different.*

[5 marks]

(c) At the end of the epidemic (which we assume is at time  $T = 100$ ) not every individual will have had the infection. It is of epidemiological interest to determine the proportion of individuals who have had the infection,  $R_{100}/N$ .

Create 50 realisations from the above model, each created by calling the function `simulate.sir` from part (a), again using the above default values for the parameters. Calculate the mean value of  $R_{100}/N$  from the simulations. [5 marks]

*If you do not get your code to work for part (a) you can use the following function instead of `simulate.sir` to answer parts (b) and (c).*

```
R 2 fake.sir <- function(N=1000, T=100) {
R 3   a <- N * runif(1)
R 4   infected <- 25 * dnorm(0:T, 45, 25) * a
```

```
R 5 | susceptible <- N - pnorm(0:T, 45, 25) * a
R 6 | data.frame(susceptible=susceptible, infected=infected, recovered=N-susceptible-infected)
R 7 | }
```

---

---

Total (Honours):	60
Total (MSc):	60