



Rによるトランスクリプトーム解析

～NGS由来塩基配列データを自在に解析する～

東京大学大学院農学生命科学研究科
アグリバイオインフォマティクス教育研究ユニット

門田 幸二(かどた こうじ)

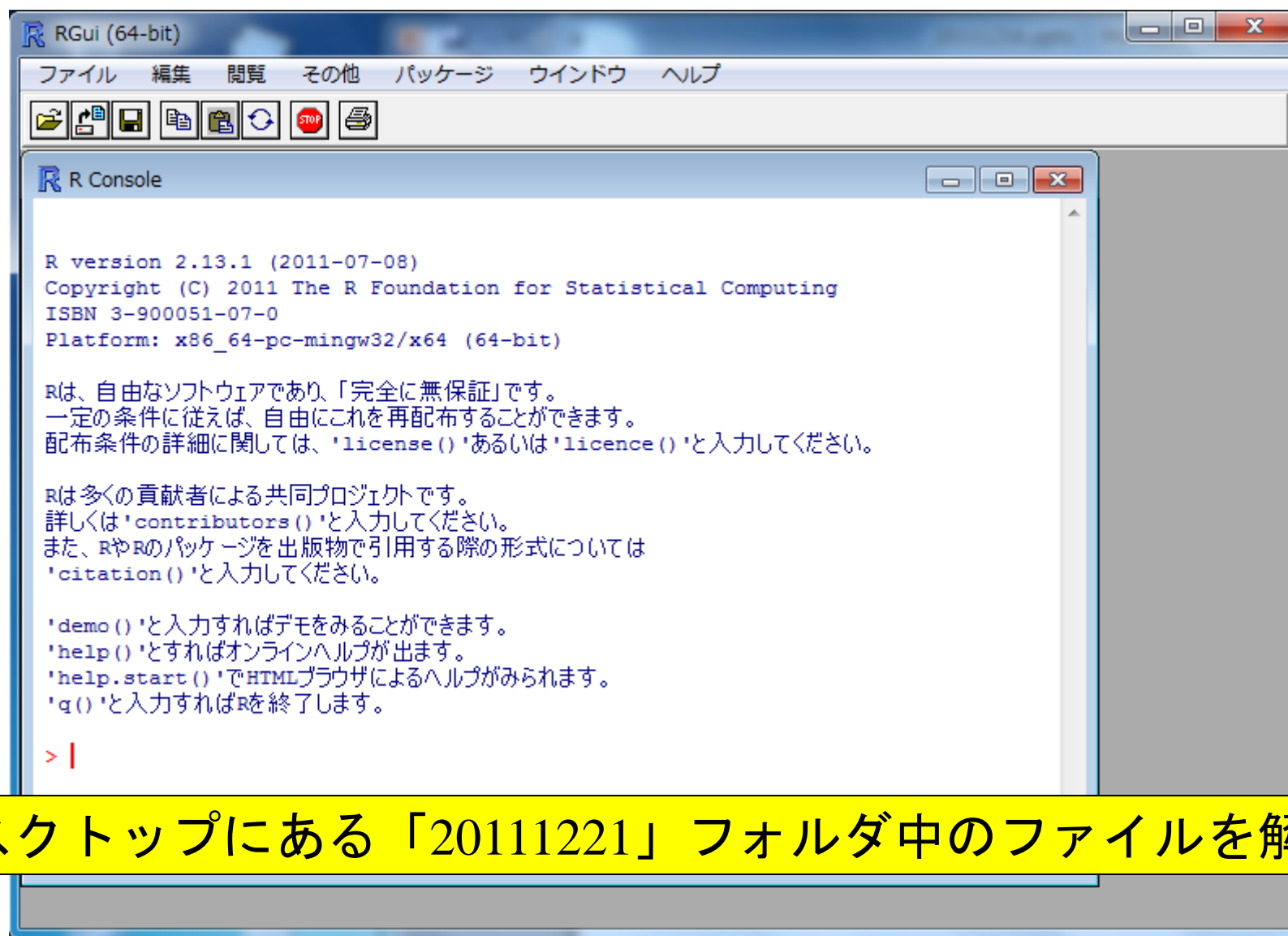
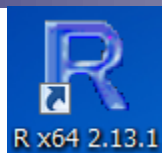
<http://www.iu.a.u-tokyo.ac.jp/~kadota/>

kadota@iu.a.u-tokyo.ac.jp

Contents

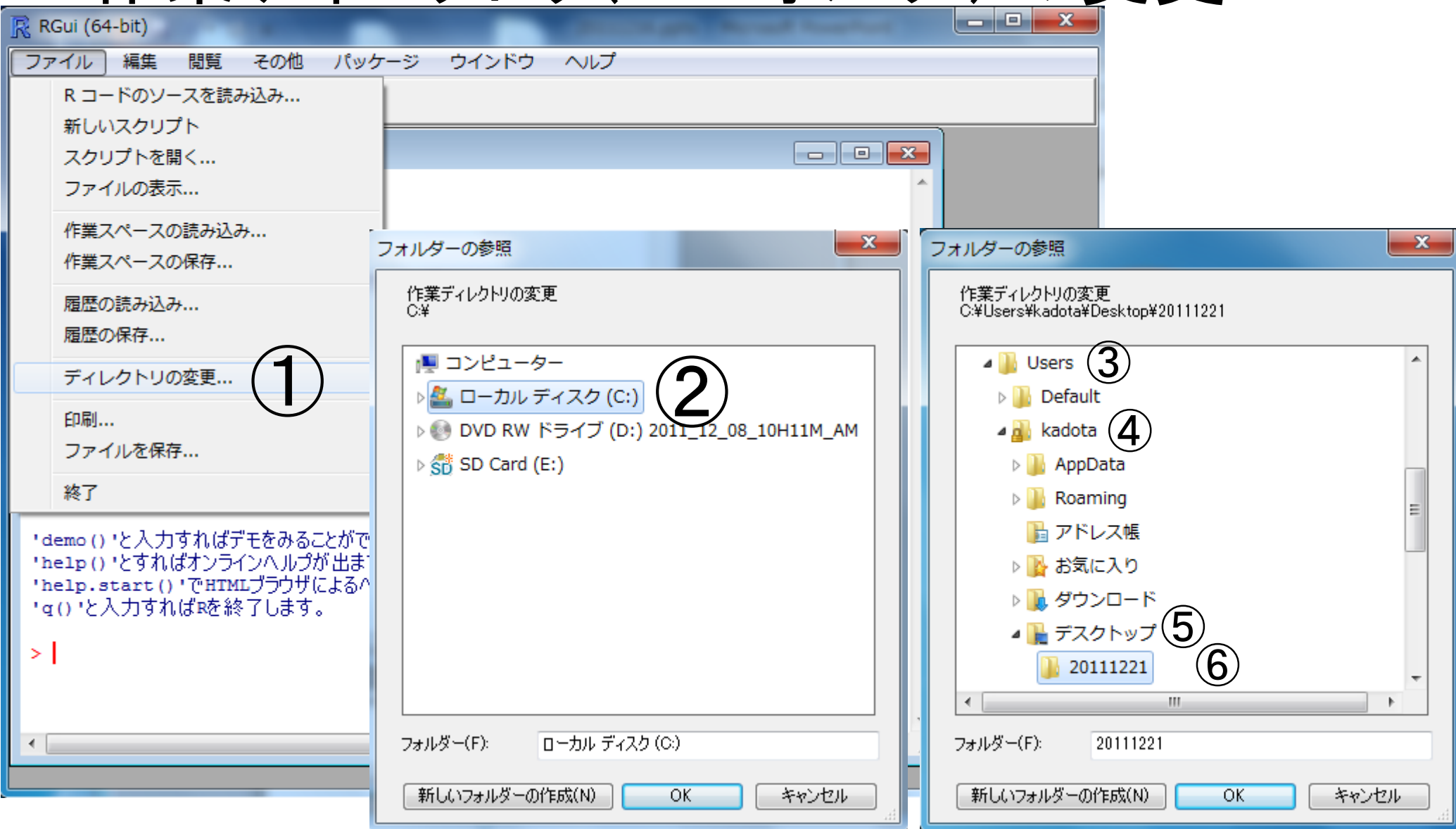
- multi-fasta形式ファイルからの情報抽出
- 比較トランスクリプトーム解析(二群間比較)
 - 各種Rパッケージ
 - 分布(ポアソン分布、負の二項分布)
 - *edgeR*パッケージを使ってみる(technical replicates)
 - MA-plot
 - 倍率変化がだめな理由をデモ
 - *edgeR*パッケージを使ってみる(biological replicates)
 - NGSデータのクオリティチェックを行う
 - 任意のRパッケージのインストールのやり方
 - サンプル間クラスタリングは基本中の基本

Rの起動

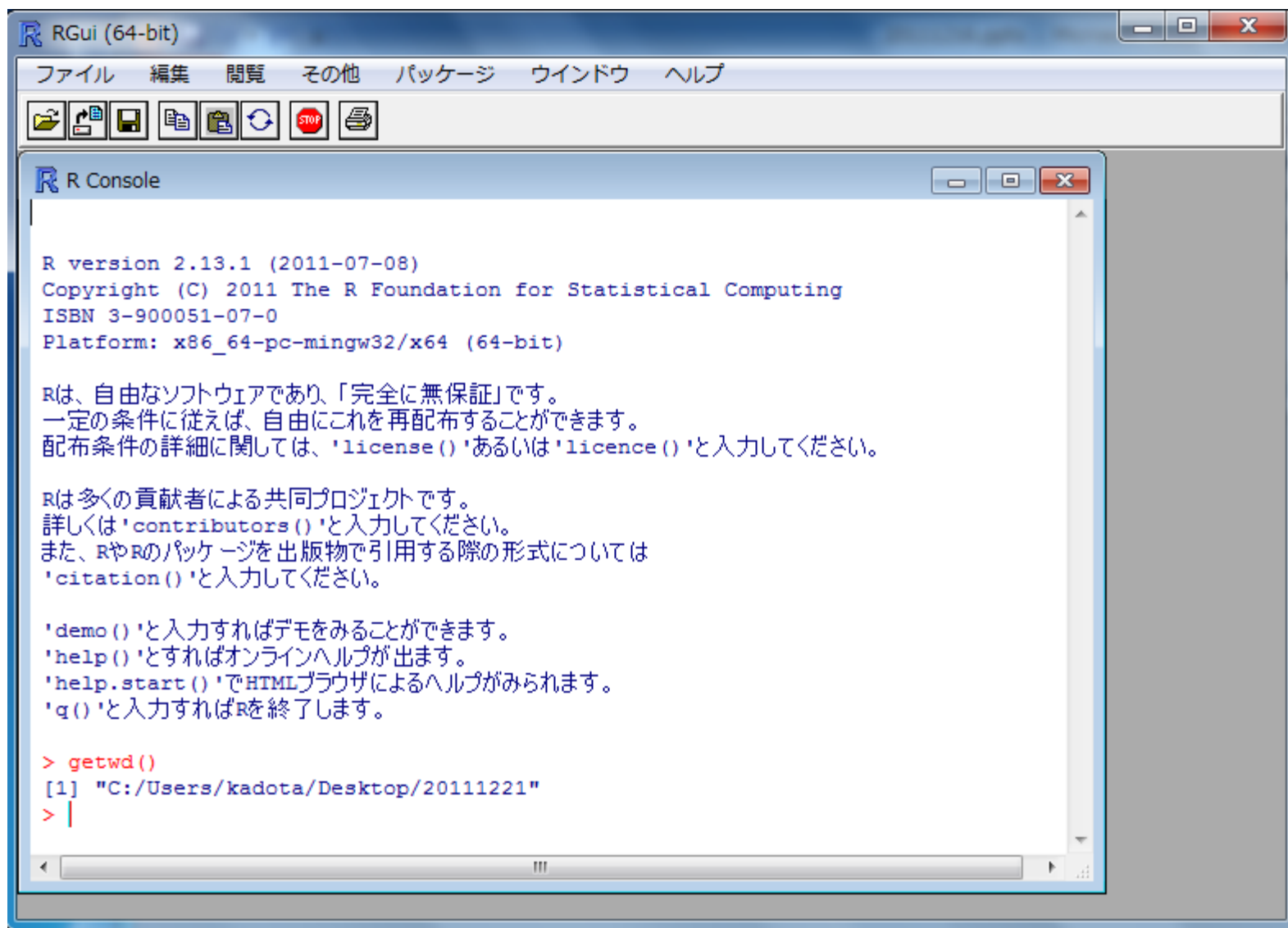


デスクトップにある「20111221」フォルダ中のファイルを解析

作業ディレクトリ(=フォルダ)の変更



「getwd()」と打ち込んで確認



The screenshot shows the RGui (64-bit) window. The title bar says 'RGui (64-bit)'. The menu bar includes 'ファイル', '編集', '閲覧', 'その他', 'パッケージ', 'ウインドウ', and 'ヘルプ'. The toolbar contains icons for file operations and execution. The R Console window is open, displaying the following text:

```
R version 2.13.1 (2011-07-08)
Copyright (C) 2011 The R Foundation for Statistical Computing
ISBN 3-900051-07-0
Platform: x86_64-pc-mingw32/x64 (64-bit)

Rは、自由なソフトウェアであり、「完全に無保証」です。
一定の条件に従えば、自由にこれを再配布することができます。
配布条件の詳細に関しては、'license()'あるいは'licence()'と入力してください。

Rは多くの貢献者による共同プロジェクトです。
詳しくは'contributors()'と入力してください。
また、RやRのパッケージを出版物で引用する際の形式については
'citation()'と入力してください。

'demo()'と入力すればデモをみることができます。
'help()'とすればオンラインヘルプが出ます。
'help.start()'でHTMLブラウザによるヘルプがみられます。
'q()'と入力すればRを終了します。

> getwd()
[1] "C:/Users/kadota/Desktop/20111221"
> |
```

参考ウェブページ

(Rで)塩基配列解析(主に次世代シーケンサーのデータ) by 門田幸二 (last modified 2012/02/06)

What's new?

