

Removing batch effects in gene expression array study

Chung-Hsing Chen

National Cancer Institute, National Health
Research Institutes

Data background (1)

- Hedenfalk et al. measured **3226** genes in seven BRCA1 and eight BRCA2 mutation-positive tumor samples.
- The goal of the study was to identify genes that showed differential expression across breast cancer tumor subtypes defined by these germline mutations.
- Several genes with apparent outliers were removed. This left **3170** genes.

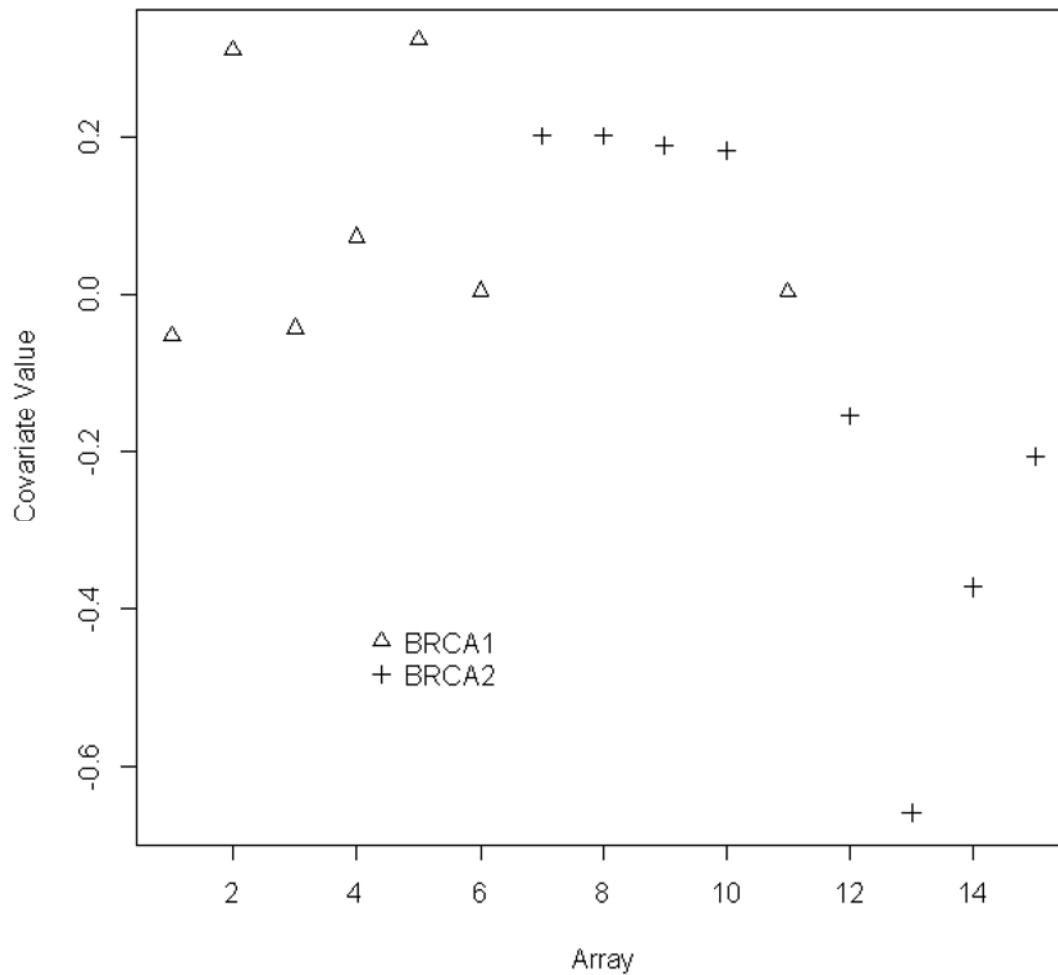
Data background (2)

```
> SVA.pheno
  NEJM-PatientID Mutation
V7          1    BRCA1
V8          5    BRCA1
V9          3    BRCA1
V10         7    BRCA1
V11         2    BRCA1
V12         4    BRCA1
V13        10    BRCA2
V14         9    BRCA2
V15         8    BRCA2
V16        10    BRCA2
V17         6    BRCA1
V18        13    BRCA2
V19        14    BRCA2
V20        11    BRCA2
V21        12    BRCA2
> head(SVA.exp.preprocessed)
   q-value    p-value fold-change (log base 2) s1996 s1822 s1714 s1224 s1252 s1510 s1900 s1787 s1721 s1486 s1905
[1,] 0.089529 0.01223344           1.203  0.15  0.22  0.30  0.26  1.22  0.44  0.35  1.10  1.07  1.46  0.38
[2,] 0.213752 0.07611987          -0.521  1.54  1.27  0.76  0.85  1.27  0.64  0.90  0.64  0.78  0.55  0.61
[3,] 0.672438 0.99530284           0.002  1.72  1.57  2.13  1.09  1.98  0.74  1.71  1.16  1.33  1.46  2.43
[4,] 0.163987 0.04212934           0.690  0.71  1.24  1.69  2.23  1.16  0.82  1.44  2.03  3.60  1.20  2.08
[5,] 0.637670 0.84745741           0.053  0.94  1.53  1.87  1.19  1.16  1.54  1.05  0.91  0.85  1.22  1.01
[6,] 0.377451 0.25437224          -0.227  0.80  0.95  1.53  1.37  1.02  1.22  0.78  0.96  0.65  1.02  1.09
   s1816 s1616 s1063 s1936
[1,] 0.73  0.63  0.77  0.66
[2,] 0.71  0.30  0.62  1.00
[3,] 1.71  1.26  1.41  3.00
[4,] 3.24  2.41  1.56  2.56
[5,] 3.25  2.20  1.09  1.29
[6,] 0.66  1.40  1.32  1.13
> head(exp.annotation)
  PlatePosition ImageCloneID Title
4      HK1A1     21652  catenin (cadherin-associated protein), alpha 1 (102kD)
5      HK1A2     22012  ADP-ribosylation factor 3
6      HK1A4     22293  uroporphyrinogen III synthase (congenital erythropoietic porphyria)
7      HK1A5     22493  ribosomal protein L26
8      HK1A6    23019  guanine nucleotide binding protein (G protein), alpha stimulating activity polypeptide 1
9      HK1A7     23132  pre-mRNA splicing factor SF3a (120 kDa subunit), similar to S. cerevisiae PRP21
```

Substructure detected

- Hierarchical clustering of the data reveals notable substructure within the BRCA2 samples.
- Applied **SVA** (Surrogate Variable Analysis) to identified a single surrogate variable that appears to capture this trend.

