

BIOINFORMATICS MSc PROBABILITY AND STATISTICS

SPLUS SHEET 1

A data set containing a segment of human chromosome 13 containing the BRCA2 breast cancer gene; it was obtained from the National Center for Biotechnology Information (NCBI) website

<http://www.ncbi.nlm.nih.gov/>

and is held at the Probability and Statistics course website

<http://stats.ma.ic.ac.uk/~das01/BioinformaticsMSc/>

Using a Web browser, go to this link and then click on the link

BRCA2 segment direct from the NCBI site

in the **Data Sets** segment. Two processed versions of this data set are also available;

- *A processed version with base characters* has just the 127079 base characters
- *Another processed version* contains a numerical representation of the sequence. This is the file to download.

Click on the link, and when the file has fully downloaded, use **File -> Save** to save the file as **z74739.txt** in the **c:\TEMP** directory on the hard drive. Of course, you may ultimately save the data to a floppy disk, or to your home directory. It is possible (and quite straightforward) to read in data from different types of file apart from plain text files, and from files with different formats. The numerical data in the file **z74739.txt** are coded so that

1 = A 2 = C 3 = G 4 = T

1. STARTING AN SPLUS SESSION

Double click the **SPLUS** icon, or start the program from the **Start** menu. When the program has opened up the *Object Browser* window, start a *Commands Window* by using the **Window -> Commands Window** pull down menus.

2. LOADING THE DATA INTO SPLUS

There are two ways to load the data. First using the **File -> Import Data -> From File** pull down menus, find the **c:\TEMP\z74739.txt** file using the dialog box (remember that the file extension is .txt, so the **SPLUS** dialog box may not find the file initially). However, when you find the file and click **Open**, the data will be downloaded into an **SPLUS** data frame called **z74739** - if you return to the *Commands Window* and type

```
>length(z74739[,1])
```

(always press **Enter** or **Return** at the end of a line command) you get the response

```
[1] 127079
```

which is the length of the sequence. A line command that achieves the same outcome is

```
>z74739<-importData('c:\\TEMP\\z74739.txt',type='ASCII')
```

To read the data into a vector, type the following at the command line:

```
>brca2<-z74739[,1]
```

which creates the vector **brca2** containing the data. The other way to read in the data is to type a command at the command line

```
>brca2<-scan('c:\\TEMP\\z74739.txt')
```

which creates the same **brca2** vector. To check the steps have worked correctly, type

```
>length(brca2)
```

you should again get the response

```
[1] 127079
```

3. SUMMARY ANALYSIS

For a simple summary of the sequence, type

```
>table(brca2)
```

You should get the response

```
 1     2     3     4
38514 24631 25685 38249
```

which is the breakdown of the sequence by base, that is there are 38514 As, 24631 Cs, 25685 Gs and 38249 Ts. To analyse a sub-sequence, say the first 50000 bases, type

```
>table(brca2[1:50000])
```

that is, look at only the positions 1 to 50000 in the **brca2** vector. You then get the response

```
 1     2     3     4
14901 10397 9908 14794
```

Therefore, to obtain the sample probabilities, of each base, type

```
>brca2.p <- table(brca2)/length(brca2)
```

```
>brca2.p
```

which calculates the probability vector **brca2.p** vector by dividing the values obtained by the **table** command by the total length of the sequence. You should get the response

```
 1          2          3          4
0.3030713 0.1938243 0.2021184 0.300986
```

Note that the probabilities are not (even approximately) equal. Note also that by typing

```
>brca2sub.p <- table(brca2[1:50000])/length(brca2[1:50000])
```

```
>brca2sub.p
```

which calculates the probability vector **brca2sub.p** vector by repeating the calculation for the first 50000 bases, you get the response

```

      1      2      3      4
0.29802 0.20794 0.19816 0.29588

```

indicating a similar collection of probabilities.

4. ADJACENT BASE PAIRS ANALYSIS

Suppose now that an analysis of adjacent pairs in the sequence is required. Type the following sequence of commands (take care not to miss any brackets):

```

>brca2.mat <- matrix(0,4,4)
>brca2.pmat <- matrix(0,4,4)
>brca2.totals <- table(brca2)
>for(i in 2:length(brca2))
+ {brca2.mat[brca2[i-1],brca2[i]] <- brca2.mat[brca2[i-1],brca2[i]]+1}
>for(i in 1:4) {brca2.pmat[i,] <- brca2.mat[i,]/sum(brca2.mat[i,])}

```

(the + sign in line 4 appears automatically when you hit return from line 3). These commands construct a matrix counting the number of adjacent pairs of each type are present in the sequence, that is a matrix

		Next base				
		A	C	G	T	
Base	A	n_{AA}	n_{AC}	n_{AG}	n_{AT}	n_A
	C	n_{CA}	n_{CC}	n_{CG}	n_{CT}	n_C
	G	n_{GA}	n_{GC}	n_{GG}	n_{GT}	n_G
	T	n_{TA}	n_{TC}	n_{TG}	n_{TT}	n_T

where the row totals are just the total numbers of bases of that type (as calculated in step 2 to produce vector **brca2.totals**). Don't worry too much about the exact meaning of each command at this stage.