・2/22の横浜理研の[Rでつなぐ次世代オミックス情報統合解析研究会](#)では(しゃべり時間50分しかないので)デモのみ行います。また、3/9のお台場での[HPCI チュートリアルセミナー](#)では(4時間ほどあるので)実際に手を動かしてもらいます。また、時間があれば私の最新の手法の使用法も伝授できればと思っています。興味ある方はどうぞ。(2012/02/03) **NEW**

・最新のパッケージ中の使用法の記述への更新が相当遅れていますので、ご利用時はもとのパッケージ中のマニュアルを各自チェックしてください。(2011/12/26)

・R2.14.1がリリースされていたのでこれに変更しました。(2011/12/26)

- [はじめに](#) (last modified 2011/07/19)
- [Rのインストールと起動](#) (last modified 2012/01/04)
- [サンプルデータ](#) (last modified 2011/02/03)
- イントロダクション | NGS | [各種覚書](#) (last modified 2010/12/10)
- イントロダクション | NGS | [様々なプラットフォーム](#) (last modified 2011/07/15)
- イントロダクション | NGS | [リファレンス配列取得\(マップされる側\)](#) (last modified 2011/02/03)
- イントロダクション | NGS | [リファレンス配列取得後の各種情報抽出\(特にRefSeq\)](#) (last modified 2011/03/20)
- イントロダクション | NGS | [リファレンス配列取得後の各種情報抽出2\(readFASTA関数の利用\)](#) (last modified 2011/04/07)
- イントロダクション | NGS | [アノテーション情報取得\(refFlatファイル\)](#) (last modified 2010/12/07)
- イントロダクション | NGS | [アノテーション情報取得\(BioMart and biomaRt\)](#) (last modified 2011/08/26)
- イントロダクション | 一般 | [配列取得](#) (last modified 2010/7/7)
- イントロダクション | 一般 | [指定した範囲の配列を取得](#) (last modified 2012/01/05)
- イントロダクション | 一般 | [翻訳配列\(translate\)を取得](#) (last modified 2011/07/27)
- イントロダクション | 一般 | [相補鎖\(complement\)を取得](#) (last modified 2011/07/27)
- イントロダクション | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2011/07/27)
- イントロダクション | 一般 | [逆鎖\(reverse\)を取得](#) (last modified 2011/07/27)
- イントロダクション | 一般 | [二連続塩基の出現頻度情報を取得](#) (last modified 2011/07/25)
- イントロダクション | 一般 | [三連続塩基の出現頻度情報を取得](#) (last modified 2011/07/25)
- イントロダクション | 一般 | [multi-fastaファイルから指定した配列長のもののみ抽出](#) (last modified 2012/02/06) **NEW**
- イントロダクション | NGS | [NGSデータ取得](#) (last modified 2011/07/19)

multi-fasta形式ファイルからの情報抽出1

• イントロダクション | NGS | アセンブル後のmulti-fastaファイルからN50などの基本情報を取得

Trinityなどのアセンブルプログラムを実行したあとのファイルからを想定して、Total lengthやaverage lengthなどを示します。ここでは130MB程度のmulti-fastaファイル「ファイル」-「ディレクトリの変更」でファイル

```
----- ここから -----
in_f <- "h_rna.fasta"
out_f <- "hoge.txt"

library(Biostrings)
reads <- read.DNAStringSet(in_f, format="fasta")

Total_length <- sum(width(reads))
Number_of_contigs <- length(reads)
Average_length <- mean(width(reads))
Median_length <- median(width(reads))
Max_length <- max(width(reads))
Min_length <- min(width(reads))

#N50計算のところ
sorted_length <- rev(sort(width(reads)))
N50 <- sorted_length[sum(sorted_length)]

#GC含量(GC content)計算のところ
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
GC_content <- sum(CG)/sum(ACGT)

#出力用に結果をまとめている
tmp <- NULL
tmp <- rbind(tmp, c("Total length (bp)", Total_length))
tmp <- rbind(tmp, c("Number of contigs", Number_of_contigs))
tmp <- rbind(tmp, c("Average length", Average_length))
tmp <- rbind(tmp, c("Median length", Median_length))
tmp <- rbind(tmp, c("Max length", Max_length))
tmp <- rbind(tmp, c("Min length", Min_length))
tmp <- rbind(tmp, c("N50", N50))
tmp <- rbind(tmp, c("GC content", GC_content))
write.table(tmp, out_f, sep="\t", append=TRUE)
```

RGui (64-bit)

ファイル 編集 閲覧 その他 パッケージ ウィンドウ ヘルプ

R Console

R version 2.13.1 (2011-07-08)
Copyright (C) 2011 The R Foundation for Statistical Computing
ISBN 3-900051-07-0
Platform: x86_64-pc-mingw32/x64 (64-bit)

Rは、自由なソフトウェアであり「完全に無保証」である。特定の条件下で、自由にこれを再配布する。配布条件の詳細に関しては、'license()' または 'help(license)' を入力してください。

Rは多くの貢献者による共同プロジェクトです。詳しくは 'contributors()' を入力してください。また、RやRのパッケージを出版物で引用する際は 'citation()' を入力してください。

'demo()' と入力すればデモをみることができます。'help()' とすればオンラインヘルプが出ます。'help.start()' でHTMLブラウザによるヘルプを見られます。'q()' と入力すればRを終了します。

> getwd()
[1] "C:/Users/kadota/Desktop/"
> |

①一連のコマンド群をコピーして
②R Console画面上でペースト

multi-fasta形式ファイルからの情報抽出1

```

RGui (64-bit)
ファイル 編集 閲覧 その他 パッケージ ウィンドウ ヘルプ

R Console
> #N50計算のところ
> sorted_length <- rev(sort(width(reads)))
> N50 <- sorted_length[cumsum(sorted_length) >= Total_length/2][1]
>
> #GC含量(GC content)計算のところ
> count <- alphabetFrequency(reads)
> CG <- rowSums(count[,2:3])
> ACGT <- rowSums(count[,1:4])
> GC_content <- sum(CG)/sum(ACGT)
>
> #出力用に結果をまとめている
> tmp <- NULL
> tmp <- rbind(tmp, c("Total length (bp)", Total_length))
> tmp <- rbind(tmp, c("Number of contigs", Number_of_contigs))
> tmp <- rbind(tmp, c("Average length", Average_length))
> tmp <- rbind(tmp, c("Median length", Median_length))
> tmp <- rbind(tmp, c("Max length", Max_length))
> tmp <- rbind(tmp, c("Min length", Min_length))
> tmp <- rbind(tmp, c("N50", N50))
> tmp <- rbind(tmp, c("GC content", GC_content))
> write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F) #tmpの中$
>
> |

```

hoge.txt - Microsoft Excel

	A	B
1	V1	V2
2	Total length (bp)	123433471.0
3	Number of contigs	46261.0
4	Average length	2668.2
5	Median length	2128.0
6	Max length	101516.0
7	Min length	33.0
8	N50	3688.0
9	GC content	0.48732425

出力ファイル名として指定したもの (hoge.txt) が「20111221」フォルダ中に作成される (はず)

練習

■ 「20111221」中にある**practice1.txt**中の記述を変更して、**Trinity.fasta**ファイルに対して同様の解析を行い、結果を**hoge1.txt**に出力せよ

```
practice1.txt - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
in_f <- "h_rna.fasta"
out_f <- "hoge.txt"

library(Biostrings)
reads <- read.DNAStringSet(in_f, format="fasta")

Total_length <- sum(width(reads))
Number_of_contigs <- length(reads)
Average_length <- mean(width(reads))
Median_length <- median(width(reads))
Max_length <- max(width(reads))
Min_length <- min(width(reads))

#N50計算のところ
sorted_length <- rev(sort(width(reads)))
N50 <- sorted_length[cumsum(sorted_length) >= Total_length/2]

#GC含量(GC content)計算のところ
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
GC_content <- sum(CG)/sum(ACGT)
```



```
practice1.txt - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
in_f <- "Trinity.fasta"
out_f <- "hoge1.txt"

library(Biostrings)
reads <- read.DNAStringSet(in_f, format="fasta")

Total_length <- sum(width(reads))
Number_of_contigs <- length(reads)
Average_length <- mean(width(reads))
Median_length <- median(width(reads))
Max_length <- max(width(reads))
Min_length <- min(width(reads))

#N50計算のところ
sorted_length <- rev(sort(width(reads)))
N50 <- sorted_length[cumsum(sorted_length) >= Total_length/2]

#GC含量(GC content)計算のところ
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
GC_content <- sum(CG)/sum(ACGT)
```

	A	B
1	V1	V2
2	Total length (bp)	2993
3	Number of contigs	4
4	Average length	748.3
5	Median length	784
6	Max length	888
7	Min length	537
8	N50	886
9	GC content	0.524

multi-fasta形式ファイルからの情報抽出2

• 解析 | 一般 | GC含量 (GC contents)

GC含量 (GC contents)の計算の仕方を書きます。ここでは、[ファイルの読み込み \(FASTA形式\)](#)で読み込んだ250 readsからなる `test1.fasta` ファイルを入力として、250 readsの各配列に対して「description」「C,Gの総数」「A,C,G,Tの総数」「配列長」「%GC含量」をファイルに出力するやり方を例示します。

尚、ここでは%GC含量の計算を「CGの総数/ACGTの総数」で計算していますので、もしGC含量を計算したい配列中にNなどが含まれる場合でNなどを含めた「配列長」を分母にしたい場合にはGC含量を計算する数式中の「CG/ACGT*100」を「CG/width(reads)*100」に変更してください。

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

```
----- ここから -----
in_f <- "test1.fasta"
out_f <- "hoge.txt"
library(Biostrings)
reads <- read.DNAStringSet(in_f, format="fasta")
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
out <- CG/ACGT*100

tmp <- cbind(names(reads), CG, ACGT, width(reads), out)
colnames(tmp) <- c("description", "CG", "ACGT", "Length", "%GC_contents")
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F, col.names=T)
----- ここまで -----
```

#読み込みたいFASTA形式のファイル名を指定してin_fに格納
 #出力ファイル名を指定
 #パッケージの読み込み
 #in_fで指定したファイルの読み込み
 #A,C,G,T,...の数を各配列ごとにカウントした結果をcountに格納
 #C,Gの総数を計算してCGに格納
 #A,C,G,Tの総数を計算してACGTに格納
 #%GC含量を計算してoutに格納
 #ファイルに出力したい情報を連結してtmpに格納
 #列名情報を与えている
 #tmpの中身をout_fで指定したファイル名で保存。

[BioconductorのBiostringsのwebページ](#)

配列ごとのGC含量を計算したいとき

練習

- 「20111221」中にある**practice2.txt**中の記述を変更して、**Trinity.fasta**ファイルに対して同様の解析を行い、結果を**hoge2.txt**に出力せよ

```
practice2.txt - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
in_f <- "test1.fasta"
out_f <- "hoge.txt"
library(Biostrings)
reads <- read.DNAStringSet(in_f, format="fasta")
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
out <- CG/ACGT*100

tmp <- cbind(names(reads), CG, ACGT, width(reads))
colnames(tmp) <- c("description", "CG", "ACGT", "Length")
write.table(tmp, out_f, sep="t", append=F, quot
```

hoge2.txt - Microsoft Excel

	A	B	C	D	E
	description	CG	ACGT	Length	%GC_contents
1	comp59_c0_seq1 le	266	537	537	49.53445065
2	comp371_c0_seq1	577	886	886	65.1241535
3	comp26_c0_seq1 le	289	682	682	42.37536657
4	comp8729_c0_seq1	437	888	888	49.21171171

multi-fasta形式ファイルからの情報抽出3

• インタロダクション | 一般 | multi-fastaファイルから指定した配列長のもののみ抽出

RefSeqのhuman mRNAのmulti-fasta形式のファイル ([h_rna.fasta](#))が手元にあったとして、任意の配列長 (例: ≥ 200 bp) のもののみ抽出するやり方を示します。