Substructure detected



Significance analysis

- Included this surrogate variable in a significance analysis comparing BRCA1 and BRCA2 tumors.
- The number of genes differentially expressed between BRCA1 and BRCA2 before and after adjusting for surrogate variables.

Analysis Type	q-Value Threshold			
	0.01	0.025	0.05	0.10
Unadjusted	1	19	96	275
SVA adjusted	0	10	48	190

R package: sva (1)

- Create the full model matrix - including both the adjustment variables and the variable of interest (BRCA1/BRCA2).
- The null model contains only the adjustment variables.

```
> mod=model.matrix(~as.factor(Mutation),data=SVA.pheno)
> mod0=model.matrix(~1,data=SVA.pheno)
> svaObj=sva(log2(SVA.exp.preprocessed[,-1:-3]),mod,mod0)
Number of significant surrogate variables is:  4
Iteration (out of 5 ):1  2  3  4  5  > |
```

R package: sva (2)

- The f.pvalue function can be used to calculate parametric F-test p-values for each gene. The F-test compares the models mod and mod0.

```
> pValues.before=f.pvalue(log2(SVA.exp.preprocessed[,-1:-3]),mod,mod0)
> qValuesObj.before<-qvalue(pValues.before)
> qsummary(qValuesObj.before)

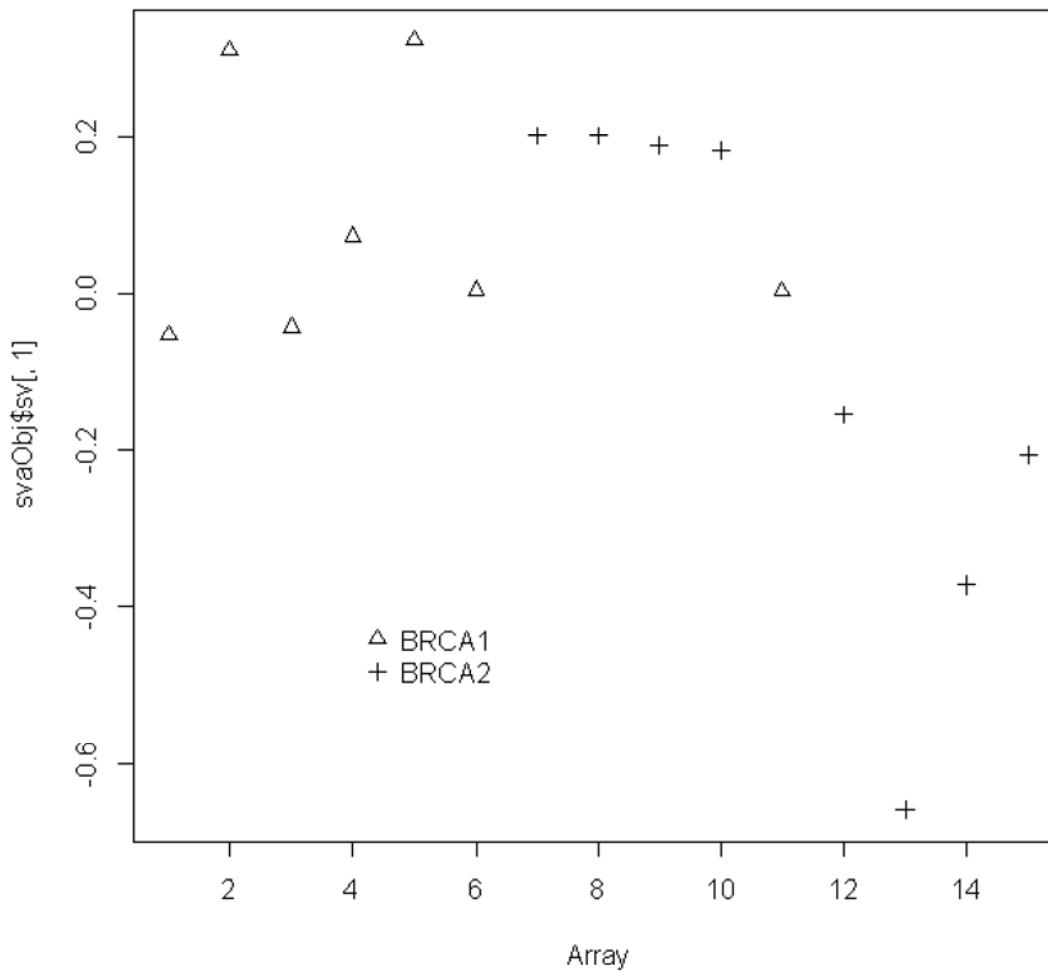
Call:
qvalue(p = pValues.before)

pi0:    0.6786508

Cumulative number of significant calls:

      <1e-04 <0.001 <0.01 <0.025 <0.05 <0.1   <1
p-value     9       62     228     392     565     832 3170
q-value     0       0       1      19      96    275 3170
```

R package: sva (3)



R package: sva (4)