The third step (which essentially counts the number of adjacent AA, AC, AG etc pairs) may take a while. The fourth step divides each row in the matrix of counts by the row totals, to produce the final matrix of probabilities:

```

> brca2.pmat
      [,1]      [,2]      [,3]      [,4]
[1,] 0.3451732 0.1628239 0.23978294 0.2522200
[2,] 0.3525638 0.2431895 0.05440299 0.3498031
[3,] 0.3011485 0.2029200 0.23262605 0.2633054
[4,] 0.2300975 0.1871421 0.23880363 0.3439567

```

Note that the row sums in this matrix are 1 (by construction). From this analysis we see that “transitions” from C to G in the sequence are relatively rare.

Two things to consider:

- (a) are the transition probabilities approximately equal in different segments of the entire sequence ?
- (b) are the coding regions/exons/introns (that can be identified from the file from the Bioinformatics MSc.page (or the NCBI site) fundamentally different in terms of their composition by base ?.

5. PROBABILITY DISTRIBUTION CALCULATIONS

SPLUS has many facilities for carrying out calculations for probability distributions. There are functions that calculate the probability mass function f_X and the discrete cumulative distribution function F_X for DISCRETE random variables, and the probability density function f_X and the continuous cumulative distribution function F_X for CONTINUOUS random variables. Also, **SPLUS** has functions that allows you to simulate random numbers from many standard probability distributions.

The probability distributions that **SPLUS** has specially written functions for include the following:

DISCRETE DISTRIBUTIONS

Binomial
Geometric
Negative Binomial
Poisson

CONTINUOUS DISTRIBUTIONS

Uniform
Exponential
Gamma
Chi-squared
Beta
Normal
Student-t
Cauchy
F

The task today is to use the **SPLUS** functions to carry out probability calculations for these distributions

There are four basic functions that are used for probability calculations: using the binomial distribution as an example, the four functions are

`dbinom`
`pbinom`
`qbinom`
`rbinom`

notice the first letters d,p,q and r; these letters determine the type of operation that is being carried out. For different distributions, these first letters will always indicate the same type of operation; the last part of the function name determines which distribution is to be used. Specifically

`dbinom` :

computes the probability mass function

`pbinom` :

computes the discrete cumulative distribution function

`qbinom` :

computes the inverse cumulative distribution function

`rbinom` :

simulates a random sample from the distribution

If you type

```
>help(dbinom)
```

at the command line, a help screen explains how each function is used. Each function takes a number of arguments that will be described below.

We will initially concentrate on the binomial distribution, but the commands issued are essentially the same for each distribution we wish to use. Recall that the binomial distribution has the following probability mass function:

$$f_X(x) = \binom{n}{x} \theta^x (1 - \theta)^{n-x} = \frac{n!}{x!(n-x)!} \theta^x (1 - \theta)^{n-x} \quad x \in \{0, 1, \dots, n\}$$

for parameters n (a positive integer) and θ (a probability lying between 0 and 1). In **SPLUS**, the parameter θ is written as p

6. MASS AND DENSITY FUNCTION CALCULATIONS.

The probability mass function for the binomial distribution is obtained using the function

```
dbinom(x, size, prob)
```

where x is the vector of points at which we wish to evaluate the mass function, $size$ is the parameter n , and $prob$ is the parameter p (or θ). Hence to evaluate all the probabilities in a *Binomial*(10,0.3) distribution, we issue the following sequence of commands at the command line:

```
>x <- c(0:10)
>n <- 10
>p <- 0.3
>dbinom(x,n,p)
```

for which the response is

```
[1] 0.0282475249 0.1210608210 0.2334744405 0.2668279320 0.2001209490 0.1029193452
[7] 0.036756909 0.0090016920 0.0014467005 0.0001377810 0.0000059049
```

What this calculation has done is to create a vector of the integers from 0 to 10 (the range of this distribution), then specified $n = 10$, then specified $p = 0.3$, and then evaluated the mass function at each value in the vector. That is we have evaluated

$$f_X(0) = P[X = 0] = \binom{10}{0} (0.3)^0 (1 - 0.3)^{10-0} = 0.0282475249$$

$$f_X(1) = P[X = 1] = \binom{10}{1} (0.3)^1 (1 - 0.3)^{10-1} = 0.1210608210$$

and so on. We can easily assign the probabilities to a vector, and plot the resulting function:

```
>y <- dbinom(x,n,p)
>plot(x,y)
```

for which the response is a point plot of the mass function.

