

ゲノム・トランスクリプトームの 各種解析をRで行う

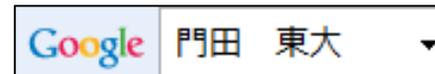
東京大学・大学院農学生命科学研究科

アグリバイオインフォマティクス教育研究プログラム

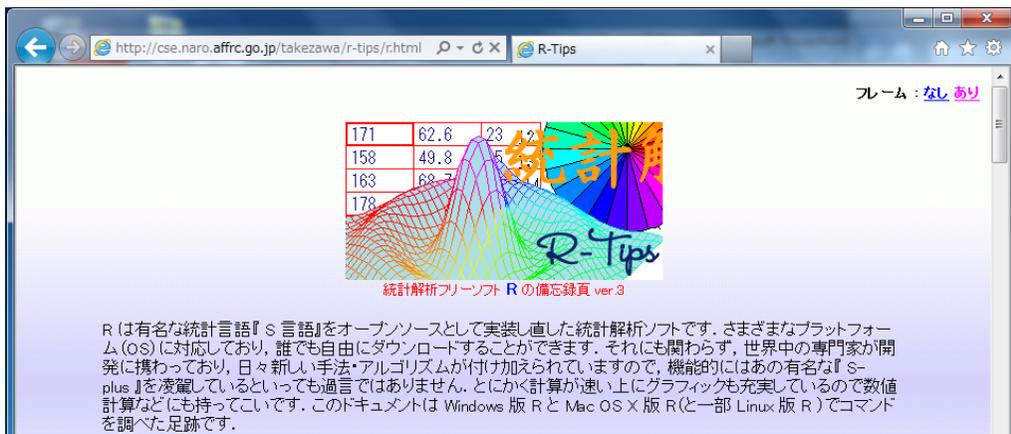
門田幸二(かどた こうじ)

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

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



参考ウェブページ



「(Rで)塩基配列解析」のウェブページを主に利用します

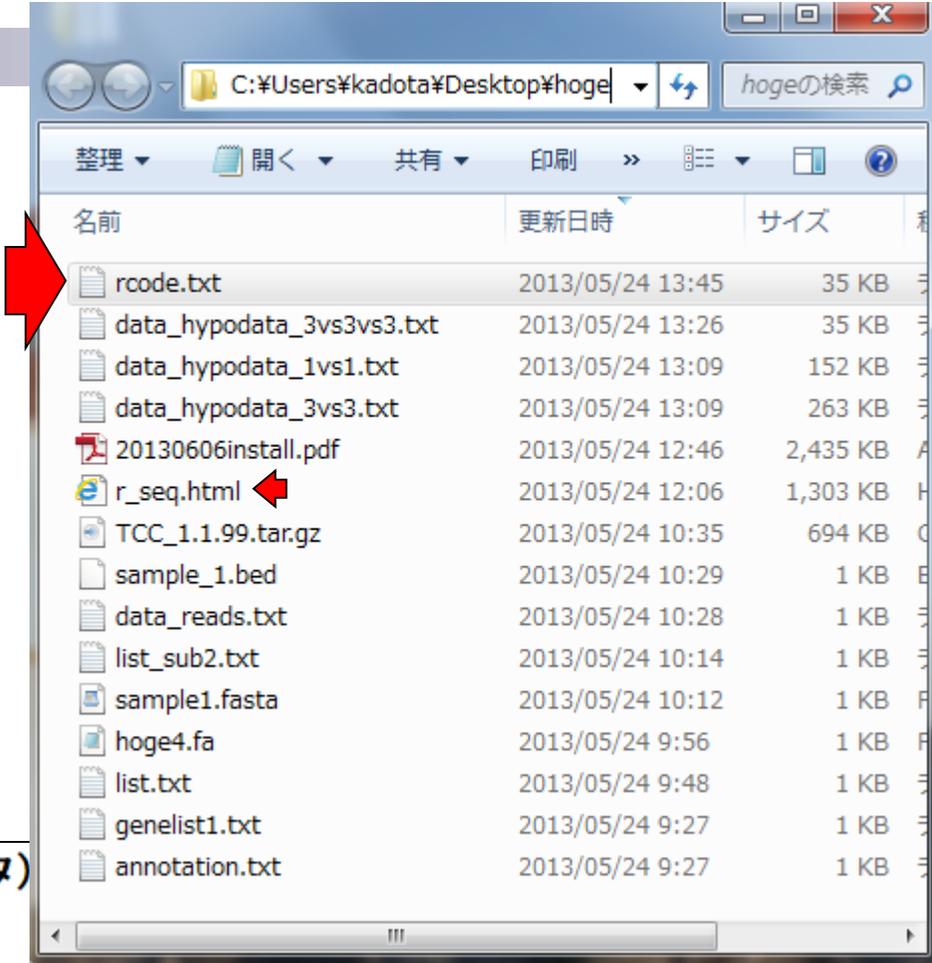
- リンク R-Web (ウェブ版でお試し; なかま)
- リンク RjeWiki (岡田先生)
- リンク PDF版 R-Tips (200頁・3Mb)
- 第01節 Rのセットアップ+参考文献
- 第02節 Rの起動と終了
- 第03節 簡単な計算
- 第04節 R用エディタ

(Rで)塩基配列解析(主に次世代シーケンサーのデータ)

What's new?