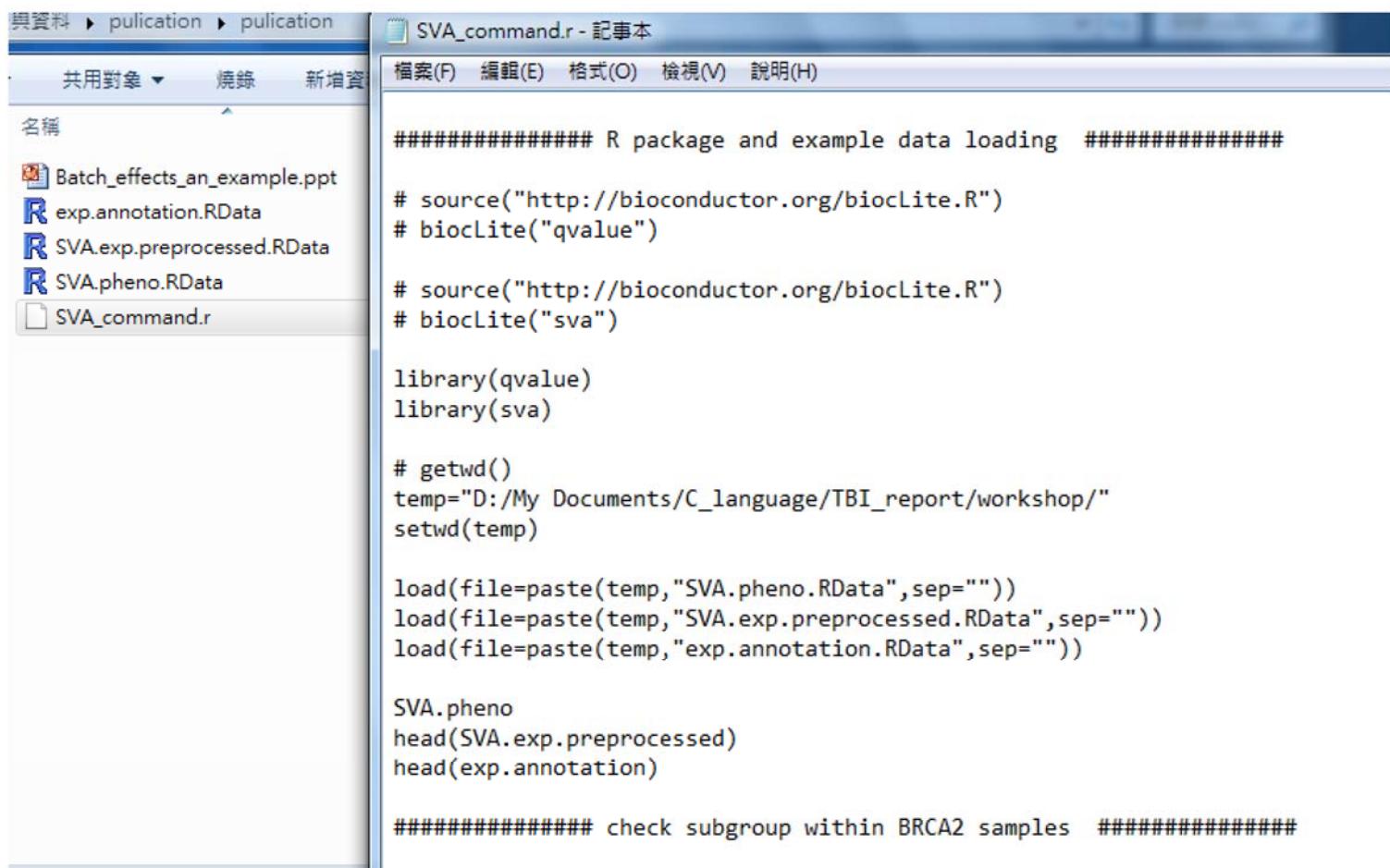
- Now we can perform the same analysis, but adjusting for surrogate variables. The first step is to include the surrogate variables in both the null and full models.

```
> svaObj=sva(log2(SVA.exp.preprocessed[,-1:-3]),mod,mod0)
Number of significant surrogate variables is: 4
Iteration (out of 5 ):1 2 3 4 5 >
> modSv=cbind(mod,svaObj$sv[,1])
> mod0Sv=cbind(mod0,svaObj$sv[,1])
```

- Then P-values and Q-values can be computed as before.

Step by Step: Open the text file

Use the Notepad to open the file named “SVA_command.r” (text file).



The screenshot shows a Windows Notepad window titled "SVA_command.r - 記事本". The left pane shows a file list with several RData files and a command file. The right pane displays the content of the "SVA_command.r" file, which contains R code for loading packages, setting the working directory, and loading data files.

```
##### R package and example data loading #####
# source("http://bioconductor.org/biocLite.R")
# biocLite("qvalue")

# source("http://bioconductor.org/biocLite.R")
# biocLite("sva")

library(qvalue)
library(sva)

# getwd()
temp="D:/My Documents/C_language/TBI_report/workshop/"
setwd(temp)

load(file=paste(temp,"SVA.pheno.RData",sep=""))
load(file=paste(temp,"SVA.exp.preprocessed.RData",sep=""))
load(file=paste(temp,"exp.annotation.RData",sep=""))

SVA.pheno
head(SVA.exp.preprocessed)
head(exp.annotation)

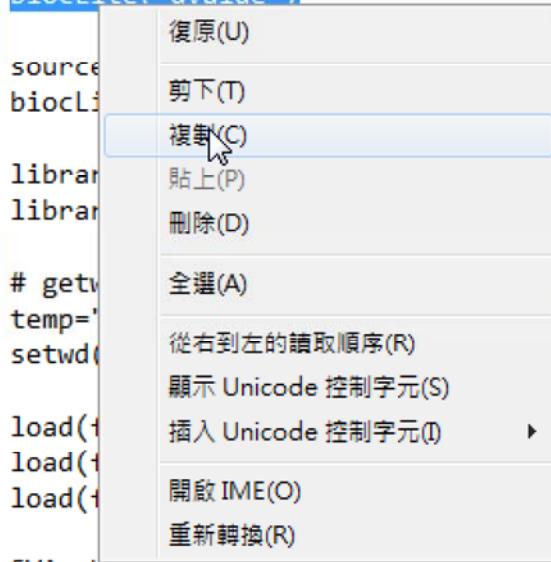
##### check subgroup within BRCA2 samples #####

```

Step by Step: Install R packages (1)

