# MAS2602: Computing for Statistics

Newcastle University lee.fawcett@ncl.ac.uk

Semester 1, 2019/20

<ロト < 部 ト < 注 ト < 注 ト 三 三 のへで</p>

- Statistics classes run in teaching weeks 5–9
- Lecturer is Dr. Lee Fawcett (lee.fawcett@ncl.ac.uk)
- Schedule: 4 lectures, 4 practicals, 1 revision class, 1 class test

<ロト 4 目 ト 4 三 ト 4 三 ト 9 0 0 0</p>

 Office hours: Tuesdays 3-4; room 2.07 Herschel (also MAS3902: Tuesdays 2-3 and Thursdays 3-4)

# Schedule

Week 5	Mon 28 Oct	11–12	Lecture (Herschel LT1)
	Fri 1 Nov	10–12	Practical (Herschel PC)
Week 6	Mon 4 Nov	11–12	Lecture (Herschel LT1)
	Wed 6 Nov	11–1	Practical (Herschel PC)
Week 7	Mon 11 Nov	11–12	Lecture (Herschel LT1)
	Thu 14 Nov	11–1	Practical (Herschel PC)
Week 8	Mon 18 Nov	11–12	Lecture (Herschel LT1)
	Thu 21 Nov	11–1	Practical (Herschel PC)
	<i>Thu 21 Nov</i>	<b>3pm</b>	<i>Assignment due</i>
	Fri 22 Nov	1–2	Revision (Herschel LT2)

Week 9 Tue 26 Nov 9–10 Test (Herschel PC)

### Assessment

### Assignment (a.k.a. "mini-project")

- Due Thursday 21 November, 3pm
- Worth 10% of the module marks
- **Class test** 
  - Tuesday 26 November, 9am one hour long
  - Worth 30% of module marks
  - Open-book test
  - You will write one or two short R programs during the test

Help available:

- Office hours
- Demonstrators in practical sessions
- Books see recommendations in booklet

- Email Lee, or just pop in!
- Blackboard and dedicated webpage

In this module, deadline extensions can be requested for the final project (by means of submitting a PEC form), and work submitted within 7 days of the deadline without good reason will be marked for reduced credit. This module also contains tests worth more than 10% for which rescheduling can be requested (by means of submitting a PEC form). There are mini-projects (worth 10% each) for which it is not possible to extend deadlines and for which no late work can be accepted. For details of the policy (including procedures in the event of illness etc.) please look at the School web site:

http://www.ncl.ac.uk/maths/students/teaching/homework/

For problems with deadlines, speak to your personal tutor and prepare a PEC form

# Lecture 1: Introduction and Simulation of Random Variables

<□▶ < □▶ < 三▶ < 三▶ = 三 のへぐ

# Introduction

In this part of the module we will do statistics with R:

- R is the foremost tool in modern computational statistics
- Using R teaches general concepts in programming
- It can be used to illustrate mathematical ideas in probability and statistics

Today's lecture: simulating random variables

1 Simulating random variables seen in MAS1604

2 Using this to simulate more complicated probability models

<ロト 4 目 ト 4 三 ト 4 三 ト 9 0 0 0</p>

Friday's practical:

- 1 Revision of R from MAS1802
- 2 Putting today's material into practice

## Introduction

In this part of the module we will do statistics with R:

- R is the foremost tool in modern computational statistics
- Using R teaches general concepts in programming
- It can be used to illustrate mathematical ideas in probability and statistics

Today's lecture: simulating random variables

- 1 Simulating random variables seen in MAS1604
- 2 Using this to simulate more complicated probability models

<ロト 4 目 ト 4 三 ト 4 三 ト 9 0 0 0</p>

Friday's practical:

- 1 Revision of R from MAS1802
- 2 Putting today's material into practice

# Introduction

In this part of the module we will do statistics with R:

- R is the foremost tool in modern computational statistics
- Using R teaches general concepts in programming
- It can be used to illustrate mathematical ideas in probability and statistics

Today's lecture: simulating random variables

- 1 Simulating random variables seen in MAS1604
- 2 Using this to simulate more complicated probability models

#### Friday's practical:

- 1 Revision of R from MAS1802
- 2 Putting today's material into practice

# The binomial distribution

If  $X \sim Bin(n, p)$  then X has PMF (probability mass function) given by

$$p_X(k) = \binom{n}{k} p^k (1-p)^{n-k}, \text{ for } k = 0, 1, \dots, n.$$

There is no closed formula for the CDF (cumulative distribution function).

R commands:

- dbinom calculate PMF
- pbinom calculate CDF
- rbinom generate random sample

# The binomial distribution

If  $X \sim Bin(n, p)$  then X has PMF (probability mass function) given by

$$p_X(k) = \binom{n}{k} p^k (1-p)^{n-k}, \text{ for } k = 0, 1, \dots, n.$$

There is no closed formula for the CDF (cumulative distribution function).