RパッケージTCCの最新版は1.1.99です。6月6日に開催されるNAIST植物グローバル教育プロジェクト・平成25年度ワークショップでは、R(ver. 3.0.1)とTCC(ver. 1.1.99)をベースに話を予定しています。Rのバージョン自体は2.15.3でもたぶん大丈夫だとは思いますが。。。Rのインストールと起動を実行したあとにTCC(ver. 1.1.99)のインストールも忘れずに行っておいてください。Windowsのヒト用の詳細なインストール方法は[こちら](#)です。また、当日はhoge.zipという圧縮ファイルを解凍して得られる“hoge”というフォルダがデスクトップ上にあるという前提でセミナーを行いますのでこの圧縮ファイルもダウンロードと解凍をやっておいてください。(2013/05/24)NEW

- 平成25年6月27日、7月3, 4日にこのウェブページ関連の実習を含む講義(農学生命情報科学特論)を行います。東大生以外の外部の方も受講可能です。詳しくは事務局までお問い合わせください。(2013/05/23)NEW
- R3.0.1がリリースされていたのでこれに変更しました。(2013/05/17)NEW





自己紹介

- 2002年3月
 - 東京大学・大学院農学生命科学研究科・応用生命工学専攻 博士課程修了
 - 学位論文:「cDNAマイクロアレイを用いた遺伝子発現解析手法の開発」
(指導教官:清水謙多郎教授)
- 2002/4/1~
 - 産総研・生命情報科学研究センター(CBRC) 産総研特別研究員
 - マイクロアレイ解析手法開発
- 2003/11/1~
 - 放医研・先端遺伝子発現研究センター 研究員
 - 一次元電気泳動波形解析手法開発
- 2005/2/16~
 - 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プログラム
 - マイクロアレイ解析手法開発
 - RNA-seqデータ解析手法開発

(トランスクリプトーム解析周辺の)手法開発系のヒトです



- + 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 <small>セミナー・討論形式 研究指導</small>	農学生命情報科学特別演習			
	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
方法論 <small>講義・実習を一体化</small>	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 <small>講義・実習を一体化</small>	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	

カテゴリー	科目名	学期・単位	実施曜日
基礎	1. 生物配列解析基礎	夏・1	火曜
	生命科学のためのデータベースの利用と基本的な解析手法について講義します。データベースの基礎、配列データベース、機能データベース、ホモロジー検索、モチーフ解析などの基本的な手法について解説します。		
	2. ゲノム情報解析基礎		

講義風景



ゲノム・トランスクリプトーム解析

■ 例: シロイヌナズナ (*Arabidopsis thaliana*)

□ ゲノム配列決定 (chr1-5, chrC, chrM)

- 1番染色体: Theologis et al., *Nature*, **408**: 816-820, 2000
- 2番染色体: Lin et al., *Nature*, **402**: 761-768, 1999
- 3番染色体: Salanoubat et al., *Nature*, **408**: 820-822, 2000
- ...

□ トランスクリプトーム配列 (cDNA配列) 決定

- アノテーション: Seki et al., *Science*, **296**: 141-145, 2002
- ...

□ まとめサイト

- The Arabidopsis Information Resource (TAIR)
- Lamesch et al., *Nucleic Acids Res.*, **40**(Database issue): D1202-1210, 2011
- <http://www.arabidopsis.org/>

	Length	GC contents
chr1	28.76MB	35.80%
chr2	19.60MB	35.80%
chr3	23.17MB	35.40%
chr4	17.40MB	36.02%
chr5	25.95MB	34.50%

Rは論文に書かれている
ことの一部を再現?!できます



ゲノムデータ取得

http://www.arabidopsis.org/ TAIR - Home Page

Home Help Contact About Us Login/Register

Search Browse Tools Portals Download Submit News

The Arabidopsis Information Resource

The Arabidopsis Information Resource (TAIR) maintains a database of biology data for the model higher plant *Arabidopsis thaliana*. Data include the complete genome sequence along with gene structure, gene function, metabolism, gene expression, DNA and seed stocks, genome maps, markers, publications, and information about the Arabidopsis research community data submissions. Gene structures are updated 1-2 times per week using computational and manual methods as well as community submissions. TAIR also provides extensive linkouts from our data pages to other resources.

The Arabidopsis Biological Resource Center at The Ohio State University preserves and distributes seed and DNA resources of *Arabidopsis thaliana*. Stock information and ordering for the ABRC are fully integrated into the TAIR database.

TAIR is located at the Carnegie Institution for Science, Department of Plant Biology and funded by the National Science Foundation with additional support from TAIR sponsors.

Updates on TAIR funding are available [here](#).

ftp://ftp.arabidopsis.org/home/tair/Sequences/

FTP ディレクトリ /home/tair/Sequences/ / ftp

エクスプローラーでこの FTP サイトを表示するには、Alt キーを押して、**FTP サイトを開く**をクリックしてください。

1 階層上のディレクトリへ

11/09/2010 12:00午前	ディレクトリ	ATH cDNA EST sequences FASTA
08/23/2011 12:00午前	ディレクトリ	Polymorphism datasets
11/18/2010 12:00午前	ディレクトリ	blast datasets
08/10/2006 12:00午前	ディレクトリ	chloroplast clones
08/10/2006 12:00午前	ディレクトリ	markers
08/10/2006 12:00午前	ディレクトリ	mitochondrial
02/22/2011 12:00午前	ディレクトリ	whole chromosomes

ゲノムデータ取得

ftp://ftp.arabidopsis.org/home/tair/Sequences/whole_chromosomes/ftp.arabidopsis.org

FTP ディレクトリ /home/tair/Sequences/whole_chromosomes/ftp.arabidopsis.org

エクスプローラーでこの FTP サイトを表示するには、Alt キーを押した状態で「FTP サイトを開く」をクリックしてください。

[1階層上のディレクトリへ](#)