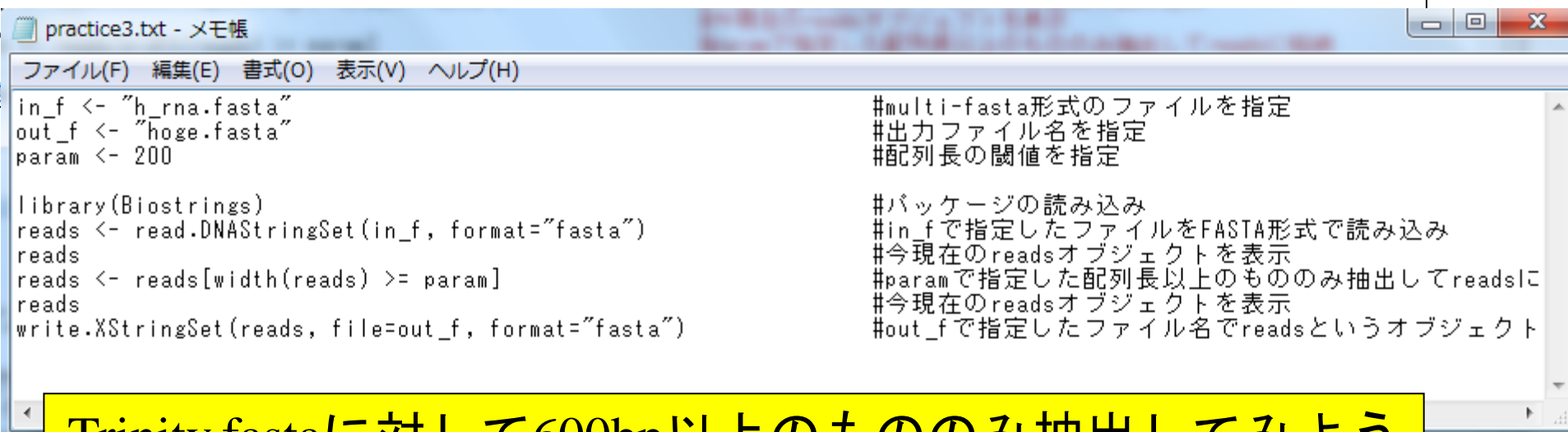
「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

```
----- ここから -----
in_f <- "h_rna.fasta"
out_f <- "hoge.fasta"
param <- 200

library(Biostrings)
reads <- read.DNAStringSet(in_f, format="fasta")
reads
reads <- reads[width(reads) >= param]
reads
write.XStringSet(reads, file=out_f, format="fasta")
```

#multi-fasta形式のファイルを指定
#出力ファイル名を指定
#配列長の閾値を指定

#パッケージの読み込み
#in_fで指定したファイルをFASTA形式で読み込み
#今現在のreadsオブジェクトを表示
#paramで指定した配列長以上のもののみ抽出してreadsに格納
#今現在のreadsオブジェクトを表示
#out_fで指定したファイル名でreadsというオブジェクトをfasta形式で保存



```
practice3.txt - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)

in_f <- "h_rna.fasta"
out_f <- "hoge.fasta"
param <- 200

library(Biostrings)
reads <- read.DNAStringSet(in_f, format="fasta")
reads
reads <- reads[width(reads) >= param]
reads
write.XStringSet(reads, file=out_f, format="fasta")

#multi-fasta形式のファイルを指定
#出力ファイル名を指定
#配列長の閾値を指定

#パッケージの読み込み
#in_fで指定したファイルをFASTA形式で読み込み
#今現在のreadsオブジェクトを表示
#paramで指定した配列長以上のもののみ抽出してreadsに
#今現在のreadsオブジェクトを表示
#out_fで指定したファイル名でreadsというオブジェクト
```

Trinity.fastaに対して600bp以上のもののみ抽出してみよう

multi-fasta形式ファイルからの情報抽出4

• 前処理 | Trinity出力ファイルからFPKM値を取得

2011年10月20日時点で最もお手軽にトランスクリプトーム配列のde novo assembleをしてくれるのはおそらくTrinity(参考文献1)です。[アセンブルプログラム\(転写物\)](#)のところでは、kをいろいろやって...など書いてますが、Trinityは(trans-ABYSSなどと違ってkの値を複数振ってコンティグの和集合を得てから重複を取り除いていくというような作業ではなくk=25だけでアセンブルを実行しているということもあるのでしょうか、とにかく早い(こちらの環境で数週間→3日程度)設定ファイルの記述など面倒なことはほとんどありません。

ここでは、手元にTrinityを実行して(どこかでやってもらって)得られた[Trinity.fasta](#)ファイルがあるという前提で、description部分に記述されているそのコンティグの発現レベル(FPKM値)の情報などを抽出するやり方を示します。

尚、このファイルは基本的なフォーマット部分のみ人工的に作ったものですので、FPKMの記述以外のところは特に気にしないでください。

「ファイル」→「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

```
----- ここから -----
in_f <- "Trinity.fasta"
out_f <- "output.txt"

library(Biostrings)
reads <- read.DNAStringSet(in_f, format="fasta")
hoge <- strsplit(names(reads), " ", fixed=TRUE)
contigID <- unlist(lapply(hoge, "[", 1))
hoge2 <- unlist(lapply(hoge, "[", 3))
hoge3 <- strsplit(hoge2, "=", fixed=TRUE)
FPKM <- unlist(lapply(hoge3, "[", 2))
transcript_length <- width(reads)
tmp <- cbind(contigID, FPKM, transcript_length)
write.table(tmp, out_f, sep="t", append=F, quote=F, row.names=F)
```

#読み込みたいファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納

#パッケージの読み込み
#in_fで指定したファイルをFASTA形式で読み込み
#names(reads)中の文字列をスペース(" ")で区切った結果をリスト形式
#hogeのリスト中の1番目の要素(コンティグID部分に相当)のみ抽出してcon
#hogeのリスト中の3番目の要素(FPKM部分に相当)のみ抽出してhoge2に格納
#hoge2中の文字列を"="で区切った結果をリスト形式でhoge3に格納
#hoge3のリスト中の2番目の要素(FPKMの実際の値部分に相当)のみ抽出して
#配列長情報をtranscript_lengthに格納
#「コンティグID」、「FPKM値」、「配列長」を結合してtmpに格納
#tmpの中身をout_fで指定したファイル名で保存。

----- ここまで -----

(参考文献1(Trinityの原著論文: Grabherr et al., Nature Biotechnol, 2011))

FPKM値：配列長補正ずみの発現量に相当する値

multi-fasta形式ファイルからの情報抽出4

Trinity.fasta

```
>comp59_c0_seq1 len=537 ~FPKM=305.1 path=[0:0-536]
TCATGCCAAAAGGCAGCAATAAGTGCCTTTTCTTCCCTTCAGAATACATGGACAATCCA
AAGCTCTATTAGTCTATTATTCAGAATGAAAAGTGTACAAATATTCGTCTCTTACTCC
TCAGTATGTGAGACTGTTCTCGTAGCAGGTAAATTTCTTCGAATTCAAAACTCCTCAT
GGAAGCATCTGTTTTGTGATCAAGGAGGGGGCTGTATGTGGAATTGCAAGGCCAAAGAC
ATCTCGGGTCAACTCTTCTCAGGACCAAGTCCAGTCCGCGAGTGAAGACACATTCAGG
CAGCCTCCACCAAGGCGCTGCTCAGGAGGAGGCTCCTGTTTATGTGG
GCCCTTGTTCCTCAGCGGGCAGTTGGGGGTCTGGAAGCTAGGAAAGCAAG
CACTCCTGCTTCCTTCTTCCCTGCAGTTGAGACGGGAGTCTTACTTTGTTC
GGTCTCAAACTCCTGGCTTCAAGCAATCCTTCCACTTTGGCCTTCCAAAGTC
>comp371_c0_seq1 len=886 ~FPKM=42 path=[27:0-88 53:8
GCTTCAGTCCAGCACCTTTCTCGGGTCAAGGCTCCTCCTGGCTCCAGGAC
AGGCAGAGGCAGGCTTCCTACACCCCTACTCCTGTGCTCCAGGCTCGAC
GCACTCGAGCACTGAGTCTCTGAGGTCACTTCAAGGTGCTCTCGGCTCACT
TGGACCAAGTGAAGGAGAGGGGTGGGGGCTCGGCTGAGGCACTCCTGCGCT
CTTGTCTACCTCTTGCCCCCGAAGGGTTAGTGTGAGGCTCACTCCAGCATC
TCCTGGTGGCTTGCAGCCCCACAAACCCGAGGTTAAAGCCAGGTACAAC
GACACACCAAGGATGGAGATGTTCCAGGGGCTGCTGCTGTTGCTGCTGCTG
GGGACATGGGCAATCAAGGAGCGGCTTCGGCCACGGTGGCGCCCATCAATC
GCTGTGGAGAGGAGGGCTGCCCGGTGTGATCAAGGTCAACACCAACATCT
TACTGCCCAACATGACCCGCTGCTGCAGGGGGTCTGCCGGGCTGCTGCTG
TCCAACTACCCGCTCTCCGCTTCCAGTCAATCCGCTGCTGCTGCTGCTGCTG
```

output.txt - Microsoft Excel

	A	B	C
1	contigID	FPKM	transcript_length
2	comp59_c0_seq1	305.1	537
3	comp371_c0_seq1	42	886
4	comp26_c0_seq1	4.8	682
5	comp8729_c0_seq1	10.5	888

FPKM値をもとにサンプル内の転写物間の発現レベルの大小を議論可能
 サンプル間の比較には使えない（といわれている）

```
TAGCATTACCAAGGATGAAGTGAAGCAGGATCTGTCTCACCATACACTGAGAACTGTA
```

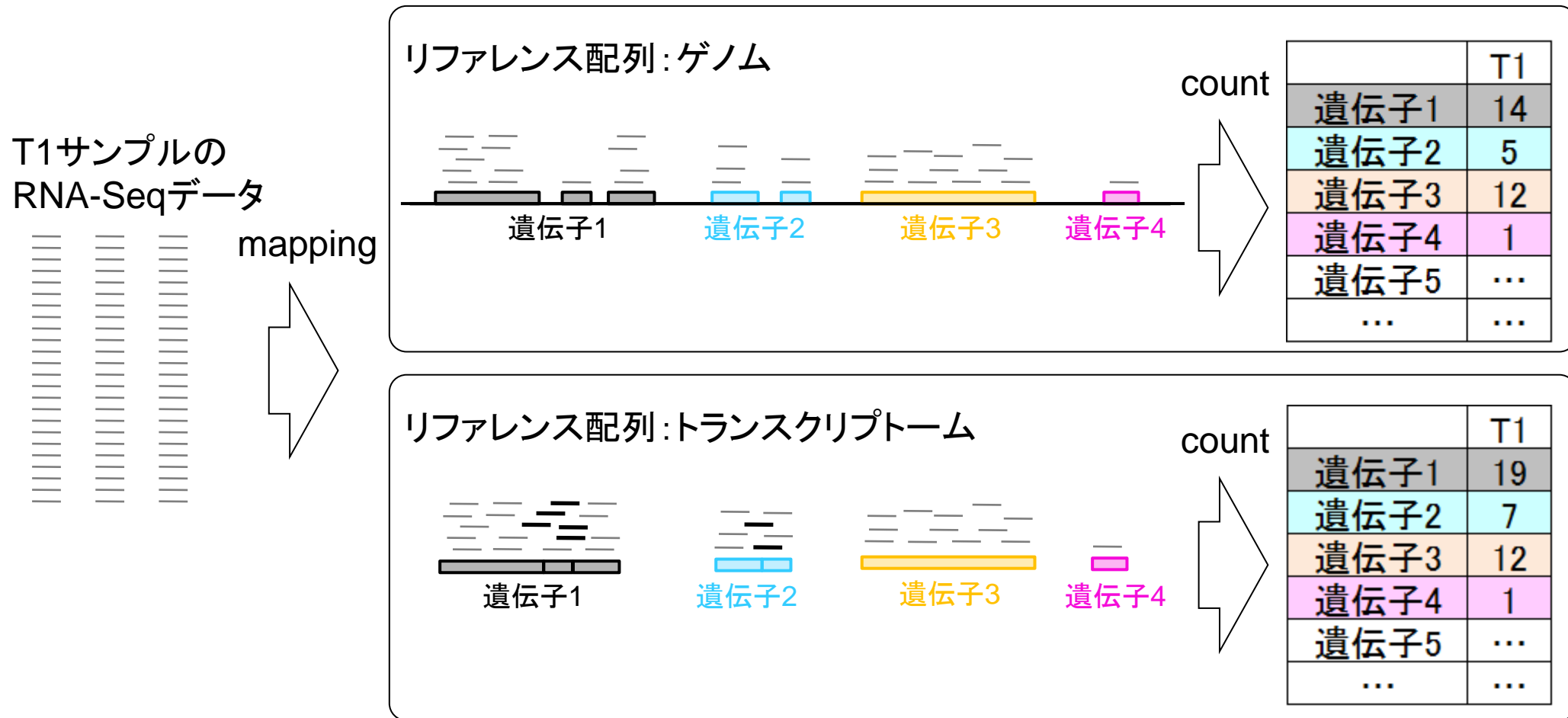
利用可能なRパッケージたち

- *DEGseq* (Wang et al., *Bioinformatics*, **26**: 136-138, 2010)
 - ポワソン分布 (variance = mean) を仮定しているためばらつきを過少評価
- *edgeR* (Robinson et al., *Bioinformatics*, **26**: 139-140, 2010)
 - 正規化法: TMM法
 - 負の二項分布 (variance > mean) を仮定。meanのみのパラメータを用いて現実のばらつきを表現
- *DESeq* (Anders and Huber, *Genome Biol.*, **11**: R106, 2010)
 - 正規化法: RLE法 (relative log expression)
 - edgeRのモデルをさらに拡張(しているらしい)
- *baySeq* (Hardcastle and Kelly, *BMC Bioinformatics*, **11**:422, 2010)
 - 正規化法: RPM (たぶん)
 - 配列の長さ情報を与えるオプションがある
 - データセット中に占めるDEGの割合 (P_{DEG}) を一意に返す
- *NBPSeq* (Di et al., *SAGMB*, **10**:24, 2011)

入力: 生のリードカウントからなる遺伝子発現行列
出力: 遺伝子ごとの発現変動の度合い (p値など)

生のリードカウント？！


■ 基本的なマッピングプログラム (bowtieなど) を用いた場合



理想的な実験デザイン(二群間比較)

■ サンプルA vs. Bの比較(Kidney vs. Liver; 腎臓 vs. 肝臓)

□ 生のリードカウントのデータ(整数値)



Gene ID	A1	A2	A3	A4	...	B1	B2	B3	B4	...
Gene1										
Gene2										
Gene3										
Gene4										
Gene5										
Gene6										
Gene7										
...										