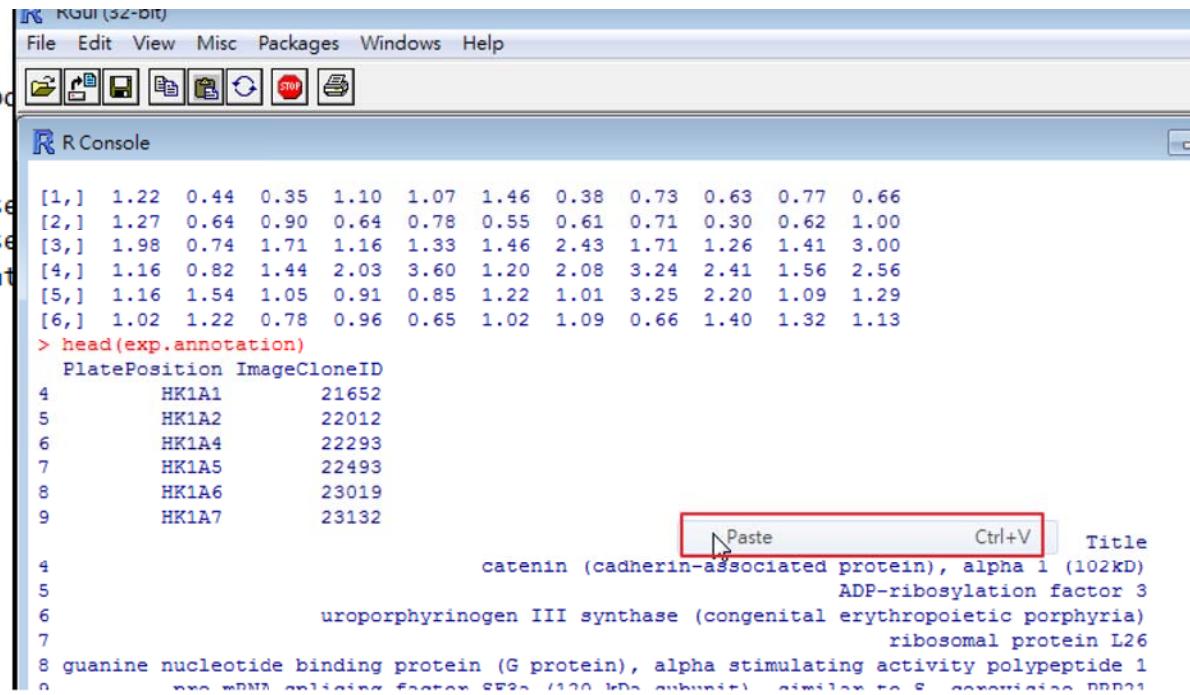
```
##### R package and example data loading #####
```

```
source("http://bioconductor.org/biocLite.R")
biocLite("qvalue")
```



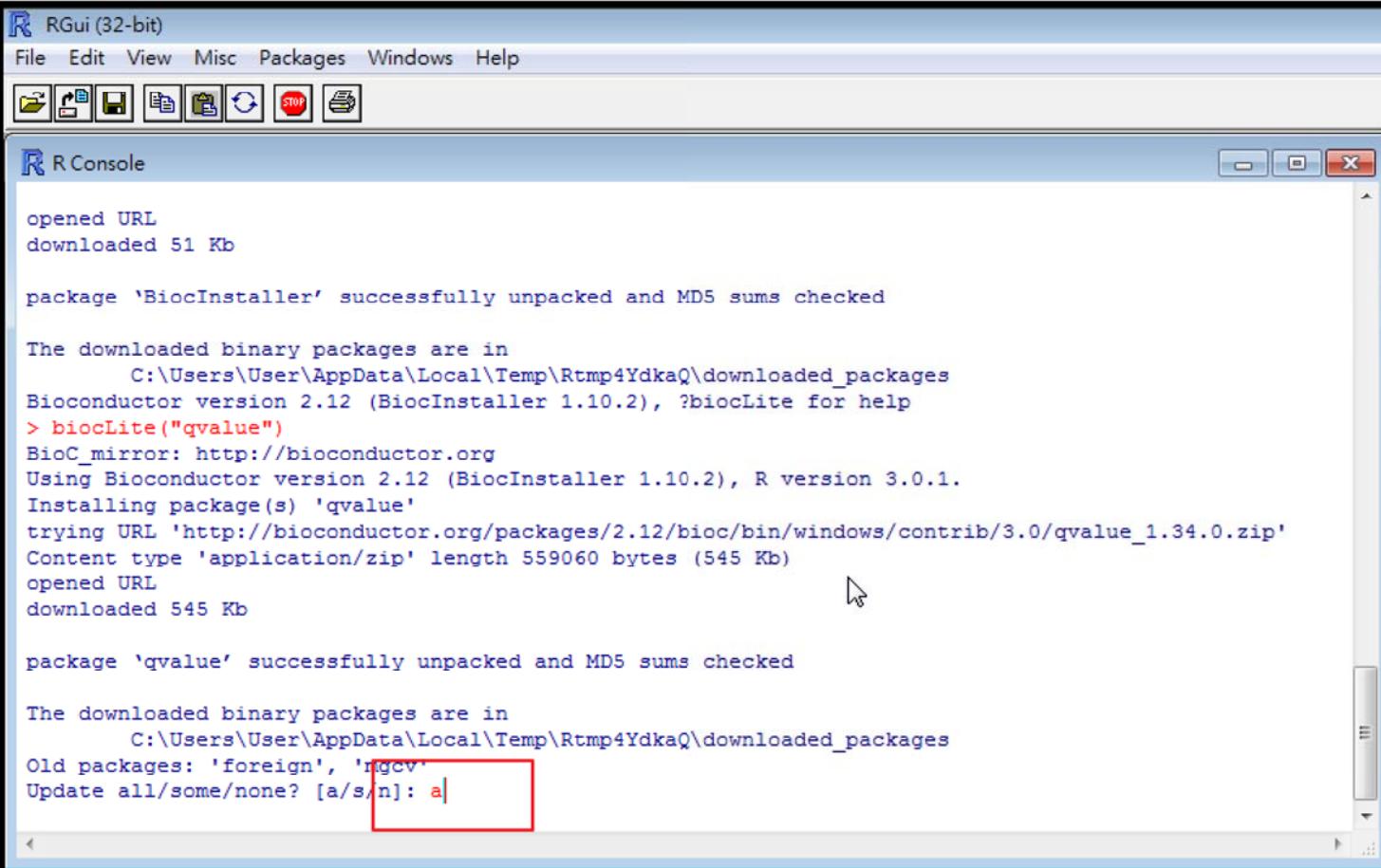
```
source
biocLi
library
library
# getw
temp='
setwd(
load(
load(
load(
SVA.pheno
head(SVA.exp.preprocessed)
head(exp.annotation)
```