08/10/2006 12:00午前	ディレクトリ	GFF TIGR OLD
02/22/2011 12:00午前	ディレクトリ	OLD
02/22/2011 12:00午前	811	README whole chromosomes.txt
02/22/2011 12:00午前	30,812,908	TAIR10 chr1.fas
02/22/2011 12:00午前	19,947,711	TAIR10 chr2.fas
02/22/2011 12:00午前	23,756,866	TAIR10 chr3.fas
02/22/2011 12:00午前	18,820,386	TAIR10 chr4.fas
02/22/2011 12:00午前	27,317,041	TAIR10 chr5.fas
02/22/2011 12:00午前	156,517	TAIR10 chrC.fas
02/22/2011 12:00午前	371,653	TAIR10 chrM.fas
04/23/2008 12:00午前	140,800	TAIR8 Assembly updates.xls
06/20/2009 12:00午前	123,904	TAIR9 assembly updates.xls
06/20/2009 12:00午前	123,904	TAIR9 assembly updates relative
06/22/2009 12:00午前	10,422	tair9 Assembly gaps.gff



```
arab.fasta - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>Chr1 CHROMOSOME dumped from ADB: Jun/20/09 14:53; last updated: 2009-02-02
CCCTAAACCCTAAACCCTAAACCCTAAACCCTCTGAATCCTTAATCCCTAAATCCCTAAAT
CTTTAAATCCTACATCCATGAATCCCTAAATACCTAATTCCTAAACCCGAAACCGGTTT
CTCTGGTTGAAATCATTGTGTATATAATGATAATTTTATCGTTTTATGTAATTGCTTA
...
>Chr2 CHROMOSOME dumped from ADB: Jun/20/09 14:54; last updated: 2009-02-02
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
...
>Chr3 CHROMOSOME dumped from ADB: Jun/20/09 14:54; last updated: 2009-02-02
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
ACCCTAAACCCTAAACCCTAAACCCTAAACCCTAAATCCATAAATCCCTAAAACCATAAT
...
>Chr4 CHROMOSOME dumped from ADB: Jun/20/09 14:54; last updated: 2009-02-02
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
...
>Chr5 CHROMOSOME dumped from ADB: Jun/20/09 14:54; last updated: 2009-02-02
TATACCATGTACCCTCAACCTTAAACCCTAAACCCTATACTATAAATCTTTAAACCTA
TACTCTAAACCATAGGGTTTGTGAGTTTGCATAAAGTGTACGATAAGTGTCTTCTAACA
TGTGAGTTTGCATAAAGAGTCTCGACTATGTGTTTGTTCAAAAGTGACGTAAGTGTTAGA
...
1行、1列
```

Rで配列長とGC含量計算

	Length	GC contents
chr1	28.76MB	35.80%
chr2	19.60MB	35.80%
chr3	23.17MB	35.40%
chr4	17.40MB	36.02%
chr5	25.95MB	34.50%

```
arab.fasta - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>Chr1 CHROMOSOME dumped from ADB: Jun/20/09 14:53; last updated: 2009-02-02
CCCTAAACCCCTAAACCCCTAAACCCCTAAACCCCTCTGAATCCTTAATCCCTAAATCCCTAAAT
CTTTAAATCCTACATCCATGAATCCCTAAATACCTAATTCCTAAACCCGAAACCGTTT
CTCTGGTTGAAAATCATTTCTCTATAATAATGATAATTTTTCCTTTTTTATCTAATTCCTTA
```

```
R Console
> in_f <- "arab.fasta" #multi-fa$
> out_f <- "result_arab.txt" #出力ファ$
> #必要なパッケージなどをロード #パッケージ$
> library(Biostrings)
> #染色体ごとのファイルを読み込み
> reads <- readDNAStringSet(in_f, #パッケージ$)
> #GC含量(GC content)計算のところ
> count <- alphabetFrequency(reads)
> CG <- rowSums(count[,2:3])
> ACGT <- rowSums(count[,1:4])
> out <- CG/ACGT*100
> #出力用に結果をまとめている
> tmp <- cbind(names(reads), CG, ACGT)
> colnames(tmp) <- c("description", "CG", "ACGT")
> write.table(tmp, out_f, sep="\t", as.is=T)
> |
```

result_arab.txt - Microsoft Excel

	A	B	C	D	E
1	description	CG	ACGT	Length	%GC_contents
2	Chr1 CHROMOSOME	10856525	30263312	30427671	35.87
3	Chr2 CHROMOSOME	7063739	19695728	19698289	35.86
4	Chr3 CHROMOSOME	8521037	23453853	23459830	36.33
5	Chr4 CHROMOSOME	6727440	18582024	18585056	36.20
6	Chr5 CHROMOSOME	9691012	26965224	26975502	35.94

最新のゲノム配列を読み込んでGC含量を計算することができます



Contents

■ Rでゲノム解析

- アノテーションファイル(行列形式)からの情報抽出
 - 条件を満たす行のみの抽出。ハッシュとかgrep周辺がRで可能
- multi-fastaファイルからの情報抽出
 - コンティグごとのGC含量計算。一定の長さ以上のものを抽出とか。

■ Rで(比較)トランスクリプトーム解析

- 研究目的別留意点(サンプル内比較、サンプル間比較)
- リードマッピング → 数値化(カウントデータの取得)
- 発現変動解析
 - 既存のRパッケージの弱点(unbiased DE: ○、biased DE: ×)
 - TCCパッケージを用いて二群間比較(複製あり、複製なし)



アノテーションファイル?!

The Arabidopsis Information Resource (TAIR) maintains a database of biology data for the model higher plant *Arabidopsis thaliana*. Data includes the complete genome sequence along with gene structure, gene metabolism, gene expression, DNA and seed stocks, genome markers, publications, and information about the Arabidopsis research community data submissions. Gene structures are updated 1-2 times per week from the latest published computational and manual methods as well as community submitted genes. TAIR also provides extensive linkouts from our data pages to other resources.

The Arabidopsis Biological Resource Center at The Ohio State University preserves and distributes seed and DNA. Stock information and ordering for the Arabidopsis Biological Resource Center is available at [ABRC](#).

TAIR is located at the Carnegie Institution for Science, Biology and funded by the National Science Foundation.