▲ロト ▲ □ ト ▲ 三 ト ▲ 三 ト ○ ○ ○ ○ ○ ○

R commands:

- dbinom calculate PMF
- pbinom calculate CDF
- rbinom generate random sample

# R commands for the binomial distribution

1dbinom (5, 10, 0.7) $\# \Pr(X = 5), X \sim Bin(10, 0.7)$ 2pbinom (4, 10, 0.7) $\# \Pr(X \le 4), X \sim Bin(10, 0.7)$ 3rbinom (50, 10, 0.7)# Sample a Bin(10, 0.7) distribution 50 times

### Creating a bar plot



Х

# The geometric distribution

If  $Y \sim Geom(p)$  then Y has PMF and CDF given by

$$p_Y(k) = (1-p)^{k-1}p,$$
  
 $F_Y(k) = 1 - (1-p)^k, \text{ for } k = 1, 2, \dots.$ 

Note that Y takes values in  $1, 2, 3, \ldots$ :

- R uses a slightly different definition
- We use the definition that Y is the number of Bernoulli trials with up to and including first success
- R counts number of trials up to, but not including the first success, so in R geometric random variables take values 0, 1, 2, . . ..

### Adjust the arguments to account for different definition:

1	dgeom(4, 0.2)		$\# \Pr(Y = 5), Y \sim \textit{Geom}(0.2)$
2	<b>pgeom</b> (2, 0.2)		$\# \Pr(Y \leq 3), Y \sim \textit{Geom}(0.2)$
3	1 + rgeom(100,	0.2)	# Sample a $Geom(0.2)$ distribution 100 times

Here's a function to replace dgeom with our definition of the geometric distribution:

```
mydgeom = function(x, p) {
    dgeom(x-1, p)}
```

◆□ > ◆□ > ◆豆 > ◆豆 > ・豆 - つへ⊙

If  $Z \sim Po(\lambda)$  then Z has PMF given by

$$p_Z(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$
, for  $k = 0, 1, 2, \dots$ 

There is no closed formula for the CDF (cumulative distribution function).

▲ロト ▲ □ ト ▲ 三 ト ▲ 三 ト ○ ○ ○ ○ ○ ○

R commands:

- dpois calculate PMF
- ppois calculate CDF
- rpois generate random sample

### R commands for the Poisson distribution

1dpois (5, 3.5) $\# \Pr(Z = 5), Z \sim Po(3.5)$ 2ppois (2, 3.5) $\# \Pr(Z \le 2), Z \sim Po(3.5)$ 3rpois (100, 3.5) $\# \operatorname{Sample}$  a Po(3.5) distribution 100 times

# Summary

Distribution	Binomial	Poisson	Geometric 🙎
PMF	dbinom()	dpois()	dgeom()
CDF	<pre>pbinom()</pre>	ppois()	pgeom()
sample	rbinom()	<pre>rpois()</pre>	rgeom()

R has functions for the uniform, exponential and normal distributions:

Distribution	Uniform	Exponential	Normal
PDF	dunif()	dexp()	dnorm()
CDF	<pre>punif()</pre>	pexp()	pnorm()
quantile	qunif()	qexp()	qnorm()
sample	runif()	<pre>rexp()</pre>	<pre>rnorm()</pre>

<ロト < 部 ト < 注 ト < 注 ト 三 三 のへで</p>



For the standard uniform (U(0,1)) and standard normal (N(0,1)) distributions you don't need to provide the parameters a = 0, b = 1 and  $\mu = 0$ ,  $\sigma = 1$  respectively. For example:

- コント 4 日 > ト 4 日 > ト 4 日 > - シックク

1	runif(20)	# Samples a $U(0,1)$ distribution 20 times
2	<b>pnorm</b> (1.96)	$\# \Pr(Z < 1.96), Z \sim N(0, 1)$
3	[1] 0.9750021	

# Quantiles

The quantile functions qunif, qexp, qnorm solve equations like

 $F_X(\alpha) = p$ 

\*ロ \* \* @ \* \* ミ \* ミ \* ・ ミ \* の < @

for  $\alpha$  given a probability  $\textbf{\textit{p}}.$  For example:

```
1 qnorm(0.9750021)
2 [1] 1.96
```

# Quantiles

#### Example

Suppose annual maximum wave heights observed off the coast at a flood-prone town are assumed Normally distributed, with mean 2 metres and standard deviation 0.5 metres.