1. Copy the command lines for installing the R package “qvalue”.



2. Paste to the R window.

Step by Step: Install R packages (2)



The screenshot shows the RGui (32-bit) application window. The title bar reads "RGui (32-bit)". The menu bar includes "File", "Edit", "View", "Misc", "Packages", "Windows", and "Help". Below the menu is a toolbar with various icons. The main window is titled "R Console". The console output is as follows:

```
opened URL
downloaded 51 Kb

package 'BiocInstaller' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
  C:\Users\User\AppData\Local\Temp\Rtmp4YdkaQ\downloaded_packages
Bioconductor version 2.12 (BiocInstaller 1.10.2), ?biocLite for help
> biocLite("qvalue")
BioC_mirror: http://bioconductor.org
Using Bioconductor version 2.12 (BiocInstaller 1.10.2), R version 3.0.1.
Installing package(s) 'qvalue'
trying URL 'http://bioconductor.org/packages/2.12/bioc/bin/windows/contrib/3.0/qvalue_1.34.0.zip'
Content type 'application/zip' length 559060 bytes (545 Kb)
opened URL
downloaded 545 Kb

package 'qvalue' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
  C:\Users\User\AppData\Local\Temp\Rtmp4YdkaQ\downloaded_packages
Old packages: 'foreign', 'ngcv'
Update all/some/none? [a/s/n]: a
```

A red box highlights the command "a" at the bottom of the console window.

3. Type “a” for updating all old packages.

Step by Step: Install R packages (3)

4. Repeat step 1-3. The same way to install the package "sva".

The screenshot shows a Windows context menu (right-clicked) over a portion of an R script in a Notepad window titled "SVA_command.r - 記事本". The menu items are: 復原(U), 剪下(T), 複製(C), 貼上(P), 刪除(D), 全選(A), 從右到左的讀取順序(R), 顯示 Unicode 控制字元(S), 插入 Unicode 控制字元(I), 開啟 IME(O), and 重新轉換(R). The "複製(C)" item is highlighted with a blue selection bar.

```
##### R package and example data loading #####
source("http://bioconductor.org/biocLite.R")
biocLite("qvalue")

source("http://bioconductor.org/biocLite.R")
biocLite("sva")
library(qvalue)
library(sva)

# getwd()
temp="D:/My
setwd(temp)

load(file=pa
load(file=pa
load(file=pa

SVA.pheno
head(SVA.ex
head(exp.anr
```

Step by Step : Load R packages

The screenshot shows the RGui interface. On the left is the R Console window, which contains R code for loading packages and setting working directory. A red box highlights the command `library(qvalue)`. On the right is a script editor window titled "SVA_command.r - 記事本" containing the same R code. A red box highlights the command `library(sva)`.

```
> source("http://biocLite("sva")
BioC_mirror: http://
Using Bioconduct
Installing packa
also installing

trying URL 'http://
Content type 'ap
opened URL
downloaded 81 Kb

trying URL 'http://
Content type 'ap
opened URL
downloaded 200 Kb

package 'corpcor'
package 'sva' si
The downloaded b
C:\Users\

> | SVA.pheno
head(SVA.exp.preprocessed)
head(exp.annotation)

# getwd()
temp="D:/My Documents/C_language/T
setwd(temp) 複製(C)

load(file=paste(temp,"SVA.pheno.R"))
load(file=paste(temp,"SVA.exp.prepro
load(file=paste(temp,"exp.annotation

SVA.pheno
head(SVA.exp.preprocessed)
head(exp.annotation)
```