	A	B	C	D	E	F	G	H	I	J
1	AT1G01010	gene:2200934	AT1G01010.1	located in	nucleus	GO:0005634	537	C	nucleus	ISM
2	AT1G01010	gene:2200934	AT1G01010.1	involved in	regulation of transcription, DNA-	GO:0006355	7461	P	transcript IEA	
3	AT1G01010	gene:2200934	AT1G01010.1	has	DNA binding	GO:0003677	961	F	DNA or RNA IEA	
4	AT1G01010	locus:2200935	AT1G01010	involved in	multicellular organismal develop	GO:0007275	5590	P	developm	ISS
5	AT1G01010	locus:2200935	AT1G01010	has	sequence-specific DNA binding tr	GO:0003700	4449	F	transcript	ISS
6	AT1G01010	locus:2200935	AT1G01010	involved in	amino acid import	GO:0043090	18041	P	transport	RCA
7	AT1G01010	locus:2200935	AT1G01010	involved in	ER to Golgi vesicle-mediated tran	GO:0006888	4768	P	other cell	RCA
8	AT1G01010	gene:2200934	AT1G01010.1	involved in	regulation of transcription, DNA-	GO:0006355	7461	P	other cell	IEA
9	AT1G01010	gene:2200934	AT1G01010.1	involved in	regulation of transcription, DNA-	GO:0006355	7461	P	other cell	IEA
10	AT1G01010	locus:2200935	AT1G01010	involved in	ER to Golgi vesicle-mediated tran	GO:0006888	4768	P	transport	RCA
11	AT1G01020	locus:2200940	AT1G01020	involved in	sterol metabolic process	GO:0016125	7324	P	other met	IMP
12	AT1G01020	locus:2200940	AT1G01020	located in	membrane	GO:0016020	453	C	other mer	ISS
13	AT1G01020	locus:2200940	AT1G01020.1	located in	mitochondrion	GO:0005739	486	C	mitochond	ISS
14										
15										
16										

遺伝子(転写物)ごとに、どの染色体上に存在するのか、どんなGene Ontology IDが割り当てられているのか、などの情報を含むファイル

[1階層上のディレクトリへ](#)

04/30/2013 06:24午前	ディレクトリ	Gene Ontology
08/10/2006 12:00午前	ディレクトリ	OLD TAIR Ontology
04/30/2013 06:23午前	ディレクトリ	Plant Ontology
08/28/2006 12:00午前		7,609 tair2po mapping temporal-060210.txt
08/28/2006 12:00午前		1,203 tair2po mapping temporal-README.txt

[1階層上のディレクトリへ](#)

07/01/2009 12:00午前		3,496 ATH GO-README.txt
04/30/2013 06:24午前		67,144,392 ATH GO GOSLIM.txt
04/30/2013 06:24午前		5,424,109 ATH GO GOSLIM.txt.gz
09/09/2008 12:00午前	ディレクトリ	OLD
03/10/2012 12:00午前		3,727 TAIR GO slim categories.txt

アノテーションファイルからの情報抽出

ATH_GO_GOSLIM.txt

	A	B	C	D	E	F	G	H	I	J
1	AT1G01010	gene:2200934	AT1G01010.1	located in	nucleus	GO:0005634	537	C	nucleus	ISM
2	AT1G01010	gene:2200934	AT1G01010.1	involved in	regulation of transcription, DNA-	GO:0006355	7461	P	transcripti	IEA
3	AT1G01010	gene:2200934	AT1G01010.1	has	DNA binding	GO:0003677	961	F	DNA or RN	IEA
4	AT1G01010	locus:2200935	AT1G01010	involved in	multicellular organismal develop	GO:0007275	5590	P	developm	ISS
5	AT1G01010	locus:2200935	AT1G01010	has	sequence-specific DNA binding tr	GO:0003700	4449	F	transcripti	ISS
6	AT1G01010	locus:2200935	AT1G01010	involved in	amino acid import	GO:0043090	18041	P	transport	RCA
7	AT1G01010	locus:2200935	AT1G01010	involved in	ER to Golgi vesicle-mediated tran	GO:0006888	4768	P	other cell	RCA
8	AT1G01010	gene:2200934	AT1G01010.1	involved in	regulation of transcription, DNA-	GO:0006355	7461	P	other met	IEA
9	AT1G01010	gene:2200934	AT1G01010.1	involved in	regulation of transcription, DNA-	GO:0006355	7461	P	other cell	IEA
10	AT1G01010	locus:2200935	AT1G01010	involved in	ER to Golgi vesicle-mediated tran	GO:0006888	4768	P	transport	RCA
11	AT1G01020	locus:2200940	AT1G01020	involved in	sterol metabolic process	GO:0016125	7324	P	other met	IMP
12	AT1G01020	locus:2200940	AT1G01020	located in	membrane	GO:0016020	453	C	other mer	ISS
13	AT1G01020	gene:2200939	AT1G01020.1	located in	mitochondrion	GO:0005739	486	C	mitochon	ISM
14	AT1G01020	locus:2200940	AT1G01020	involved in	sterol metabolic process	GO:0016125	7324	P	other met	IMP
15	AT1G01020	locus:2200940	AT1G			GO:0003674	3226	F	unknown	ND
16	AT1G01020	locus:2200940	AT1G			GO:0005783	268	C	ER	IDA

「核」に存在する遺伝子のみからなるリストを得たいときにもRが利用可能



アノテーションファイルからの情報抽出

genelist1.txt

	A
1	gene1
2	gene7
3	gene9

annotation.txt

	A	B	C	D
1	genename	accession	description	subcellular_location
2	gene1	hoge01	plasma_mem	nuclear
3	gene2	hoge02	hohinu	membrane
4	gene3	hoge03	agribio	endoplasmic
5	gene4	hoge04	genesis	endoplasmic
6	gene5	hoge05	kamo	membrane
7	gene6	hoge06	netteba	humei
8	gene7	hoge07	tebasaki	nuclear
9	gene8	hoge08	biiru	nuclear
10	gene9	hoge09	nihonshu	nuclear
11	gene10	hoge10	agene1	membrane
12	gene11	hoge11	iyaaaa	endoplasmic



hoge1.txt

	A	B	C	D
1	gene1	hoge01	plasma_mem	nuclear
2	gene7	hoge07	tebasaki	nuclear
3	gene9	hoge09	nihonshu	nuclear

目的: アノテーションファイル (`annotation.txt`) 中の第1列目に対して、リストファイル (`genelist1.txt`) 中の文字列と一致する行を抜き出して、`hoge1.txt` というファイル名で出力したい

(Rで)塩基配列解析(主に次世代シーケンサーのデータ) (last modified 2013/05/24, since 2010)