A1: ある生物の腎臓
A2: 同じ生物種の別個体の腎臓
A3: 同じ生物種のさらに別個体の腎臓
...
B1: ある生物の肝臓
B2: 同じ生物種の別個体の肝臓
...

Biological replicatesのデータ
生物学的なばらつき(個体間の違い)を考慮すべし

分布の話

■ 例題: Marioni et al., *Genome Res.*, **18**: 1509-1517, 2008のデータ(の一部)

	kidney (腎臓)					liver (肝臓)				
	A1	A2	A3	A4	A5	B1	B2	B3	B4	B5
EnsemblGeneID	R1 L1 Kidney	R1 L3 Kidney	R1 L7 Kidney	R2 L2 Kidney	R2 L6 Kidney	R1 L2 Liver	R1 L4 Liver	R1 L6 Liver	R1 L8 Liver	R2 L3 Liver
ENSG00000146556	0	0	0	0	0	0	0	0	0	0
ENSG00000197194										
ENSG00000197490										
ENSG00000205292										
ENSG00000177693										
ENSG00000209338										
ENSG00000196573										
ENSG00000177799										
ENSG00000209341										
ENSG00000209342										
ENSG00000209343										
ENSG00000209344										
ENSG00000209346										
ENSG00000209349	0	0	0	0	0	0	0	0	0	0
ENSG00000209350	4	7	3	6	7	35	32	31	29	34
ENSG00000209351	0	0	0	0	0	0	0	0	0	0
ENSG00000209352	0	0	1	1	0	2	0	0	0	0
ENSG00000212679	110	131	149	112	118	177	135	141	148	145
ENSG00000212678	12685	13204	12403	13031	13268	9246	9312	8746	8496	9070
ENSG00000185097	0	0	0	0	0	0	0	0	0	0

Technical replicatesのデータ

サンプル内の技術的なばらつき(例: レーン間の違い)の度合いを調べるためのデータであり、このようなデータで二群間比較し、発現変動遺伝子がどの程度あるかといった数に関する議論は無意味

解析例: アリエナイ?! 数(50%とか)が発現変動遺伝子として検出される

理由: Biological variation > Technical variation

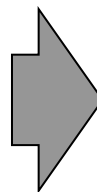
分布の話

■ 例題: Marioni et al., *Genome Res.*, **18**: 1509-1517, 2008のデータ(の一部)

 kidney (腎臓)

EnsemblGeneID	A1	A2	A3	A4	A5
ENSG00000146556	0	0	0	0	0
ENSG00000197194	0	0	0	0	0
ENSG00000197490	0	0	0	0	0
ENSG00000205292	0	0	0	0	0
ENSG00000177693	0	0	0	0	0
ENSG00000209338	0	0	0	0	0
ENSG00000196573	0	0	0	0	0
ENSG00000177799	0	0	0	0	0
ENSG00000209341	0	0	0	0	0
ENSG00000209342	0	0	2	4	3
ENSG00000209343	0	0	0	0	0
ENSG00000209344	0	0	0	0	0
ENSG00000209346	0	0	0	0	0
ENSG00000209349	0	0	0	0	0
ENSG00000209350	4	7	3	6	7
ENSG00000209351	0	0	0	0	0
ENSG00000209352	0	0	1	1	0
ENSG00000212679	110	131	149	112	118
ENSG00000212678	12685	13204	12403	13031	13268
ENSG00000185097	0	0	0	0	0
...
総リード数	1804977	1855190	1742426	1927517	1963420

RPM
正規化



EnsemblGeneID	A1	A2	A3	A4	A5
ENSG00000209342	0.0	0.0	1.1	2.1	1.5
ENSG00000209350	2.2	3.8	1.7	3.1	3.6
ENSG00000209352	0.0	0.0	0.6	0.5	0.0
ENSG00000212679	60.9	70.6	85.5	58.1	60.1
ENSG00000212678	7027.8	7117.3	7118.2	6760.5	6757.6
ENSG00000197049	0.0	0.0	0.0	0.5	0.0
ENSG00000177757	1.1	0.0	1.1	0.5	1.5
ENSG00000177750	0.6	2.2	1.7	1.6	3.6
ENSG00000177741	0.6	0.5	0.0	3.1	0.0
ENSG00000198907	3.3	0.0	3.4	1.0	0.0
ENSG00000187634	27.1	23.2	23.5	21.8	23.9
ENSG00000188976	40.4	41.5	39.0	36.3	41.8
ENSG00000187961	8.3	8.1	7.5	6.2	7.6
ENSG00000187583	0.6	0.5	1.7	0.0	1.5
ENSG00000187642	2.2	2.7	6.9	4.7	4.6
ENSG00000188290	5.0	5.4	6.9	5.2	6.6
ENSG00000187608	6.6	5.9	4.0	8.3	6.6
ENSG00000188157	227.1	223.2	200.9	239.7	240.4
ENSG00000131591	5.5	4.9	4.0	6.2	8.1
ENSG00000215916	5.5	4.9	4.6	6.7	8.7
...
総リード数	1000000	1000000	1000000	1000000	1000000

$$\boxed{12,685} \times \frac{1,000,000}{1,804,977} = \boxed{7027.8}$$

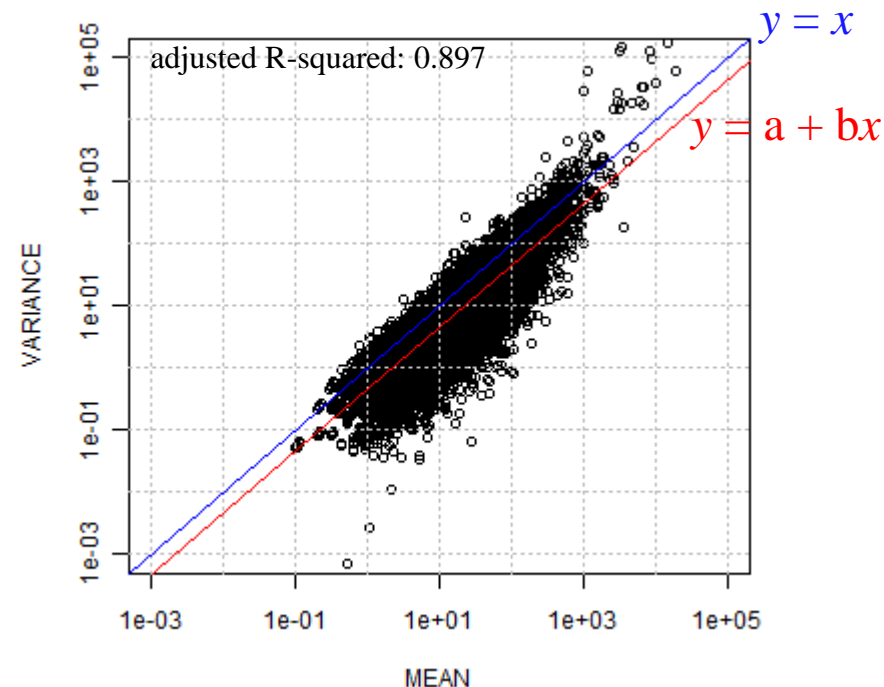
分布の話

■ 例題: Marioni et al., *Genome Res.*, **18**: 1509-1517, 2008のデータ(の一部)



kidney(腎臓)

EnsemblGeneID	A1	A2	A3	A4	A5	MEAN	VARIANCE
ENSG00000209342	0.0	0.0	1.1	2.1	1.5	1.0	0.9
ENSG00000209350	2.2	3.8	1.7	3.1	3.6	2.9	0.8
ENSG00000209352	0.0	0.0	0.6	0.5	0.0	0.2	0.1
ENSG00000212679	60.9	70.6	85.5	58.1	60.1	67.1	129.8
ENSG00000212678	7027.8	7117.3	7118.2	6760.5	6757.6	6956.3	33770.4
ENSG00000197049	0.0	0.0	0.0	0.5	0.0	0.1	0.1
ENSG00000177757	1.1	0.0	1.1	0.5	1.5	0.9	0.4
ENSG00000177750	0.6	2.2	1.7	1.6	3.6	1.9	1.2
ENSG00000177741	0.6	0.5	0.0	3.1	0.0	0.8	1.7
ENSG00000198907	3.3	0.0	3.4	1.0	0.0	1.6	2.9
ENSG00000187634	27.1	23.2	23.5	21.8	23.9	23.9	3.9
ENSG00000188976	40.4	41.5	39.0	36.3	41.8	39.8	5.0
ENSG00000187961	8.3	8.1	7.5	6.2	7.6	7.5	0.7
ENSG00000187583	0.6	0.5	1.7	0.0	1.5	0.9	0.5
ENSG00000187642	2.2	2.7	6.9	4.7	4.6	4.2	3.4
ENSG00000188290	5.0	5.4	6.9	5.2	6.6	5.8	0.8
ENSG00000187608	6.6	5.9	4.0	8.3	6.6	6.3	2.4
ENSG00000188157	227.1	223.2	200.9	239.7	240.4	226.3	258.8
ENSG00000131591	5.5	4.9	4.0	6.2	8.1	5.8	2.5
ENSG00000215916	5.5	4.9	4.6	6.7	8.7	6.1	2.8
...
総リード数	1000000	1000000	1000000	1000000	1000000		

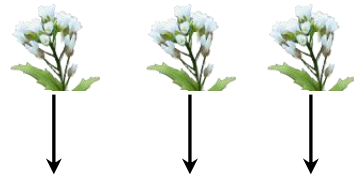


Technical replicatesのデータは:

- ・ (遺伝子の) VARIANCEはそのMEANで説明可能である
- ・ $VARIANCE \approx MEAN$
- ・ ポアソン分布に従う
- ・ ポアソンモデルが適用可能

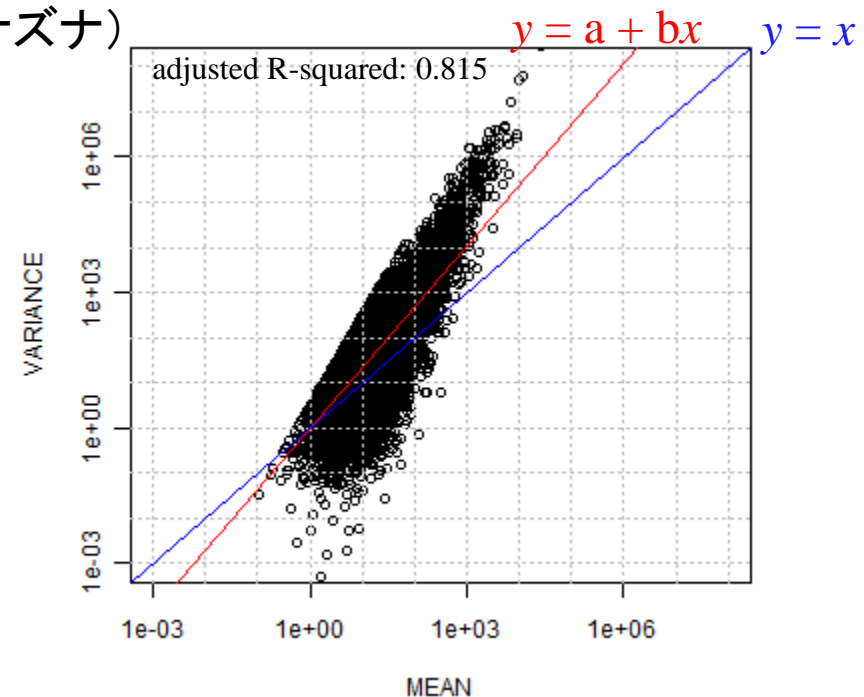
分布の話

■ 例題: Cumbie et al., *PLoS ONE*, 6: e25279, 2011のデータ(の一部)



Arabidopsis (シロイヌナズナ)

	mock1	mock2	mock3	MEAN	VARIANCE
AT1G01010	18.4	39.8	12.3	23.5	209.1
AT1G01020	22.6	23.3	9.8	18.6	57.5
AT1G01030	8.4	12.4	8.0	9.6	6.0
AT1G01040	37.9	22.2	19.6	26.6	97.1
AT1G01050	25.8	40.3	27.6	31.2	62.9
AT1G01060	0.0	7.8	0.6	2.8	18.6
AT1G01070	8.4	17.6	1.8	9.3	62.5
AT1G01080	89.4	98.8	117.2	101.8	200.2
AT1G01090	153.0	178.9	172.7	168.2	183.1
AT1G01100	59.4	64.6	75.5	66.5	67.1
AT1G01110	0.0	0.5	0.3	0.3	0.1
AT1G01120	119.9	97.7	82.8	100.1	347.3
AT1G01130	4.7	5.7	0.3	3.6	8.2
AT1G01140	95.2	62.0	43.6	66.9	683.3
...
総リード数	1000000	1000000	1000000		



Biological replicatesのデータは:

- ・ **VARIANCE > MEAN**
- ・ 負の二項 (NB) 分布に従う
- ・ NBモデルが適用可能

なぜ沢山の方法が存在しているのか？

- *DEGseq* (Wang et al., *Bioinformatics*, **26**: 136-138, 2010) $\text{VAR} = \mu$
 - ポワソン分布 (variance = mean) を仮定しているためばらつきを過少評価
- *edgeR* (Robinson et al., *Bioinformatics*, **26**: 139-140, 2010) $\text{VAR} = \mu(1 + \phi\mu)$
 - 正規化法: TMM法
 - 負の二項分布 (variance > mean) を仮定。
- *DESeq* (Anders and Huber, *Genome Biol.*, **11**: R106, 2010) $\text{VAR} = \mu(1 + \phi_\mu\mu)$
 - 正規化法: RLE法 (relative log expression)
 - edgeRのモデルをさらに拡張 (しているらしい)
- *baySeq* (Hardcastle and Kelly, *BMC Bioinformatics*, **11**: 422, 2010)
 - 正規化法: RPM (たぶん) Ans. VarianceとMeanの関係を表現する手段が沢山あるから
 - 配列の長さ情報を与えるオプションがある
 - データセット中に占めるDEGの割合 (P_{DEG}) を一意に返す
- *NBPSeq* (Di et al., *SAGMB*, **10**: 24, 2011) $\text{VAR} = \mu(1 + \phi\mu^{\alpha-1})$

edgeRを試してみる

■ 例題: Marioni et al., *Genome Res.*, **18**: 1509-1517, 2008のデータ



kidney (腎臓)



liver (肝臓)

	A1	A2	A3	A4	A5	B1	B2	B3	B4	B5
EnsemblGeneID	R1 L1 Kidney	R1 L3 Kidney	R1 L7 Kidney	R2 L2 Kidney	R2 L6 Kidney	R1 L2 Liver	R1 L4 Liver	R1 L6 Liver	R1 L8 Liver	R2 L3 Liver
ENSG00000146556	0	0	0	0	0	0	0	0	0	0
ENSG00000197194	0	0	0	0	0	0	0	0	0	0
ENSG00000197490	0	0	0	0	0	0	0	0	0	0
ENSG00000205292	0	0	0	0	0	0	0	0	0	0
ENSG00000177693	0	0	0	0	0	0	0	0	0	0
ENSG00000209338	0	0	0	0	0	0	0	0	0	0
ENSG00000196573	0	0	0	0	0	0	0	0	0	0
ENSG00000177799	0	0	0	0	0	0	0	0	0	0
ENSG00000209341	0	0	0	0	0	0	0	0	0	0
ENSG00000209342	0	0	2	4	3	0	0	0	1	0
ENSG00000209343	0	0	0	0	0	0	0	0	0	0
ENSG00000209344	0	0	0	0	0	0	0	0	0	0
ENSG00000209346	0	0	0	0	0	0	0	0	0	0
ENSG00000209349	0	0	0	0	0	0	0	0	0	0
ENSG00000209350	4	7	3	6	7	35	32	31	29	34
ENSG00000209351	0	0	0	0	0	0	0	0	0	0
ENSG00000209352	0	0	1	1	0	2	0	0	0	0
ENSG00000212679	110	131	149	112	118	177	135	141	148	145
ENSG00000212678	12685	13204	12403	13031	13268	9246	9312	8746	8496	9070
ENSG00000185097									0	0

ファイル名: SupplementaryTable2_changed.txt

内容: A群が最初の5列、B群が残りの5列のデータ

解析結果をhoge2.txtという名前でファイルに出力したい

edgeRを試してみる

• 解析 | NGS(RNA-seq) | 発現変動遺伝子 | 二群間 | edgeR (Robinson_2010)

2011/12/26にDispersionの計算方法をcommon dispersionからtag-wise dispersionに変更しました。ご注意ください。

参考文献1のedgeRパッケージを用いて解析を行います。edgeRはempirical Bayes)を実装したものです。(おそらくedgeRパッケージの使用例中ではTMM法で得られた正規化係数を用いた位置に組み込めばいいです。また、参考文献2によりますが、このような極端な例でなくても常にTMM法で得られた正規化係数は1に近い値となるので、入力ファイルは、“遺伝子発現行列”形式のもので、ここでは、[サンプルデータ2](#) (つまり[Supplementary](#)

「ファイル」→「ディレクトリの変更」で解析したいファイルを選択します。

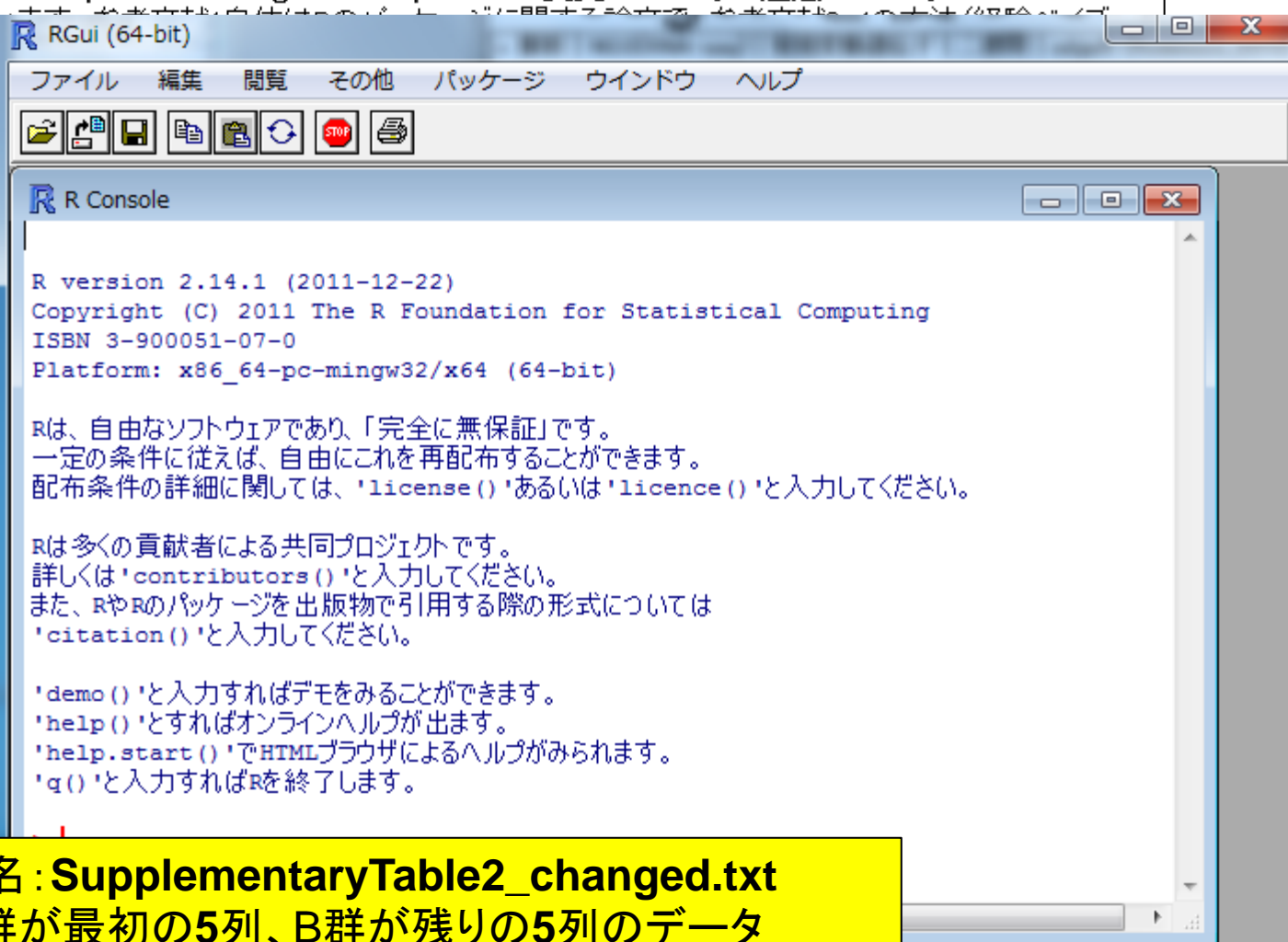
1. 基本形

```
----- ここから -----
in_f <- "SupplementaryTable2_changed.txt"
out_f <- "hoge2.txt"
param1 <- 5
param2 <- 5

library(DEGseq)
library(edgeR)
data <- read.table(in_f, header=TRUE, row.names=1)
data <- as.matrix(data)

data.cl <- c(rep(1, param1), rep(2, param2))
d <- DGEList(counts=data, group=data.cl)
d <- calcNormFactors(d)
d <- estimateCommonDisp(d)
d <- estimateTagwiseDisp(d, prop.used=0.5, grid=10)
out <- exactTest(d)
fdr <- p.adjust(out$table$p.value, method="fdr")
tmp <- cbind(row.names(data), data)
write.table(tmp, out_f, sep="\t", as.is=TRUE)
----- ここまで -----
```

Feb 22 2012



ファイル名: SupplementaryTable2_changed.txt
内容: A群が最初の5列、B群が残りの5列のデータ
解析結果をhoge2.txtという名前でファイルに出力したい

edgeRを試してみる

• 解析 | NGS(RNA-seq) | 発現変動遺伝子 | 二群間 | edgeR (Robinson_2010)

2011/12/26にDispersionの計算方法をcommon dispersionからtag-wise dispersionに変更しました。ご注意ください。

参考文献1のedgeRパッケージを用いて解析を行うempirical Bayes)を実装したものです。(おそらくパッケージの使用例中ではTMM法で得られた正規化した位置に組み込めばいいです。また、参考文献1ですが、このような極端な例でなくても常にTMM法で得られた正規化係数は1に近い値となるので、入力ファイルは、“遺伝子発現行列”形式のものここでは、[サンプルデータ2](#) (つまりSupplementa

「ファイル」-「ディレクトリの変更」で解析したい

1. 基本形

----- ここから -----

```
in_f <- "SupplementaryTable2_changed.txt"
out_f <- "hoge2.txt"
param1 <- 5
param2 <- 5
```

```
library(DEGseq)
library(edgeR)
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
data <- as.matrix(data)
```

```
data.cl <- c(rep(1, param1), rep(2, param2))
d <- DGEList(counts=data, group=data.cl)
d <- calcNormFactors(d)
d <- estimateCommonDisp(d)
d <- estimateTagwiseDisp(d, prop.used=0.5, grid.length=500)
out <- exactTest(d)
fdr <- p.adjust(out$table$p.value, method="BH")
tmp <- cbind(rownames(data), data, out$table, fdr)
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F)
```

----- ここまで -----
Feb 22 2012

RGui (64-bit)

ファイル 編集 閲覧 その他 パッケージ ウィンドウ ヘルプ Vignettes

R Console

```
> out_f <- "hoge2.txt"
> param1 <- 5
> param2 <- 5
>
> library(DEGseq)
> library(edgeR)
> data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
> data <- as.matrix(data)
>
> data.cl <- c(rep(1, param1), rep(2, param2))
> d <- DGEList(counts=data, group=data.cl)
Calculating library sizes from column totals.
> d <- calcNormFactors(d)
> d <- estimateCommonDisp(d)
> d <- estimateTagwiseDisp(d, prop.used=0.5, grid.length=500)
> out <- exactTest(d)
Comparison of groups: 2 - 1
> fdr <- p.adjust(out$table$p.value, method="BH")
> tmp <- cbind(rownames(data), data, out$table, fdr)
> write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F)
>
> |
```

R上でスクリプトをコピペ！
(エラーメッセージが出ていなければ hoge2.txtというファイルができています)

edgeRを使ってみる

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
rownames(data)	R1 L1 Kidr	R1 L3 Kidr	R1 L7 Kidr	R2 L2 Kidr	R2 L6 Kidr	R1 L2 Live	R1 L4 Live	R1 L6 Live	R1 L8 Live	R2 L3 Live	logConc	logFC	p.value	fdr
ENSG00000146556	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000197194	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000197490	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000205292	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000177693	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000209338	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000196573	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000177799	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000209341	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000209342	0	0	2	4	3	0	0	0	1	0	-21.4867	-2.44627	0.179885	0.242167
ENSG00000209343	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000209344	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000209346	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000209349	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000209350	4	7	3	6	7	35	32	31	29	34	-17.0288	3.299663	1.78E-40	5.60E-40
ENSG00000209351	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000209352	0	0	1	1	0	2	0	0	0	0	-22.0717	0.72357	1	1
ENSG00000212679	110	131	149	112	118	177	135	141	148	145	-13.6603	0.98883	4.88E-22	1.37E-21
ENSG00000212678	12685	13204	12403	13031	13268	9246	9312	8746	8496	9070	-7.35443	0.197539	6.28E-11	1.52E-10
ENSG00000185097	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104
ENSG00000209353	0	0	0	0	0	0	0	0	0	0	-50.0161	0	2.29E-104	8.09E-104