1. Copy the command lines for loading two R packages, “qvalue” and “sva”.

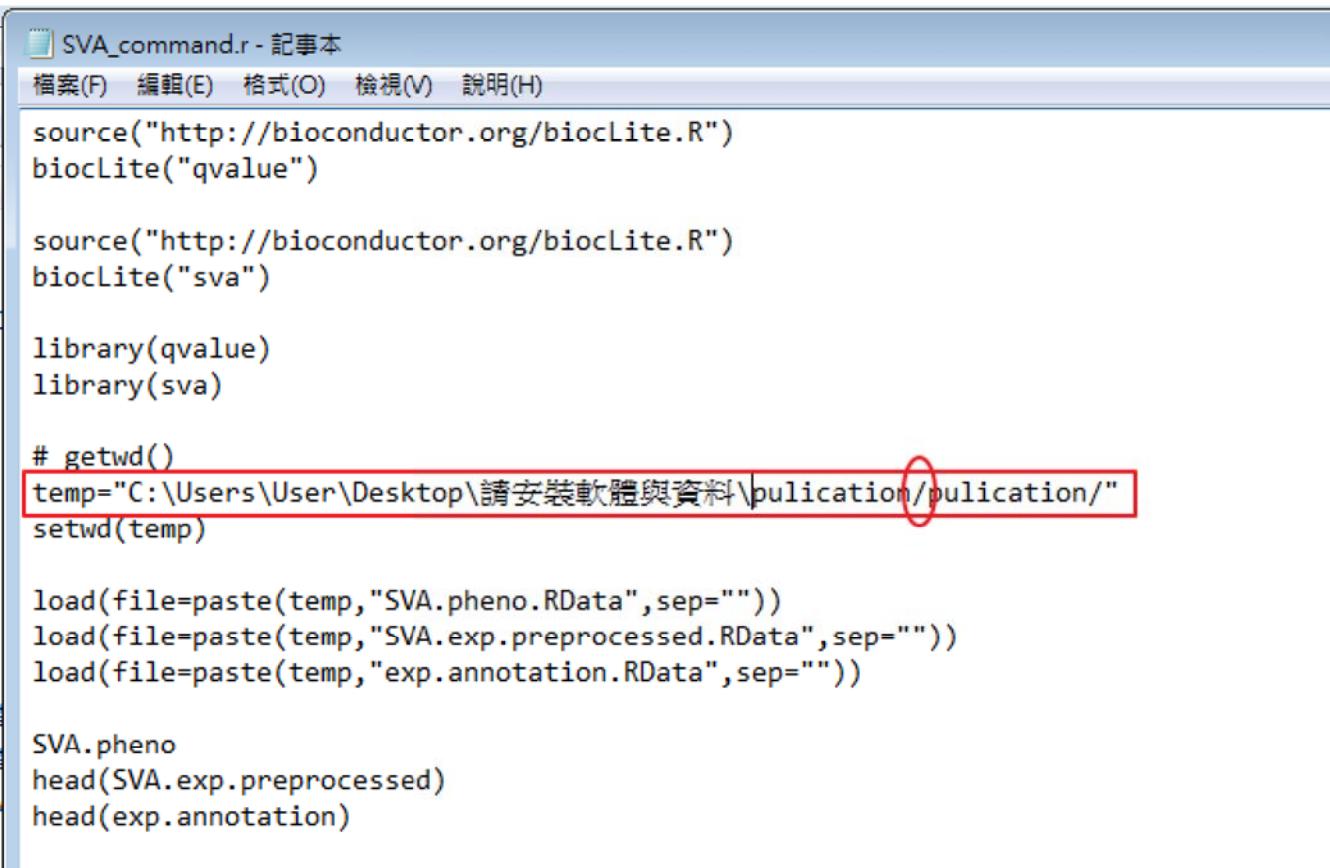
The screenshot shows the RGui interface with the R window active. A red box highlights the "Paste" option in the context menu, which is open over the R window. The menu also includes "Ctrl+V" and "Title".

```
[1] 1.22 0.44 0.35 1.10 1.07 1.46 0.38 0.73 0.63 0.77 0.66
[2] 1.27 0.64 0.90 0.64 0.78 0.55 0.61 0.71 0.30 0.62 1.00
[3] 1.98 0.74 1.71 1.16 1.33 1.46 2.43 1.71 1.26 1.41 3.00
[4] 1.16 0.82 1.44 2.03 3.60 1.20 2.08 3.24 2.41 1.56 2.56
[5] 1.16 1.54 1.05 0.91 0.85 1.22 1.01 3.25 2.20 1.09 1.29
[6] 1.02 1.22 0.78 0.96 0.65 1.02 1.09 0.66 1.40 1.32 1.13
> head(exp.annotation)
   PlatePosition ImageCloneID
4      HK1A1      21652
5      HK1A2      22012
6      HK1A4      22293
7      HK1A5      22493
8      HK1A6      23019
9      HK1A7      23132
4
catenin (cadherin-associated protein), alpha 1 (102kD)
5
ADP-ribosylation factor 3
6
uroporphyrinogen III synthase (congenital erythropoietic porphyria)
7
ribosomal protein L26
8
guanine nucleotide binding protein (G protein), alpha stimulating activity polypeptide 1
9
DNA (cytosine-5-ribonucleoside diphosphate ribose) polymerase similar to C-nucleoside polymerase
```

2. Paste to the R window.

Step by Step : R environment

1. Change the current working directory for loading the example data files.
2. Replace the symbol "\ to "/.



```
SVA_command.r - 記事本
檔案(F) 編輯(E) 格式(O) 檢視(V) 說明(H)

source("http://bioconductor.org/biocLite.R")
biocLite("qvalue")

source("http://bioconductor.org/biocLite.R")
biocLite("sva")

library(qvalue)
library(sva)

# getwd()
temp="C:\Users\User\Desktop\請安裝軟體與資料\pulation/pulation/"
setwd(temp)

load(file=paste(temp,"SVA.pheno.RData",sep=""))
load(file=paste(temp,"SVA.exp.preprocessed.RData",sep=""))
load(file=paste(temp,"exp.annotation.RData",sep=""))

SVA.pheno
head(SVA.exp.preprocessed)
head(exp.annotation)
```

Step by Step : Load example files

The screenshot shows the R GUI interface. On the left, there is a script editor window titled "SVA_command.r - 記事本" containing R code. On the right, the "R Console" window displays the output of the R code. A context menu is open over the script editor, with the "Copy" option highlighted.

1. Copy the command lines for loading example files.

```
source("http://bioconductor.org/biocLite.R")
biocLite("qvalue")

source("http://bioconductor.org/biocLite.R")
biocLite("sva")

library(qvalue)
library(sva)

# getwd()
temp="C:/Users/User/Desktop/請安裝軟體"
setwd(temp)

load(file= paste(temp, "SVA.pheno.RData"))
load(file= paste(temp, "SVA.exp.RData"))
```

2. Paste to the R window.

```
[1]  1.22  0.44  0.35  1.10  1.07  1.46  0.38  0.73  0.63  0.77  0.66
[2]  1.27  0.64  0.90  0.64  0.78  0.55  0.61  0.71  0.30  0.62  1.00
[3]  1.98  0.74  1.71  1.16  1.33  1.46  2.43  1.71  1.26  1.41  3.00
[4]  1.16  0.82  1.44  2.03  3.60  1.20  2.08  3.24  2.41  1.56  2.56
[5]  1.16  1.54  1.05  0.91  0.85  1.22  1.01  3.25  2.20  1.09  1.29
[6]  1.02  1.22  0.78  0.96  0.65  1.02  1.09  0.66  1.40  1.32  1.13

> head(exp.annotation)
  PlatePosition ImageCloneID
4       HK1A1      21652
5       HK1A2      22012
6       HK1A4      22293
7       HK1A5      22493
8       HK1A6      23019
9       HK1A7      23132
```

Paste Ctrl+V Title
catenin (cadherin-associated protein), alpha 1 (102kD)
ADP-ribosylation factor 3
uroporphyrinogen III synthase (congenital erythropoietic porphyria)
ribosomal protein L26
guanine nucleotide binding protein (G protein), alpha stimulating activity polypeptide 1
DNA (cytosine-5-ribonucleotidyl transferase) catalytic subunit similar to C-ribonucleotid

Step by Step : Command lines

The package “sva” command lines.

```
mod=model.matrix(~as.factor(Mutation),data=SVA.pheno)
mod0=model.matrix(~1,data=SVA.pheno)
svaObj=sva(log2(SVA.exp.preprocessed[,-1:-3]),mod,mod0)
```