What's new?
RパッケージTOCCの最新版は1.1.99です。6月6日に開催されるNAIST植物グローバル教育プロジェクト・平成25年度ワークショップでは、R(ver. 3.0.1)とTOC(ver. 1.1.99)をベースに話をする予定です。Rのバージョン自体は2.15.3でもたぶん大丈夫だとは思いますが。。。Rのインストールと起動を実行したあとにTOCC(ver. 1.1.99)のインストールも忘れずに行っておいてください。Windowsのヒト用の詳細なインストール方法は[こちら](#)です。また、当日はhoge.zipという圧縮ファイルを解凍して得られる"hoge"というフォルダがデスクトップ上にあるという前提でセミナーを行いますのでこの圧縮ファイルもダウンロードと解凍をやっておいてください。(2013/05/24)NEW

- ・平成25年6月27日、7月3、4日にこのウェブページ関連の実習を含む講義(農学生命情報科学特論)を行います。東大生以外の外部の方も受講可能です。詳しくは事務局までお問い合わせください。(2013/05/23)NEW
- ・R3.0.1がリリースされていたのでこれに変更しました。(2013/05/17)NEW

- ・はじめに (last modified 2012/03/29)
- ・Rのインストールと起動 (last modified 2013/05/17) NEW
- ・サンプルデータ (last modified 2013/01/23)
- ・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 アン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 アンテーション情報取得(BioMart and Ensembl)
- ・イントロダクション 一般 遺伝子の転写開始点近傍配列(上流)
- ・イントロダクション 一般 任意のキーワードを含む行を抽出
- ・イントロダクション 一般 ランダムな塩基配列を作成 (last modified 2010/12/07)
- ・イントロダクション 一般 任意の長さの可能な全ての塩基配列を作成
- ・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 2010/7/7)
- ・イントロダクション 一般 翻訳配列(translate)を取得 (last modified 2010/7/7)
- ・イントロダクション 一般 相補鎖(complement)を取得 (last modified 2010/7/7)
- ・イントロダクション 一般 逆相補鎖(reverse complement)を取得 (last modified 2010/7/7)

イントロダクション | 一般 | 任意のキーワードを含む行を抽出

例えばタブ区切りテキストファイルのannotation.txtが手元であり、この中からgenelist1.txtのようなリストファイル方を示します。
Linux (UNIX)のgrepコマンドのようなものです。perlのハッシュのようなものです。

「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。

1. 目的のタブ区切りテキストファイル(annotat ion.txt)中の第1列目をキーとして、リストファイル(genelist1.txt)中

```

-----   ここから   -----
in_f1 <- "annotation.txt"
in_f2 <- "genelist1.txt"
out_f <- "hoge1.txt"
param <- 1

#ファイルの読み込み
data <- read.table(in_f1, header=TRUE, sep="\\t", quote="")
keywords <- readLines(in_f2)
dim(data)

#本番
obj <- is.element(as.character(data[,param]), keywords)
out <- data[obj,]
dim(out)
write.table(out, out_f, sep="\\t", append=F, quote=F, row.names=F)

-----   ここまで   -----

```

#in_f1で読み込む目的のファイルの何列目のデータを抽出する

#in_f2で読み込む目的のファイル(リストファイル)を読み込んで、その中のキーワードを抽出する

#オブジェクトdataの行数と列数を表示

#in_f1で読み込んだファイル中の(param)列目のデータを抽出する

#オブジェクトoutの行数と列数を表示

#outの中身をout_fで指定したファイル名で保存する

```

1. 目的のタブ区切りテキストファイル(annotation.txt)中の第1列目をキーとして、リストファイル(genelist.txt)中のものが含まれる行全体を出力したい場合:
-----   ここから   -----
in_f1 <- "annotation.txt"
in_f2 <- "genelist.txt"
out_f <- "hoge1.txt"
param <- 1

#ファイルの読み込み
data <- read.table(in_f1, header=TRUE, sep="\\t", quote="")
keywords <- readLines(in_f2)
dim(data)

#本番
obj <- is.element(as.character(data[,param]), keywords)
out <- data[obj,]
dim(out)
write.table(out, out_f, sep="\\t", append=F, quote=F, row.names=F)
-----   ここまで   -----

```

#入力ファイル名(目的のタブ区切りテキストファイル)を指定
 #入力ファイル名(キーワードなどのリストファイル)を指定
 #出力ファイル名を指定
 #in_f1で読み込む目的のファイルの何列目のデータに対して検索したいかを指定

 #入力ファイル(目的のファイル)を読み込んでdataに格納
 #入力ファイル(リストファイル)を読み込んでkeywordsに格納
 #オブジェクトdataの行数と列数を表示

 #in_f1で読み込んだファイル中の(param)列目の文字列ベクトル中の各要素がベクトルkeywordsの要素と一致する行を抽出する
 #行列dataからobjがTRUEとなる行のみを抽出した結果をoutに格納
 #オブジェクトoutの行数と列数を表示
 #outの中身をout_fで指定したファイル名で保存。

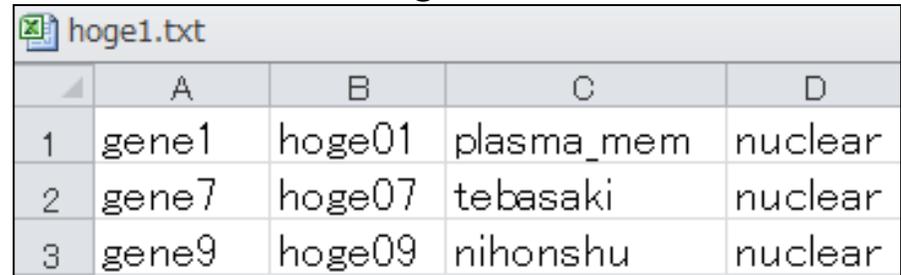
annotation.txt

	A	B	C	D
1	gene1	hoge01	plasma_mem	nuclear
2	gene2	hoge02	hohinu	membrane
3	gene3	hoge03	agribio	endoplasmic
4	gene4	hoge04	genesis	endoplasmic
5	gene5	hoge05	kamo	membrane
6	gene6	hoge06	netteba	humei
7	gene7	hoge07	tebasaki	nuclear
8	gene8	hoge08	biiru	nuclear
9	gene9	hoge09	nihonshu	nuclear
10	gene10	hoge10	agene1	membrane
11	gene11	hoge11	iyaaa	endoplasmic