The d- functions for the various distributions, for example,

`dbinom, dgeom, dnbinom, dpois, dexp, dgamma, dnorm, ...`

and so on all have this same basic syntax - the only difference is that the parameters for each distribution change. For example, for the *Gamma* distribution, which we have seen with pdf

$$f_X(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \quad x > 0$$

has an **SPLUS** function

`dgamma(x, shape, rate=1)`

where *shape* is the α parameter, and *rate* is the β parameter (which has the default value 1 in the function). Hence to plot four *Gamma* pdfs

`Gamma(2,2)` `Gamma(2,1)` `Gamma(2,0.5)` `Gamma(4,2)`

on the range we use the following commands

```
>x <- c(0:1000)/100
>y1 <- dgamma(x,2,2)
>y2 <- dgamma(x,2,1)
>y3 <- dgamma(x,2,0.5)
>y4 <- dgamma(x,4,2)
>plot(x,y1,type='l')
>lines(x,y2,lty=2)
>lines(x,y3,lty=3)
>lines(x,y4,lty=4)
```

for which the response is a series of line plots of the pdf. The x vector created is a series of 1000 points equally spaced on 0 to 10, and y_1, y_2, y_3 and y_4 are the four evaluated pdf curves. For information the **type="l"** command produces a line (rather than a point) plot, **lines** adds a line to the current plot, and the **lty=2,3,4** commands produce different line styles.

For the different probability models, you need to use the **help** command to find out the precise syntax and parameter specification for each distribution.

7. DISTRIBUTION FUNCTION CALCULATIONS.

Again we begin with the Binomial model, for which the cumulative distribution function is obtained using the function

`pbinom(q, size, prob)`

where q is the vector of points (or *quantiles*) at which we wish to evaluate the cdf, *size* is the parameter n , and *prob* is the parameter p (or θ). Hence to evaluate the cdf in a *Binomial* (10,0.3) distribution, we use the following commands:

```
>x <- c(0:10)
>n <- 10
>p <- 0.3
>pbinom(x,n,p)
```

(with x acting as q in the specification for convenience) for which the response is

```
[1] 0.02824752 0.14930835 0.38278279 0.64961072 0.84973167 0.95265101
[7] 0.98940792 0.99840961 0.99985631 0.99999410 1.00000000
```

Again, this calculation has created a vector x of the integers from 0 to 10, then specified $n = 10$, then specified $p = 0.3$, and then evaluated the cdf function at each value in the vector. That is we have evaluated

$$F_X(0) = P[X \leq 0] = P[X = 0]$$

$$F_X(1) = P[X \leq 1] = P[X = 0] + P[X = 1]$$

and so on. Again, we assign the cumulative probabilities to a vector, and plot the resulting function:

```
>y <- pbinom(x,n,p)
>plot(x,y)
```

for which the response is a point plot of the cdf.

The p - functions for the various distributions, for example,

`pbinom`, `pgeom`, `pnbinom`, `ppois`, `pexp`, `pgamma`, `pnorm`, ...

and so on all have this same basic syntax - the only difference is that the parameters for each distribution change. For continuous distributions, we proceed as above

```
>x <- c(0:1000)/100
>y <- pgamma(x,2,2)
>plot(x,y,type='l')
```

produces a plot of the required cdf

Again, for the different probability models, you need to use the **help** command to find out the precise syntax and parameter specification for each distribution.

8. INVERSE DISTRIBUTION FUNCTION CALCULATIONS.

Again we begin with the Binomial model, for which the inverse cumulative distribution function, that is the function that solves the equation

$$F_X(x) = p_0$$

for x with p_0 fixed, is obtained using the function

```
qbinom(p, size, prob)
```

where $p \equiv p_0$ is the probability at which we wish to evaluate the inverse cdf. This calculation is very important in many statistical problems. To evaluate the inverse cdf in a *Binomial* (10, 0.3) distribution, we use the following commands - we will use x to replace p as the argument of the function, to avoid confusion:

```

>p0 <- 0.265
>n <- 10
>p <- 0.3
>x_qbinom(p0,n,p)
>x

```

for which the response is

```
[1] 2
```

Again, we can use a vector argument to this function, solve for x , and plot the resulting function:

```

>p0 <- c(1:100)/100
>x <- qbinom(p0,n,p)
>plot(p0,x)
>plot(x,p0)

```

The q - functions for the various distributions, for example,

```
qbinom, qgeom, qnbinom, qpois, qexp, qgamma, qnorm, ...
```

and so on all have this same basic syntax but slightly different parameter specifications. For continuous distributions, we proceed as above for the *Gamma*(2,2) distribution:

```

>p0 <- c(0:1000)/1000
>x <- qgamma(p0,2,2)
>plot(p0,x,type='l')
>plot(x,p0,type='l')

```

produces a plot of the required inverse cdf. For the different probability models, you need to use the **help** command to find out the precise syntax and parameter specification for each distribution.