一番右側の数値がFalse Discovery Rate (FDR)
この列(O列)で昇順にソートすれば任意の閾値
を満たす遺伝子数がわかる

- ・19,785個がFDR < 0.01を満たす
- ・21,291個がFDR < 0.05を満たす



edgeRを使ってみる

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
rownames(data)	R1 L1 Kidr	R1 L3 Kidr	R1 L7 Kidr	R2 L2 Kidr	R2 L6 Kidr	R1 L2 Live	R1 L4 Live	R1 L6 Live	R1 L8 Live	R2 L3 Live	logConc	logFC	p.value	fdr
ENSG00000116285	115	144	115	143	153	1669	1753	1710	1675	1794	-11.8424	4.407151	0	0
ENSG00000049239	183	232	179	207	199	838	822	814	773	895	-12.0808	2.773867	0	0
ENSG00000186510	515	564	516	568	590	6	1	1	3	2	-15.5081	-7.00256	0	0
ENSG00000184908	484	486	463	573	512	4	2	5	4	3	-15.3378	-6.40426	0	0
ENSG00000142949	332	320	312	350	354	732	772	716	711	808	-11.7855	1.888007	0	0
ENSG00000117472	572	614	603	624	688	14	17	15	13	16	-14.1581	-4.64598	0	0
ENSG00000162366	730	782	720	832	866	4	8	7	7	6	-14.6019	-6.21592	0	0
ENSG00000121310	229	223	247	228	239	1832	1805	1812	1693	1954	-11.4022	3.686564	0	0
ENSG00000116171	542	568	545	548	601	1777	1800	1817	1663	1845	-10.7847	2.390384	0	0
ENSG00000162391	435	444	414	455	450	5	2	5	6	7	-15.1986	-5.73479	0	0
ENSG00000116133	632	681	622	733	702	3534	3396	3178	3196	3657	-10.1878	3.054915	0	0
ENSG00000169174	10	8	8	7	13	223	230	221	173	219	-15.281	5.257754	0	0
ENSG00000157131	14	11	13	7	14	1352	1405	1400	1345	1402	-13.7532	7.59514	0	0
ENSG00000021852	10	12	11	4	20	968	1002	969	982	982	-14.0249	7.151354	0	0
ENSG00000132855	82	96	86	76	90	822	874	823	821	885	-12.6749	4.020325	0	0
ENSG00000134243	919	875	849	883	937	86	93	77	75	94	-12.6438	-2.66983	0	0
ENSG00000163399	7334	7494	6959	7702	7744	284	272	272	250	243	-10.2953	-4.0931	0	0
ENSG00000134240	170	189	180	191	199	2161	2229	2166	2019	2393	-11.4316	4.284652	0	0
ENSG00000168509	9	10	7	8	8	696	710	736	666	711	-14.4847	7.112899	0	0
ENSG00000143384	582	626	568	626	618	1164	1236	1126	1134	1234	-11.0288	1.688325	0	0
ENSG00000197956	942	961	886	1071	995	64	56	58	62	60	-12.8347	-3.29241	0	0

Top-ranked geneの生リードカウントを眺めても確かに発現変動 (Kidney << Liver) していることが分かる



edgeRを試してみる

■ M-A plotを描画(FDR < 0.01を満たすものを赤色で表示)

7. MA-plotも描く場合(FDR < 0.01を満たすものを赤色で示したMA-plotをファイルに保存)

```
----- ここから -----
in_f <- "SupplementaryTable2_changed.txt"
out_f1 <- "hoge2.txt"
out_f2 <- "hoge2.png"
param1 <- 5
param2 <- 5
param3 <- 0.01
param4 <- c(800, 400)
```

#読み込みたい発現データファイルを指定してin_fに格納
#出力ファイル名を指定
#出力ファイル名を指定
#A群のサンプル数を指定
#B群のサンプル数を指定
#MA-plot描画時のFDRの閾値を指定
#MA-plotのファイル出力時の横幅(単位はピクセル)と縦幅を指定

```
library(DEGseq)
library(edgeR)
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
data <- as.matrix(data)
data.cl <- c(rep(1, param1), rep(2, param2))
d <- DGEList(counts=data, group=data.cl)
d <- calcNormFactors(d)
d <- estimateCommonDisp(d)
d <- estimateTagwiseDisp(d, prop.used=0.5, grid.length=500)
out <- exactTest(d)
fdr <- p.adjust(out$table$p.value, method="BH")
rank_edgeR <- rank(fdr)
hoge <- cbind(rownames(data), data, out$table, fdr, rank_edgeR)
tmp <- hoge[order(rank_edgeR),]
write.table(tmp, out_f1, sep="\t", append=F, quote=F, row.names=F)
```

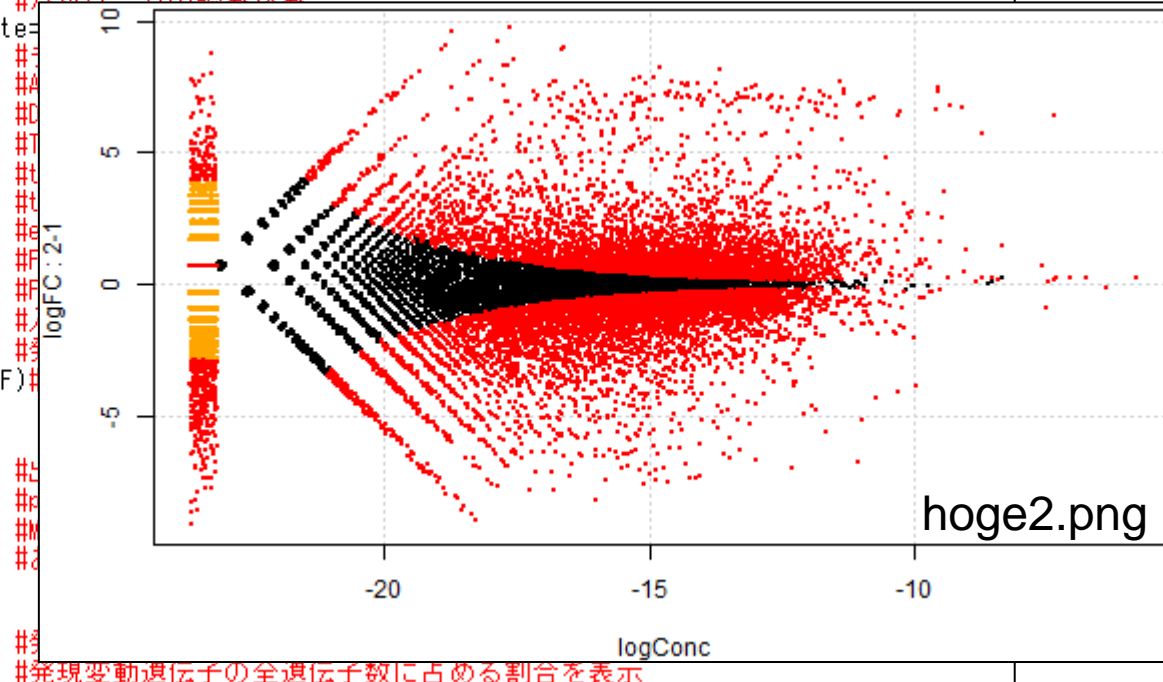
#パッケージの読み込み
#パッケージの読み込み

#MA-plotを描画

```
png(out_f2, width=param4[1], height=param4[2])
obj <- rownames(data)[fdr < param3]
plotSmear(d, de.tags=obj)
dev.off()
```

#おまけ

```
length(obj)
length(obj)/nrow(data)
```



#発現変動遺伝子の全遺伝子数に占める割合を表示

倍率変化がだめな理由をデモ

■ 例題: Marioni et al., *Genome Res.*, **18**: 1509-1517, 2008のデータ



kidney (腎臓)



liver (肝臓)

	A1	A2	A3	A4	A5	B1	B2	B3	B4	B5
EnsemblGeneID	R1 L1 Kidney	R1 L3 Kidney	R1 L7 Kidney	R2 L2 Kidney	R2 L6 Kidney	R1 L2 Liver	R1 L4 Liver	R1 L6 Liver	R1 L8 Liver	R2 L3 Liver
ENSG00000146556	0	0	0	0	0	0	0	0	0	0
ENSG00000197194	0	0	0	0	0	0	0	0	0	0
ENSG00000197490	0	0	0	0	0	0	0	0	0	0
ENSG00000205292	0	0	0	0	0	0	0	0	0	0
ENSG00000177693	0	0	0	0	0	0	0	0	0	0
ENSG00000209338	0	0	0	0	0	0	0	0	0	0
ENSG00000196573	0	0	0	0	0	0	0	0	0	0
ENSG00000177799	0	0	0	0	0	0	0	0	0	0
ENSG00000209341	0	0	0	0	0	0	0	0	0	0
ENSG00000209342	0	0	2	4	3	0	0	0	1	0
ENSG00000209343	0	0	0	0	0	0	0	0	0	0
ENSG00000209344	0	0	0	0	0	0	0	0	0	0
ENSG00000209346	0	0	0	0	0	0	0	0	0	0
ENSG00000209349	0	0	0	0	0	0	0	0	0	0
ENSG00000209350	4	7	3	6	7	35	32	31	29	34
ENSG00000209351	0	0	0	0	0	0	0	0	0	0
ENSG00000209352	0	0	1	1	0	2	0	0	0	0
ENSG00000212679	110	131	149	112	118	177	135	141	148	145
ENSG00000212678	12685	13204	12403	13031	13268	9246	9312	8746	8496	9070
ENSG00000185097	0	0	0	0	0	0	0	0	0	0

発現変動遺伝子がないデータで二群間比較を試みる

A群

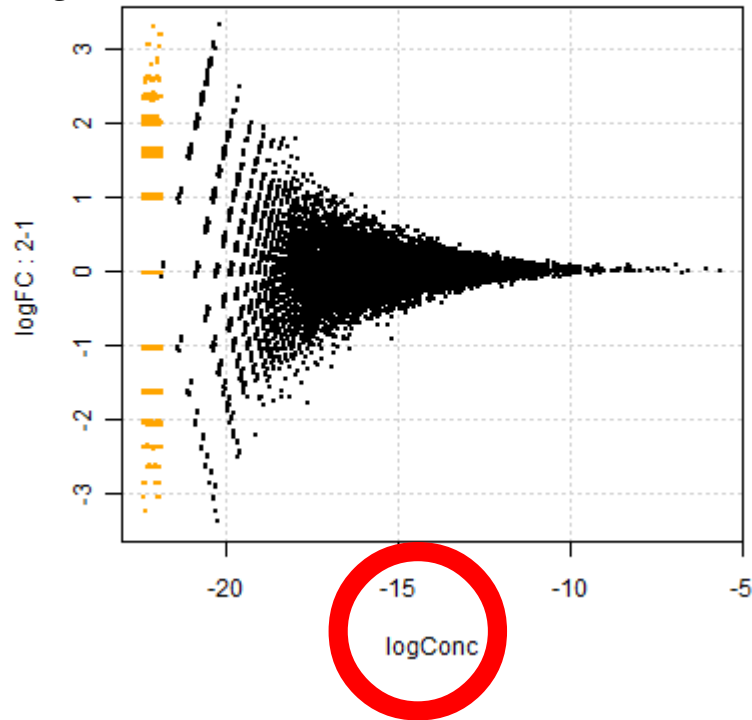
B群

倍率変化がだめな理由をデモ

■ 例題: Marioni et al., *Genome Res.*, **18**: 1509-1517, 2008のデータ(の一部)