genelist.txt

	A
1	gene1
2	gene7
3	gene9

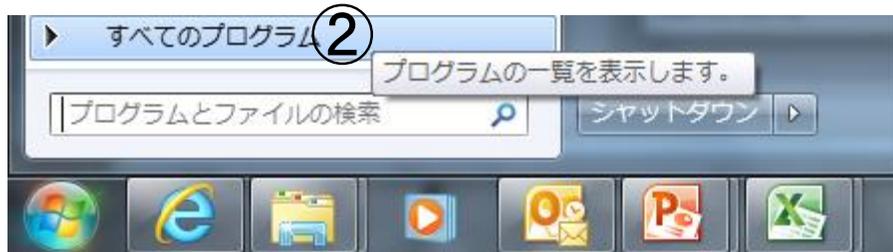
hoge1.txt



	A	B	C	D
1	gene1	hoge01	plasma_mem	nuclear
2	gene7	hoge07	tebasaki	nuclear
3	gene9	hoge09	nihonshu	nuclear

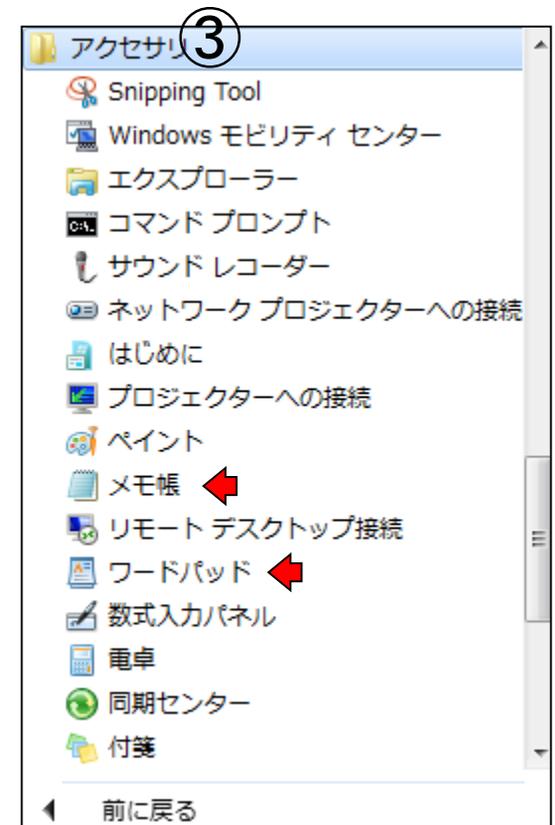
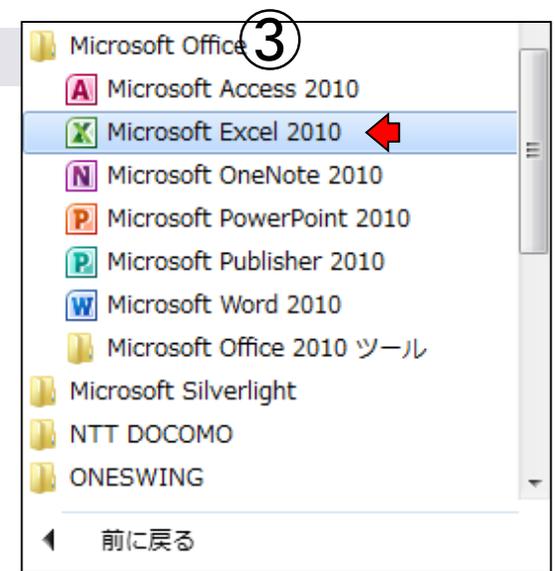
デスクトップ上に「hoge」という名前のフォルダがあり、フォルダ中にannotation.txtとgenelist.txtが存在する、という前提です。
 「メモ帳」で開くと改行コードが崩れている場合は、「ワードパッド」などで開くとよい

各種ソフトの場所



①

スタート – すべてのプログラム – ...



Rの起動および作業ディレクトリの変更

RGui (64-bit) window showing the File menu. The menu item "ディレクトリの変更..." (Change Directory...) is circled with a 1. Below the menu, a terminal window shows the command prompt with a red cursor.

作業ディレクトリの変更 dialog box (Change Working Directory) showing the file explorer view. The "ローカル ディスク (C:)" (Local Disk (C:)) is circled with a 2. A tooltip shows "空き領域: 280 GB" (Free space: 280 GB) and "合計サイズ: 453 GB" (Total size: 453 GB).

Folder (F): ローカル ディスク (C:)