- (a) Write down the R command to find the probability that the largest wave height next year will exceed 3.25 metres.
- (b) Write down the R command to estimate the height of a new sea wall such that we might expect the town to be flooded, on average, once per century. Why might our modelling assumption be invalid?



wave height

Could work out the old-fashioned way:

$$\Pr(X > 3.25) = \Pr\left(Z > \frac{3.25 - 2}{0.5}\right)$$
  
=  $\Pr(Z > 2.5)$ 

$$= 1 - \Pr(Z \le 2.5) = 1 - 0.994 = 0.006.$$

Or could just use R:

1 **1-pnorm**(3.25, 2, 0.5) 2 [1] 0.006209665

▲ロト ▲ □ ト ▲ 三 ト ▲ 三 ト ● ● ● ● ●

Could work out the old-fashioned way:

$$\Pr(X > 3.25) = \Pr\left(Z > \frac{3.25 - 2}{0.5}\right)$$
  
=  $\Pr(Z > 2.5)$ 

 $= 1 - \Pr(Z \le 2.5) = 1 - 0.994 = 0.006.$ 

Or could just use R:

Could work out the old-fashioned way:

$$\Pr(X > 3.25) = \Pr\left(Z > \frac{3.25 - 2}{0.5}\right)$$
  
=  $\Pr(Z > 2.5)$ 

$$= 1 - \Pr(Z \le 2.5) = 1 - 0.994 = 0.006.$$

Or could just use R:

1 **1-pnorm**(3.25, 2, 0.5) 2 [1] 0.006209665



€ 990

### Similarly:

1 qnorm(0.99, 2, 0.5) 2 [1] 3.163174

# A more advanced model

- Number of arrivals per day at an IT help-desk is modelled using a Poisson distribution
- Mean of the Poisson distribution might vary from day-to-day
- Suppose the number of arrivals *X* ~ *Po*(Λ)
- $\Lambda$  is itself a random variable, with  $\Lambda \sim Exp(c)$  for a constant c = 0.05

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

What can we say about the distribution of X?

- 1 What are the expectation and variance of X?
- 2 What is Pr(X > 30)?

### Arrival model – R code



### Arrival model – R code

Bar plot:





# Using simulated samples

- The sample mean is an approximation to the distribution mean
- The same applies for the sample variance
- To calculate Pr(X > 30) we count the proportion of times this occurs in the sample

We expect the approximation to improve as we increase the sample size

< ロ > < 回 > < 三 > < 三 > < 三 > < 三 > < ○ < ○</p>

# Using simulated samples



▲ロト ▲ □ ト ▲ 三 ト ▲ 三 ト ● ● ● ● ●

### Suppose

- $\mu_1, \mu_2, \sigma_1 > 0, \sigma_2 > 0$  are fixed constants,
- $w_1, w_2$  are positive constants with  $w_1 + w_2 = 1$ .

Consider the following function:

$$f(x) = w_1 f_1(x) + w_2 f_2(x)$$

where  $f_1$  and  $f_2$  are the density functions for  $Z_1 \sim N(\mu_1, \sigma_1^2)$  and  $Z_2 \sim N(\mu_2, \sigma_2^2)$  respectively.

Check that f(x) represents a valid probability density function.

# Mixtures of normal distributions

First note that  $f(x) \ge 0$  everywhere. Also

$$\int_{-\infty}^{\infty} f(x)dx = w_1 \int_{-\infty}^{\infty} f_1(x)dx + w_2 \int_{-\infty}^{\infty} f_2(x)dx$$
$$= w_1 + w_2 = 1.$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

# Mixtures of normal distributions

First note that  $f(x) \ge 0$  everywhere. Also

$$\int_{-\infty}^{\infty} f(x)dx = w_1 \int_{-\infty}^{\infty} f_1(x)dx + w_2 \int_{-\infty}^{\infty} f_2(x)dx$$
$$= w_1 + w_2 = 1.$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

If a random variable X has PDF corresponding to f(x) we say it is a mixture of normal distributions  $N(\mu_1, \sigma_1^2)$  and  $N(\mu_2, \sigma_2^2)$  with weights  $w_1, w_2$ .

Note that X is not the sum of two normal random variables i.e.

$$X \neq w_1 Z_1 + w_2 Z_2$$
 where  $Z_i \sim N(\mu_i, \sigma_i^2), i = 1, 2.$ 

<ロト 4 目 ト 4 三 ト 4 三 ト 9 0 0 0</p>

Example

For example: suppose  $\mu_1 = 3$ ,  $\sigma_1 = 1$  and  $\mu_2 = 6$ ,  $\sigma_2 = 2$  with  $w_1 = w_2 = 1/2$ .



Х

- **1** Sample a random variable  $J \in \{1, 2\}$  such that  $Pr(J = 1) = w_1$  and  $Pr(J = 2) = w_2$ .
- 2 The random variable *J* tells you which component of the mixture to sample *X* from.
- 3 If J = 1 then sample X from  $N(\mu_1, \sigma_1^2)$ , but if J = 2 then sample X from  $N(\mu_2, \sigma_2^2)$ .

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

# Sampling from a mixture distribution

```
normal.mixture = function(n, mu1, sig1, w1, mu2, sig2,
1
       w2) {
       p = c(w1, w2)
       x = vector(mode = 'numeric', length = n)
3
4
       for (i in 1:n) {
            j = sample(c(1, 2), 1, prob = p)
5
            if (i = 1) {
6
                x[i] = rnorm(1, mu1, sig1)
7
            }
8
            else {
9
                x[i] = rnorm(1, mu2, sig2)
10
            }
11
12
       return(x)
13
14
```