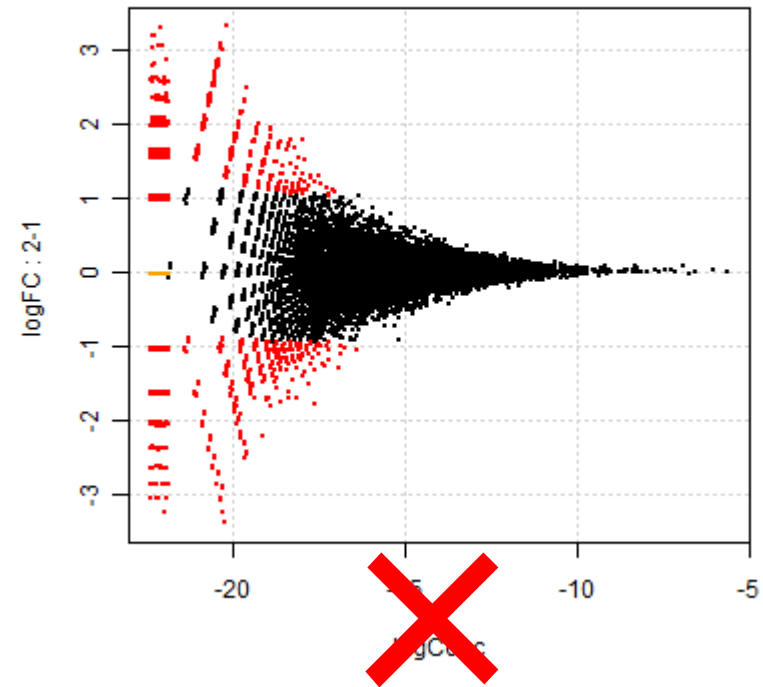
□ (A1, A2) vs. (A3, A4)の二群間比較結果

*edgeR*でFDR < 0.01を満たすものは0個



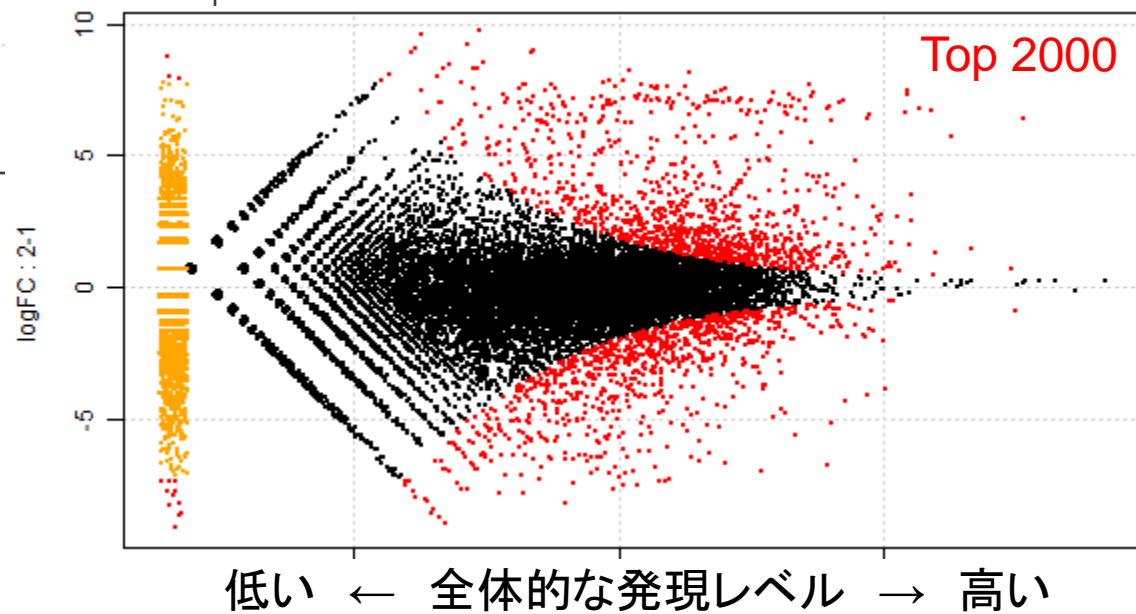
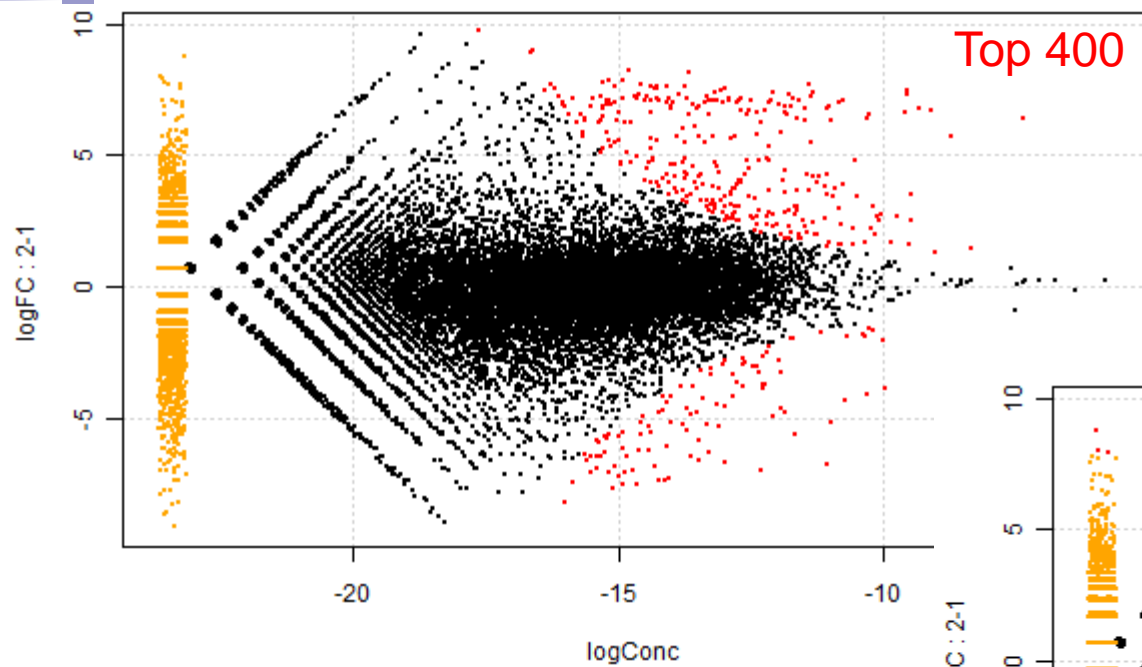
Rcode_edgeR_tech_rep_fdr001.txt

(*edgeR*で)2倍以上発現変動しているものは3814個



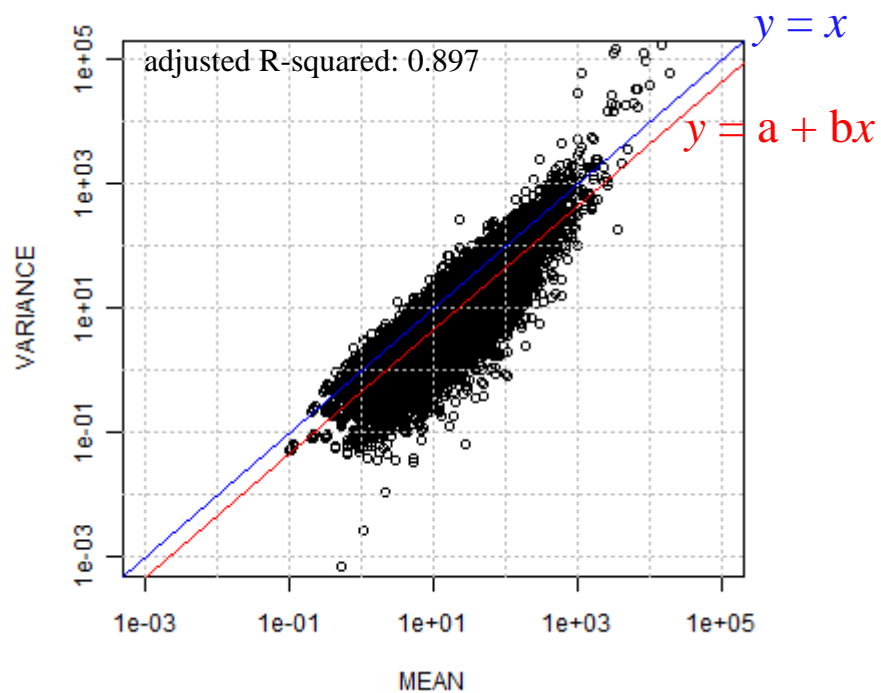
Rcode_edgeR_tech_rep_fc2.txt

低発現領域でlog比が大きくなる現象をうまくモデル化することが重要

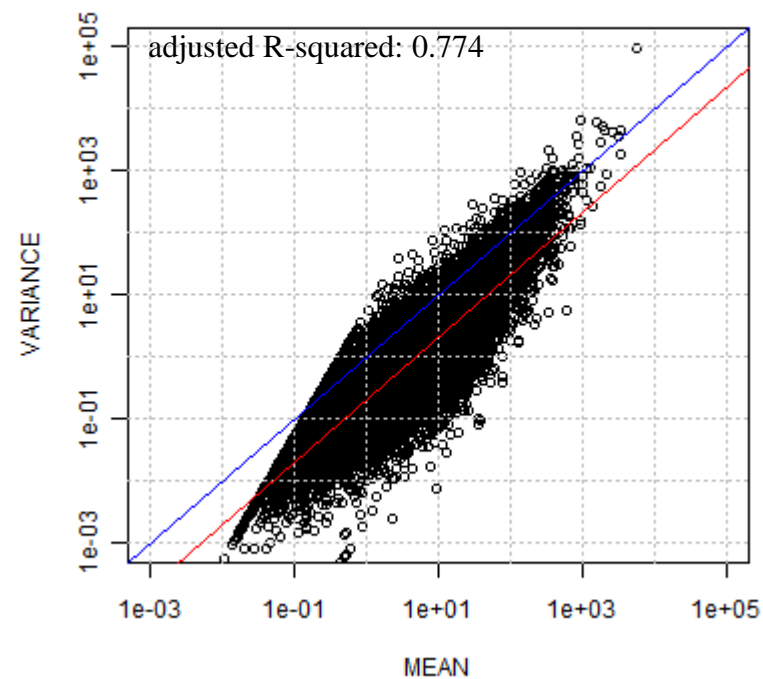


こんな感じでランキングすることが重要です

ちなみに



RPM正規化データ



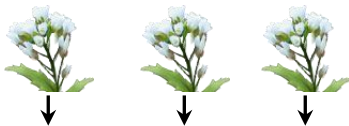

RPKM正規化データ

長さ補正をすると、仮定した分布からのずれが大きくなります...

Biological replicatesの3 vs. 3サンプル

■ 例題: Cumbie et al., *PLoS ONE*, 6: e25279, 2011のArabidopsisデータ

A群 B群

↓ ↓ ↓ ↓ ↓ ↓

	identifier	mock1	mock2	mock3	hrcc1	hrcc2	hrcc3
26,221 genes	AT1G01010	35	77	40	46	64	60
	AT1G01020	43	45	32	43	39	49
	AT1G01030	16	24	26	27	35	20
	AT1G01040	72	43	64	66	25	90
	AT1G01050	49	78	90	67	45	60
	AT1G01060	0	15	2	0	21	8
	AT1G01070	16	34	6	9	20	1
	AT1G01080	170	191	382	127	98	184
	AT1G01090	291	346	563	171	116	453
	AT1G01100	113	125	246	78	27	361
	AT1G01110	0	1	1	0	0	0

data_arab.txt

オリジナルは” AT4G32850”のものが重複して存在していたため19520行目のデータを予め除去している

edgeRをdefaultの手順(edgeR/default)で実行

```
Rcode_edgeR_default.txt - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)

in_f <- "data_arab.txt"
out_f1 <- "result_edgeR_default.txt"
out_f2 <- "result_edgeR_default.png"
param1 <- 3
param2 <- 3
param3 <- 0.05
param4 <- c(600, 400)

library(DEGseq)
library(edgeR)
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
data <- as.matrix(data)
data.cl <- c(rep(1, param1), rep(2, param2))
d <- DGEList(counts=data, group=data.cl)
d <- calcNormFactors(d)
d <- estimateCommonDisp(d)
d <- estimateTagwiseDisp(d, prop.used=0.5, grid.length=500)
out <- exactTest(d)
fdr <- p.adjust(out$table$p.value, method="BH")
rank_edgeR <- rank(fdr)
hoge <- cbind(rownames(data), data, out$table, fdr, rank_edgeR)
tmp <- hoge[order(rank_edgeR),]
write.table(tmp, out_f1, sep="\t", append=F, quote=F, row.names=F)

#MA-plotを描画
png(out_f2, width=param4[1], height=param4[2])
obj <- rownames(data)[fdr < param3]
plotSmear(d, de.tags=obj)
dev.off()
length(obj)
length(obj)/nrow(data)
```

```
#読み込みたい発現データファイルを指定してin_fに格納
#出力ファイル名を指定
#出力ファイル名を指定
#A群のサンプル数を指定
#B群のサンプル数を指定
#MA-plot描画時のFDRの閾値を指定
#MA-plotのファイル出力時の横幅(単位はピクセル)と縦幅を指定

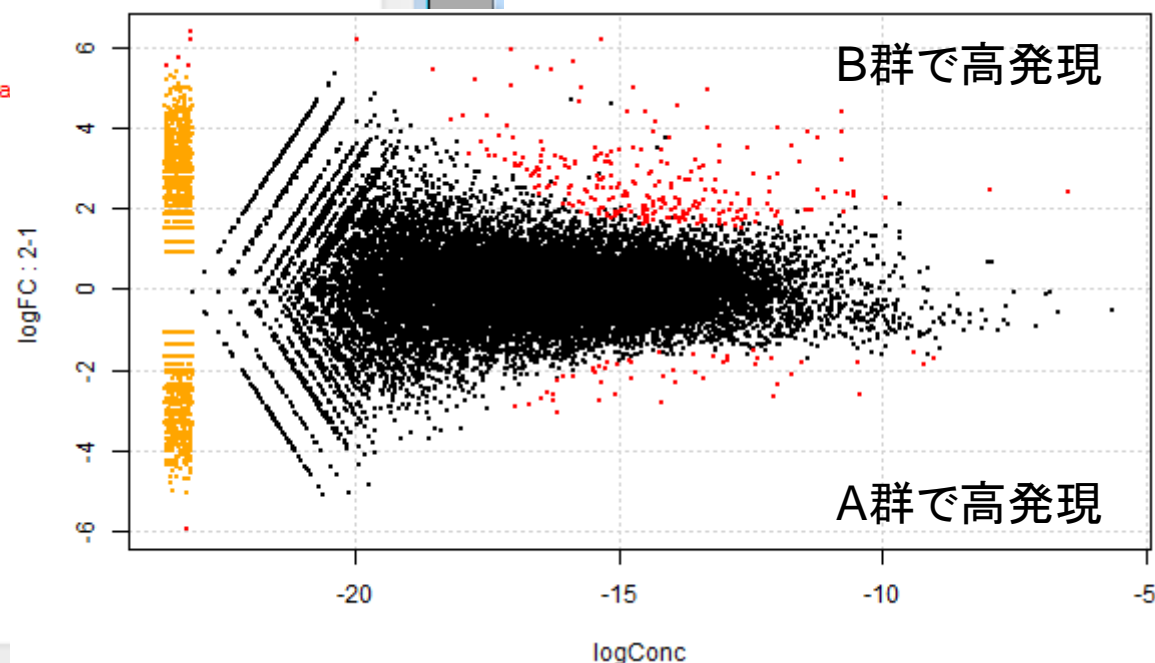
#パッケージの読み込み
#パッケージの読み込み
#発現データファイルの読み込み
#データの型をmatrixにしている
#A群を1、B群を2としたベクトルdata.clを作成
#DGEListオブジェクトを作成してdに格納
#TMM正規化(参考文献5)を実行
#the quantile-adjusted conditional maximum likelihood (qCML)
#the quantile-adjusted conditional maximum likelihood (qCML)
#exact test (正確確率検定)で発現変動遺伝子を計算した結果をou
#False Discovery Rate (FDR)を計算し、結果をfdrに格納
#FDR値でランキングした結果をrank_edgeRに格納
#入力データの右側に、「logConc (M-A plotのAに相当するもの ; ±
#発現変動の度合いでソートした結果をtmpに格納
#tmpの中身をout_f1で指定したファイル名で保存。

#出力ファイルの各種パラメータを指定
#param3で指定したFDRの閾値を満たす遺伝子名情報をobjに格納
#MA-plotの基本形に加え、発現変動遺伝子に相当する
#おまじない
#発現変動遺伝子数を表示
#発現変動遺伝子の全遺伝子数に占める割合を表示
```

edgeRをdefaultの手順(edgeR/default)で実行

```
RGui (64-bit)
ファイル 編集 閲覧 その他 パッケージ ウィンドウ ヘルプ Vignettes

R Console
> out <- exactTest(d)
Comparison of groups: 2 - 1
> fdr <- p.adjust(out$table$p.value, method="BH")
> rank_edgeR <- rank(fdr)
> hoge <- cbind(rownames(data), data, out$table, fdr, ra)
> tmp <- hoge[order(rank_edgeR),]
> write.table(tmp, out_f1, sep="\t", append=F, quote=F,
>
> #MA-plotを描画
> png(out_f2, width=param4[1], height=param4[2])
> obj <- rownames(data)[fdr < param3]
> plotSmear(d, de.tags=obj)
> dev.off()
null device
1
> length(obj)
[1] 318
> length(obj)/nrow(data)
[1] 0.01212768
>
> |
```