9. RANDOM NUMBER SIMULATION

It is often useful to be able to generate a random sample from a given probability distribution. Concentrating first on the *Binomial*(10,0.3) we use the function

```
rbinom(n, size, prob)
```

where n is the required simulated sample size. To simulate a sample of size 500 from this Binomial model and store it in vector x , and then to plot a histogram of this simulated data, we can issue the following commands:

```

>n <- 10
>p <- 0.3
>x <- rbinom(500,n,p)
>hist(x)

```

which produces a histogram. We can change the number and/or positions of bars or “bins” in the histogram by using the commands


```
>hist(x,nclass=5)
>hist(x,breaks=c(0:10))
```

The `r-` functions for the various distributions, for example,

```
rbinom, rgeom, rnbinom, rpois, rexp, rgamma, rnorm, ...
```

and so on all have this same basic syntax but slightly different parameter specifications. For continuous distributions, we proceed identically to the discrete case: for the *Gamma*(2,2) distribution:

```
>x <- rgamma(500,2,2)
>hist(x)
>hist(x,nclass=20)
```

For the different probability models, you need to use the **help** command to find out the precise syntax and parameter specification for each distribution.

10. TRANSFORMATIONS

For simulated data, generating a “transformed” sample is straightforward. If we wish to generate a sample from a continuous *Uniform*(0,1) distribution, and then to transform it using a $-\log$ transformation, we can proceed as follows:

```
>x1 <- runif(5000,0,1)
>hist(x1)
>x2 <- -log(x1)
>hist(x2)
```

EXERCISES:

1. Evaluate $P[X = 5]$ if $X \sim \text{Geometric}(0.6)$ (note: take care with the parameterization - for example check $P[X = 0]$ and compare this with the parameterization given in lecture notes; the **SPLUS** Geometric functions are based on the mass function

$$(1 - \theta)^x \theta \quad x = 0, 1, 2, \dots$$

which is a slightly different model to ours)

2. Evaluate $P[X = 15]$ if $X \sim \text{Poisson}(9)$
3. Evaluate $P[X \leq 12]$ if $X \sim \text{Binomial}(20, 0.6)$
4. Evaluate $P[X > 20]$ if $X \sim \text{Poisson}(15)$
5. Evaluate $P[30 < X \leq 45]$ if $X \sim \text{Binomial}(100, 0.35)$
6. Plot the pmf of the $\text{Poisson}(8)$ distribution on the range $0 \leq x \leq 20$
7. Plot the pdf of the $\text{Gamma}(5, 2)$ distribution on the range $0 \leq x \leq 10$
8. Plot the pdf of the $\text{Normal}(-5, 5^2)$ distribution on the range $-20 \leq x \leq 20$
9. Plot the cdf of the $\text{Normal}(0, 1)$ distribution on the range $-3 \leq x \leq 3$
10. Produce a sample of 5000 values from a $\text{Normal}(0, 1)$ distribution, plot a histogram, and then plot a histogram of the squares of these values.

It is also possible to generate a random sequence that is similar to a biological sequence using the **SPLUS** function **sample**: we proceed by issuing the following command:

```
>bases <- c('A','C','G','T')
>pvec <- c(0.25,0.25,0.25,0.25)
>x <- sample(bases,size=50,replace=T,prob=pvec)
>x
```

that will produce (something like) the following output

```
[1] 'C' 'A' 'T' 'G' 'A' 'G' 'C' 'C' 'A' 'A'
[11] 'G' 'G' 'C' 'T' 'C' 'T' 'C' 'C' 'C' 'C'
[21] 'G' 'C' 'T' 'A' 'A' 'T' 'C' 'G' 'T' 'G'
[31] 'C' 'A' 'C' 'A' 'A' 'G' 'A' 'T' 'C' 'A'
[41] 'C' 'T' 'T' 'G' 'A' 'G' 'C' 'G' 'C' 'T'
```

The commands created some base labels A , C , G and T , and then a probability for each label (in this case the probability is 0.25 for each label), and then produced a sample of size 50 independently sampled from this distribution.

The **prob** vector determines how probable each label is; in nature, it is unlikely that each base is observed with equal probability, and also that the base sequence is not independently generated (that is, the base observed at one position is influenced by bases observed in previous positions). In light of the analysis carried out for the BRCA2 sequence, how could a more realistic biological sequence be generated ?