The command for plotting the graph.

```
symbol<-rep(2,15)
symbol[BRCA2]<-3
svaObj$sv

plot(1:15,svaObj$sv[,1],xlab="Array",ylab="Covariate Value",pch=symbol)
legend(4,-0.4,c("BRCA1","BRCA2"),pch=c(2,3),bty = "n")
```

The command for computing F-test and q-value.

```
##### compare q-value before/after batch effect #####
pValues.before=f.pvalue(log2(SVA.exp.preprocessed[,-1:-3]),mod,mod0)
qValuesObj.before<-qvalue(pValues.before)

qsummary(qValuesObj.before)

modSv=cbind(mod,svaObj$sv[,1])
modOSv=cbind(mod0,svaObj$sv[,1])

pValues.after=f.pvalue(log2(SVA.exp.preprocessed[,-1:-3]),modSv,modOSv)
qValuesObj.after<-qvalue(pValues.after)

qsummary(qValuesObj.after)
```

Appendix: print screen (1)

R Gui

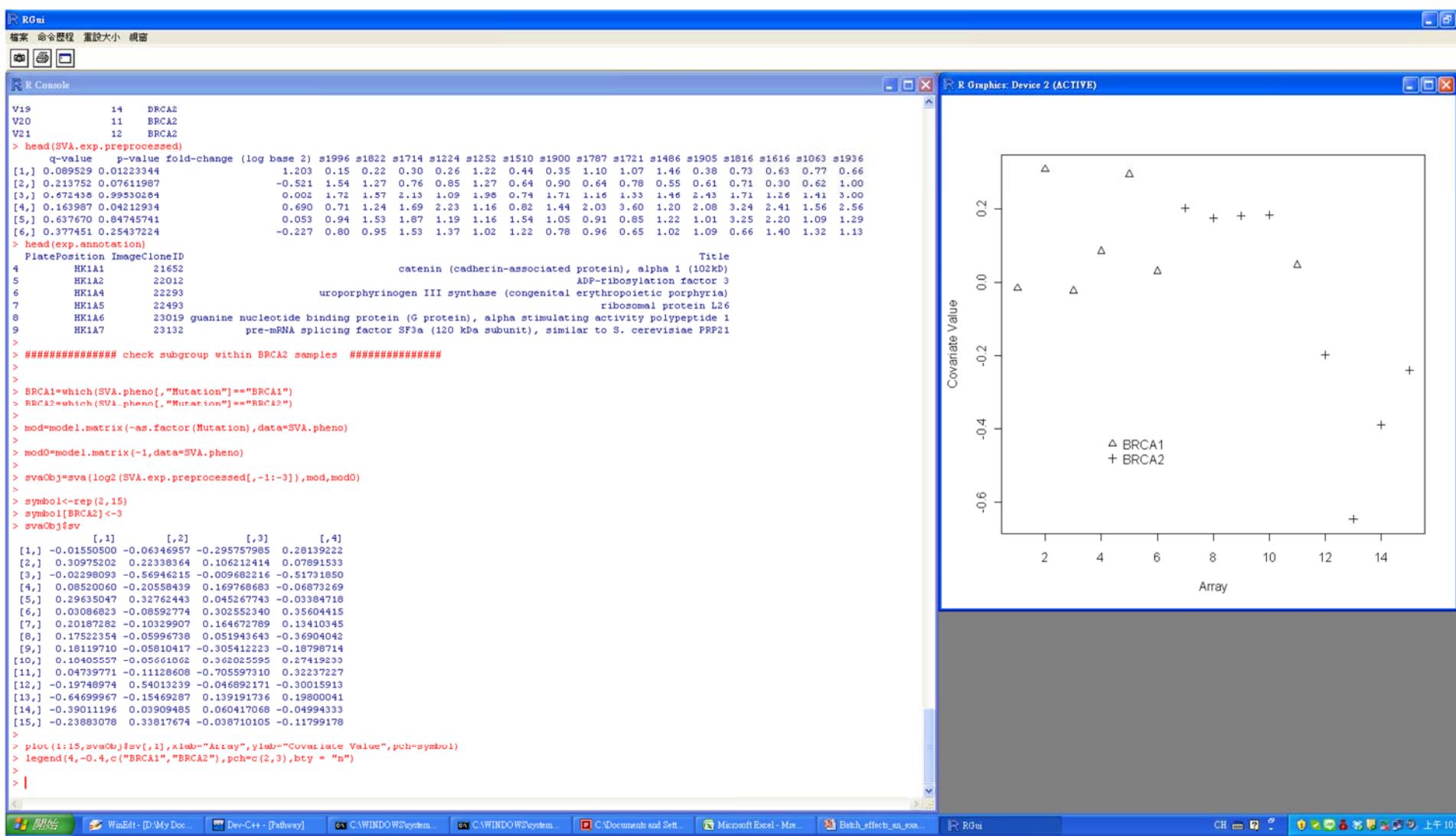
檔案 編輯 看 其它 程式套件 調查 幫助

R Console

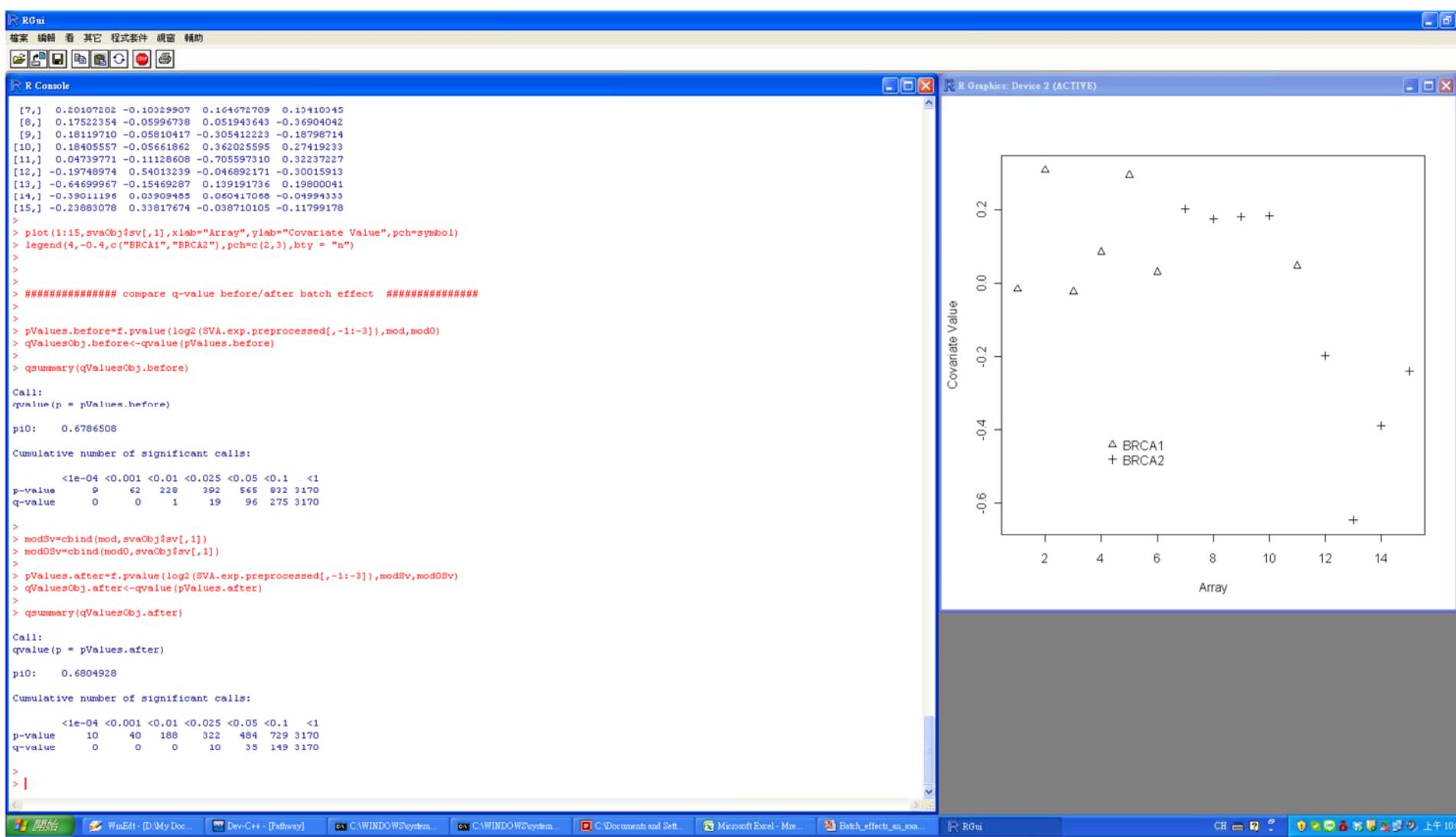
```
>
>
>
>
>
>
>
> library(qvalue)
> library(sva)
>
> # getwd()
> temp="D:/My Documents/C_language/TBI_report/workshop/"
> setwd(temp)
>
> load(file=paste(temp,"SVA.pheno.RData",sep=""))
> load(file=paste(temp,"SVA.exp.preprocessed.RData",sep=""))
> load(file=paste(temp,"exp.annotation.RData",sep=""))
>
> SVA.pheno
NEJM-PatientID Mutation
V7      1    BRCA1
V8      5    BRCA1
V9      3    BRCA1
V10     7    BRCA1
V11     2    BRCA1
V12     4    BRCA1
V13    10    BRCA2
V14     9    BRCA2
V15     8    BRCA2
V16    10    BRCA2
V17     6    BRCA1
V18    13    BRCA2
V19    14    BRCA2
V20    11    BRCA2
V21    12    BRCA2
> head(SVA.exp.preprocessed)
   q-value    p-value fold-change (log base 2) s1996 s1822 s1714 s1224 s1252 s1510 s1900 s1787 s1721 s1486 s1905 s1616 s1616 s1063 s1936