クオリティチェック | NGS(一般) | qrgc (Quick Read Quality Control)

FastQCのR版のようなものです。Sanger FASTQ形式ファイルを読み込んで、positionごとの「クオリティスコア (quality score)」、「どんな塩基が使われているのか (base frequency and base proportion)」、「リード長の分布」、「GC含量」、「htmlレポート」などを出力してくれます。ここでは [SRR037439.fastq](#) ファイルに対して解析を行う例を示します。

下記を実行すると「SRR037439-report」という名前のフォルダが作成されます。中にあるreport.htmlをダブルクリックするとhtmlレポートを見ることができます。

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

```
----- ここから -----
in_f <- "SRR037439.fastq"
```

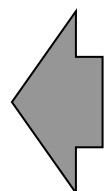
#読み込みたいFASTA形式のファイル名を指定してin_fに格納

```
library(qrgc)
reads <- readSeqFile(in_f, quality="phred")
makeReport(reads)
```

#パッケージの読み込み
#FASTQ形式ファイルの読み込み (Sanger FASTQ形式の場合は"phred", ILLUMINAの場合は"illumina")
#htmlレポートの作成

```
----- ここまで -----
```

[Bioconductorのqrgcのwebページ](#)



「そんなライブラリはない！」と文句を言われたときは、自分でインストール

Quick Read Quality Control

Bioconductor version: Release (2.9)

Quickly scans reads and gathers statistics on base and quality sequences. Produces graphical output of statistics for use in a HTML quality report. S4 SequenceSummary objects allow spe around the data collected.

Author: Vince Buffalo

Maintainer: Vince Buffalo

To install this package, start R and enter:

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

RGui (64-bit)

ファイル 編集 閲覧 その他 パッケージ ウィンドウ ヘルプ



R Console

Rは多くの貢献者による共同プロジェクトです。
詳しくは 'contributors()' と入力してください。
また、RやRのパッケージを出版物で引用する際の形式については
'citation()' と入力してください。

'demo()' と入力すればデモをみることができます。
'help()' とすればオンラインヘルプが出ます。
'help.start()' でHTMLブラウザによるヘルプがみられます。
'q()' と入力すればRを終了します。

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

```
BioC_mirror = http://bioconductor.org
```

```
Change using chooseBioCmirror().
```

```
> biocLite("qrc")
```

```
Using R version 2.13.1, biocinstall version 2.8.4.
```

```
Installing Bioconductor version 2.8 packages:
```

```
[1] "qrc"
```

```
Please wait...
```

```
Installing package(s) into 'C:/Users/kadota/Documents/R/win-library/2.13'  
(as 'lib' is unspecified)
```

```
URL 'http://bioconductor.org/packages/2.8/bioc/bin/windows/contrib/2.13/qrc'
```

```
Content type 'application/zip' length 317872 bytes (310 Kb)
```

```
開かれた URL
```

R上でコピーすれば任意のパッケージをインストールできます

```
The downloaded packages are in
```

```
C:\Users\kadota\AppData\Local\Temp\RtmpZqm3CS\downloaded_packages
```

```
> |
```

サンプル間クラスターリングも重要です

(Rで)マイクロアレイデータ解析 by 門田幸二 (last modified 2011/12/20)

What's new?

- ・ R2.14.0がリリースされ
- ・ 最新の論文([Kadota et al. 2011](#))
- ・ GSA (Efron 2007)の中
- ・ Hook (Binder 2008)を
- ・ Agilent two-color proc
- ・ 作図 | ROC曲線 (ROC
- ・ このページとは直接関
- ありますので、そっち方
- ・ 作図 | ROC曲線 (ROC
- ・ Linksのところにこのペ
- ・ ヒートマップのところに

- ・ はじめに (last modif
- ・ Rのインストールと起

2. サンプル間クラスターリングの場合(類似度:「1-相関係数」、方法:平均連結法(average)):

・ R Graphics画面上に表示したい場合:

```
----- ここから -----
in_f <- "sample3.txt"
param2 <- "average"
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
data.dist <- as.dist(1 - cor(data))
out <- hclust(data.dist, method=param2)
plot(out)
```

#入力ファイル名(発現データファイル)を指定
#方法(method)を指定
#発現データを読み込んでdataに格納。
#サンプル間の距離を計算し、結果をdata.distに格納
#階層的クラスターリングを実行し、結果をoutに格納
#樹形図(デンドログラム)の表示

----- ここまで -----
・ png形式のファイルで図の大きさを指定して得たい場合(Pearson相関係数):

```
----- ここから -----
in_f <- "sample3.txt"
param2 <- "average"
out_f <- "hoge.png"
param3 <- 500
param4 <- 400
param5 <- 14
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
data.dist <- as.dist(1 - cor(data, method="pearson"))
out <- hclust(data.dist, method=param2)
png(out_f, pointsize=param5, width=param3, height=param4)
plot(out)
dev.off()
```

#入力ファイル名(発現データファイル)を指定
#方法(method)を指定
#出力ファイル名(クラスターリング結果ファイル)を指定
#クラスターリング結果の横幅(width; 単位はピクセル)を指定
#クラスターリング結果の縦幅(height; 単位はピクセル)を指定
#クラスターリング結果の文字の大きさ(単位はpoint)を指定
#発現データを読み込んでdataに格納。
#サンプル間の距離を計算し、結果をdata.distに格納
#階層的クラスターリングを実行し、結果をoutに格納
#出力ファイルの各種パラメータを指定
#樹形図(デンドログラム)の表示
#おまじない

----- ここまで -----
・ png形式のファイルで図の大きさを指定して得たい場合(Spearman相関係数):

```
----- ここから -----
in_f <- "sample3.txt"
param2 <- "average"
out_f <- "hoge.png"
```

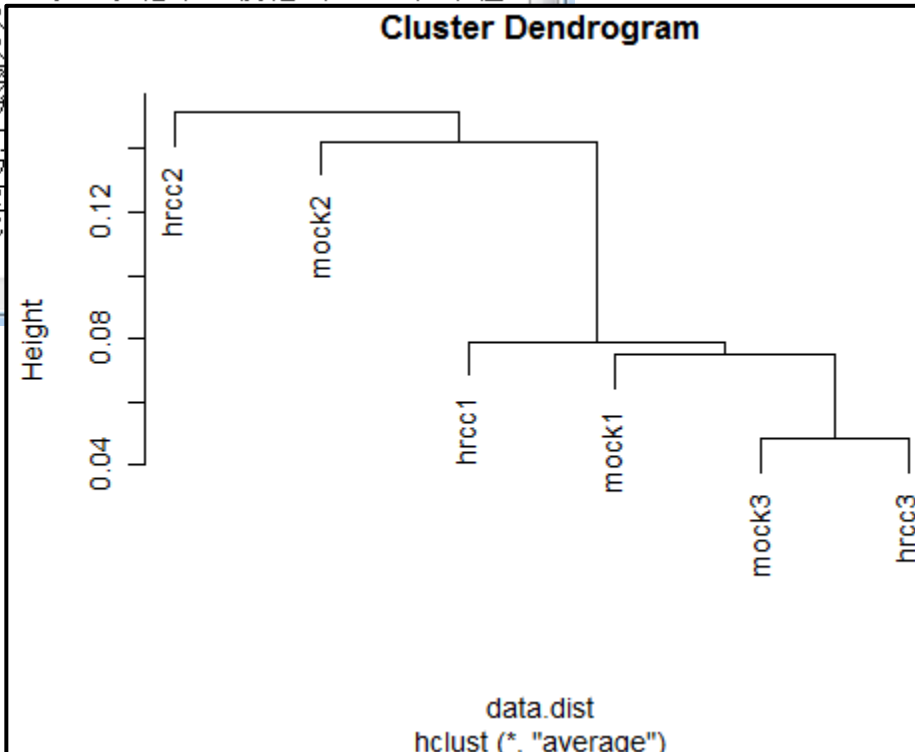
#入力ファイル名(発現データファイル)を指定
#方法(method)を指定
#出力ファイル名(クラスターリング結果ファイル)を指定

サンプル間クラスターリングも重要です

```
Rcode_clustering.txt - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)

in_f <- "data_arab.txt"
param2 <- "average"
out_f <- "result_cluster.png"
param3 <- 500
param4 <- 400
param5 <- 14
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
data.dist <- as.dist(1 - cor(data, method="spearman"))
out <- hclust(data.dist, method=param2)
png(out_f, pointsize=param5, width=param3, height=param4)
plot(out)
dev.off()
```

```
#入力ファイル名(発現データファイル)を指
#方法(method)を指定
#出力ファイル名(クラスターリング結果ファィ
#クラスターリング結果の横幅 (width; 単位:
#クラス
#クラス
#サンプル
#階層的
#出力フ
#樹形図
#おまじ
```



data.dist
hclust (*, "average")

データ中に発現変動遺伝子がありそうかどうかは
クラスターリング結果を眺めるだけでかなりわかる

まとめ(Rでできること)

- NGSで得られたFASTQ形式ファイルの読み込みから、論文でよく出てくるquality score分布やGC含量などの計算結果が得られます
- (de novo transcriptome assemblyなどで得られた)multi-fasta形式ファイルの各種解析が可能です
 - 配列長でのフィルタリング、N50の計算など
- 比較トランスクリプトーム解析用のRパッケージも充実しています
 - 但し、入力データには注意が必要です
 - サンプル間比較: 生のリードカウントデータ
 - サンプル内比較: 長さ補正を行ったデータ(RPKMやFPKMなど)
 - データの分布を考えることは重要です(発現変動遺伝子数を議論したい場合)
 - technical replicatesや**biological replicates**
 - Rパッケージを用いれば発現変動遺伝子の検出から描画まで簡単
 - 「二倍(倍率変化)じゃだめなんです。〇〇さん」

「(Rで)塩基配列解析」のウェブページを用いて...なるべく自力で解析



「予測する生命科学・医療および創薬基盤」

戦略機関 独立行政法人理化学研究所

人材養成プログラム

実施機関 独立行政法人産業技術総合研究所 生命情報工学研究センター



▶ top

▶ outline

▶ seminar

▶ workshop

▶ tutorial

▶ links

▶ contact

チュートリアル

2011年度 HPCI チュートリアルセミナー

2012年3月8日(木) ~9日(金)

- 次世代シーケンサーと共に歩む
- 次世代シーケンサー データ解析入門 -

次世代シーケンサーを用いた研究では、読み取った配列データをどのように解析するかがポイントとなります。このチュートリアルでは、先進的な事例や、誰でも使えるウェブツールやフリーソフトを用いたデータ解析テクニックを紹介します。1人1台のPCを使用しながら、シーケンサー付属ソフトとは一味違った解析が体験できます。



対象

- ・ 次世代シーケンサーを使用している、または使用する予定の方
 > 社会人、学生いずれも可
- ・ コンピュータの基本操作ができる方
 > キーボードからテキスト入力ができる、マウスでウィンドウ操作やコピー・ペーストができる
 > 実習用PCのOSは Windows と Linux です

3/9の門田の担当部分ではもっと詳細な話をする予定です



東京大学大学院農学生命科学研究科

アグリバイオインフォマティクス教育研究ユニット

Agricultural Bioinformatics Research Unit

+ サイトマップ + English



受講生の方へ



研究者の方へ

+ ホーム

+ 本ユニットについて

+ メンバー

+ 教育プログラム

+ 研究フォーラム

+ イベント

+ お問い合わせ

+ リンク

+ モバイルサイト

[ホーム](#) > [教育プログラム](#) > 各講義のページ



各講義のページ

(科目名をクリックすると各講義のページに移動します)

先端
トピックス

セミナー・
討論形式
研究指導

農学生命情報
科学特論 I

農学生命情報
科学特論 II

農学生命情報
科学特論 III

農学生命情報
科学特論 IV

農学生命情報科学特別演習

方法論

講義・実習を
一体化

生物配列統計学

システム生物学概論

知識情報処理論

オーム情報解析

機能ゲノム学

分子モデリングと分子シミュレーション

基 礎

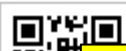
講義・実習を
一体化

ゲノム情報解析基礎

構造バイオインフォマティクス基礎

生物配列解析基礎

バイオスタティスティクス基礎論



東大生以外の方も受講可能です(来年度もやります)