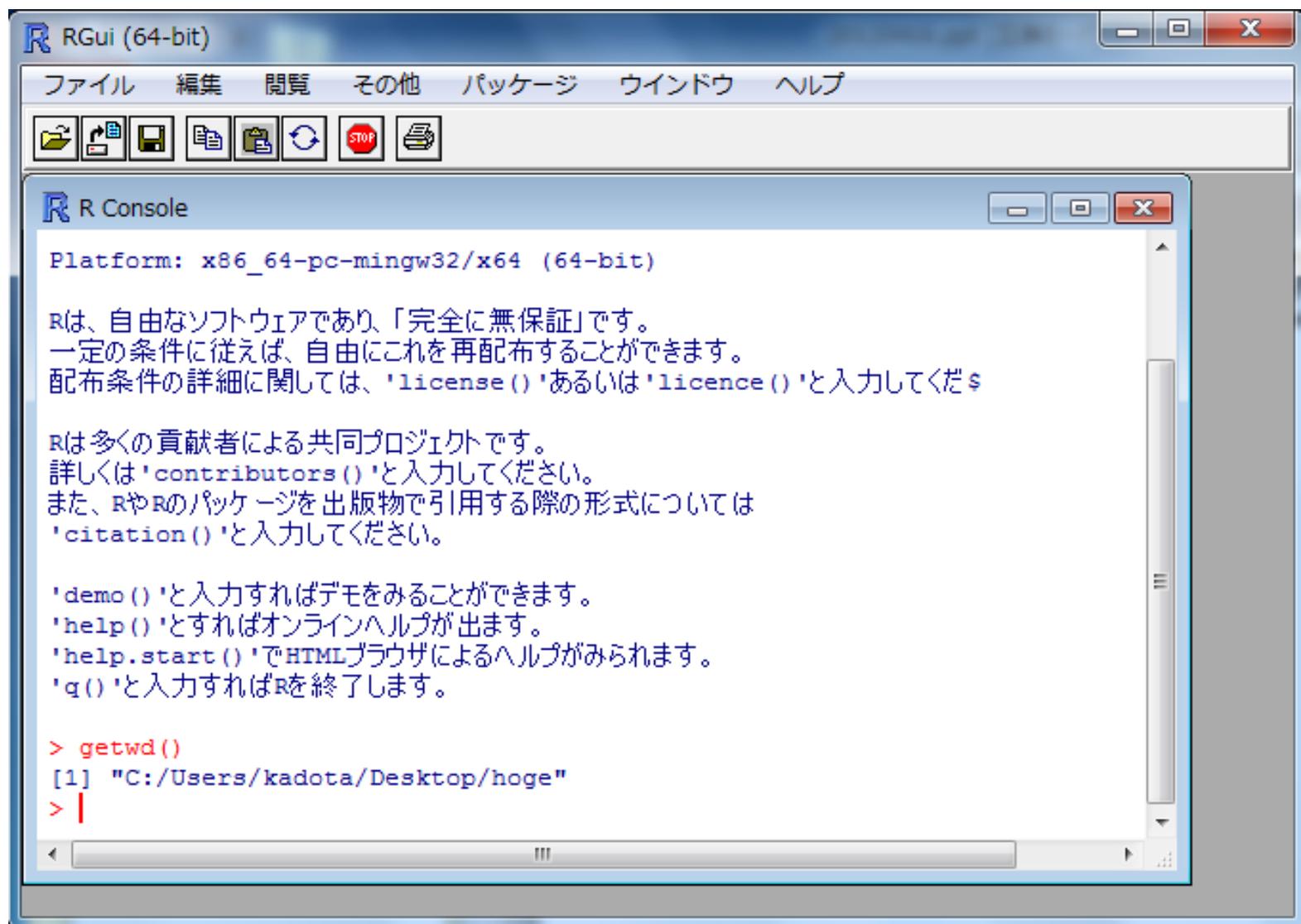
Buttons: 新しいフォルダーの作成(N) (Create new folder), OK, キャンセル (Cancel)

作業ディレクトリの変更 dialog box (Change Working Directory) showing the file explorer view. The "Users" folder is circled with a 3. The "kadota" user folder is circled with a 4. The "デスクトップ" (Desktop) folder is circled with a 5. The "hoge" folder is circled with a 6. A yellow box with a red arrow points to the "kadota" folder, containing the text "④のところは人それぞれ" (The place ④ is different for each person).

Folder (F): hoge

Buttons: 新しいフォルダーの作成(N) (Create new folder), OK (circled with a 7), キャンセル (Cancel)

「getwd()」と打ち込んで確認



The screenshot shows the RGui (64-bit) window with the R Console pane open. The console displays the following text:

```
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

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/hoge"
> |
```

基本はコピー

• イントロダクション | 一般 | 任意のキーワードを含む行を抽出

例えばタブ区切りテキストファイルの `annotation.txt` が手元があり、この中から `genelist1.txt` のようなリストファイル中の文字列を含む行を抽出するやります。

Linux (UNIX) の `grep` コマンドのようなものです。perl を使った方法もあります。

「ファイル」-「ディレクトリの変更」で解析したいファイルの場所へ移動します。

```
1 目的のタブ区切りテキストファイル(annotation.txt)
   -----   ここから   -----
in_f1 <- "annotation.txt"
in_f2 <- "genelist1.txt"
out_f  <- "hogel.txt"
param <- 1
```

```
#ファイルの読み込み
data <- read.table(in_f1, header=TRUE, sep="¥t")
keywords <- readLines(in_f2)
dim(data)
```

```
#本番
obj <- is.element(as.character(data[,param]), keywords)
out <- data[obj,]
dim(out)
write.table(out, out_f, sep="¥t", append=F, quote=F)
```

```
-----   ここまで   -----
```

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

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を見ます。
'q()'と入力すればRを終了します。

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

コピー Ctrl+C
ペースト Ctrl+V
コマンドのみペースト
コピー&ペースト Ctrl+X
ウィンドウの消去 Ctrl+L
全て選択
バッファに出力 Ctrl+W