[1,] 0.089529 0.01223344          1.203 0.15 0.22 0.30 0.26 1.22 0.44 0.35 1.10 1.07 1.46 0.38 0.73 0.63 0.77 0.66
[2,] 0.213752 0.07611987         -0.521 1.54 1.27 0.76 0.85 1.27 0.64 0.90 0.64 0.78 0.55 0.61 0.71 0.30 0.62 1.00
[3,] 0.672438 0.99530284          0.002 1.72 1.57 2.13 1.09 1.98 0.74 1.71 1.16 1.33 1.46 2.43 1.71 1.26 1.41 3.00
[4,] 0.163987 0.04212934          0.690 0.71 1.24 1.69 2.23 1.16 0.82 1.44 2.03 3.60 1.20 2.08 3.24 2.41 1.56 2.56
[5,] 0.637670 0.84745741          0.053 0.94 1.53 1.87 1.19 1.16 1.54 1.05 0.91 0.85 1.22 1.01 3.25 2.20 1.09 1.29
[6,] 0.377451 0.25437224         -0.227 0.80 0.95 1.53 1.37 1.02 1.22 0.78 0.96 0.65 1.02 1.09 0.66 1.40 1.32 1.13
> head(exp.annotation)
  PlatePosition ImageCloneID           Title
 4      HK1A1    21652   catenin (cadherin-associated protein), alpha 1 (102kD)
 5      HK1A2    22012           ADP-ribosylation factor 3
 6      HK1A4    22293 uroporphyrinogen III synthase (congenital erythropoietic porphyria)
 7      HK1A5    22493           ribosomal protein L26
 8      HK1A6  23019 guanine nucleotide binding protein (G protein), alpha stimulating activity polypeptide 1
 9      HK1A7    23132 pre-mRNA splicing factor SF3a (120 kDa subunit), similar to S. cerevisiae PRP21
>
>
>
```

開始 WmEdit - [D:\My Doc... Dev-C++ - [Pathway] CAWINDOWSystem CAWINNDSystem CADocuments and Set... Microsoft Excel - Min... Batch_effect_an_ex... R Gui CH 上午 09:49

Appendix: print screen (2)



Appendix: print screen (3)



Appendix: print screen (4)