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

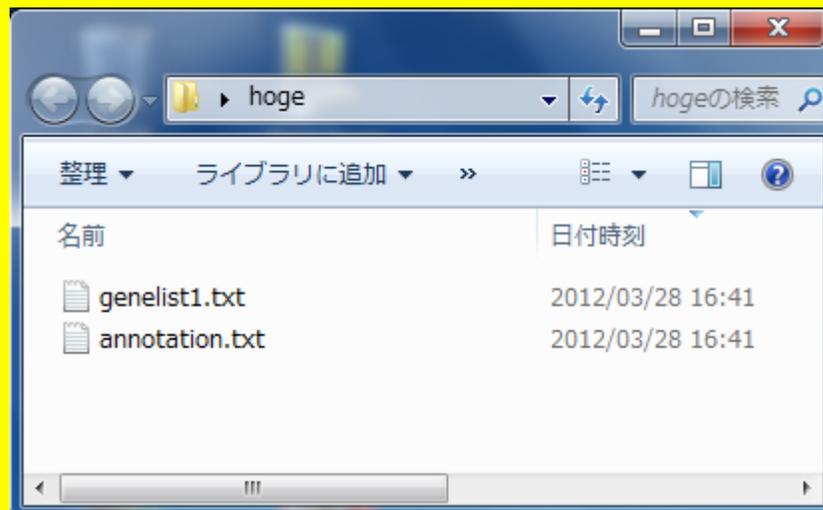
実行結果

```
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

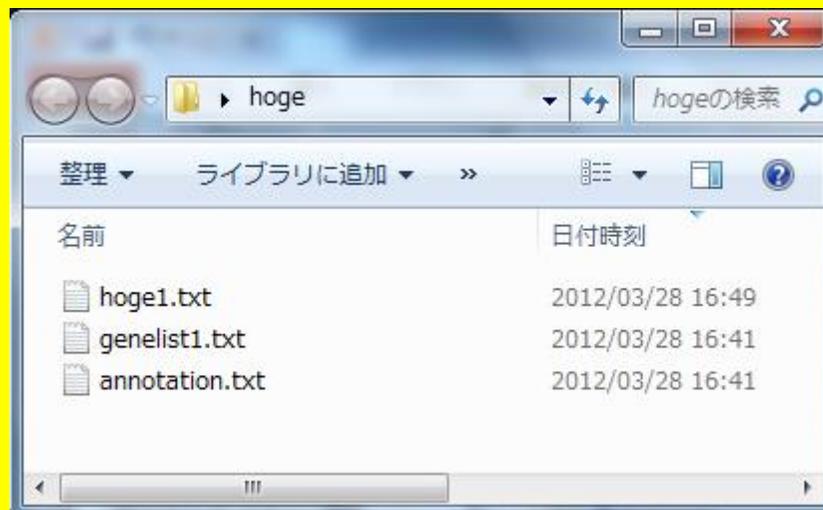
R Console
> in_f2 <- "genelist1.txt"
> out_f <- "hoge1.txt"
> param <- 1
>
> #ファイルの読み込み
> data <- read.table(in_f1, header=TRUE, sep="\t", quote="")
> keywords <- readLines(in_f2)
> dim(data)
[1] 11 4
>
> #本番
> obj <- is.element(as.character(data[,param]), keywords)
> out <- data[obj,]
> dim(out)
[1] 3 4
> write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F)
>
> |
```

	A	B	C	D
1	gene name	accession	description	subcellular_location
2	gene1	hoge01	plasma_mem	nuclear
3	gene7	hoge07	tebasaki	nuclear
4	gene9	hoge09	nihonshu	nuclear

実行前のhogeフォルダ



実行後のhogeフォルダ



色についての説明

(Rで)塩基配列解析(主に次世代シーケンサーのデータ) (last modified 2013/05/24, since 2010)

What's new?

RパッケージTCCの最新版は1.1.99です。6月6日に開催されるNAIST植物グローバル教育プロジェクト・平成25年度ワークショップでは、R(ver. 3.0.1)とTCC(ver. 1.1.99)をベースに話を予定しています。Rのバージョン自体は2.15.3でもたぶん大丈夫だとは思いますが。。。Rのインストールと起動を実行したあとにTCC(ver. 1.1.99)のインストールも忘れずに行っておいてください。Windowsのヒト用の詳細なインストール方法は[こちら](#)です。また、当日はhoge.zipという圧縮ファイルを解凍して得られる“hoge”というフォルダがデスクトップ上にあるという前提でセミナーを行いますのでこの圧縮ファイルもダウンロードと解凍をやっておいてください。(2013/05/24)NEW

- ・平成25年6月27日、7月3、4日にこのウェブページ関連の実習を含む講義(農学生命情報科学特論)を行います。東大生以外の外部の方も受講可能です。詳しくは[事務局](#)までお問い合わせください。(2013/05/23)NEW
- ・R3.0.1がリリースされていたのでこれに変更しました。(2013/05/17)NEW

	はじめに (last modified 2012/03/29)
	・ Rのインストールと起動 (last modified 2013/05/17) NEW
	・ サンプルデータ (last modified 2012/01/02)
	・ イントロダクション NGS

このページ内で用いる色についての説明:

コメント

特にやらなくてもいいコマンド

プログラム実行時に目的に応じて変更すべき箇所

色についての説明

このページ内で用いる色についての説明:

コメント

特にやらなくてもいいコマンド

プログラム実行時に目的に応じて変更すべき箇所

```
----- ここから -----  
in_f1 <- "annotation.txt"  
in_f2 <- "genelist1.txt"  
out_f <- "hoge1.txt"  
param <- 1  
  
#ファイルの読み込み  
data <- read.table(in_f1, header=TRUE, sep="\t", quote="")  
keywords <- readLines(in_f2)  
dim(data)
```

#本番

```
obj <- is.element(as.character(data[,param]), keywords)  
out <- data[obj,]  
dim(out)  
write.csv(out, "out.csv")  
-----
```

#入力ファイル名(目的のタブ区切りテキストファイル)
#入力ファイル名(キーワードなどのリストファイル)
#出力ファイル名を指定
#in_f1で読み込む目的のファイルの何列目のデータに

#入力ファイル(目的のファイル)を読み込んでdataに
#入力ファイル(リストファイル)を読み込んでkeyword
#オブジェクトdataの行数と列数を表示

#in_f1で読み込んだファイル中の(param)列目の文字
#行列dataからobjがTRUEとなる行のみを抽出した結果
#行数と列数を表示

上記は1列目でキーワード検索する場合

	A	B	C	D
1	gene name	accession	description	subcellular_location
2	gene1	hoge01	plasma_mem	nuclear
3	gene2	hoge02	hohinu	membrane
4	gene3	hoge03	agribio	endoplasmic
5	gene4	hoge04	genesis	endoplasmic
6	gene5	hoge05	kamo	membrane
7	gene6	hoge06	netteba	humei
8	gene7	hoge07	tebasaki	nuclear
9	gene8	hoge08	biiru	nuclear
10	gene9	hoge09	nihonshu	nuclear
11	gene10	hoge10	agene1	membrane
12	gene11	hoge11	iyaaaa	endoplasmic



4列目でキーワード検索したいときは?

	A	B	C	D
1	gene name	accession	description	subcellular_location
2	gene1	hoge01	plasma_mem	nuclear
3	gene2	hoge02	hohinu	membrane
4	gene3	hoge03	agribio	endoplasmic
5	gene4	hoge04	genesis	endoplasmic
6	gene5	hoge05	kamo	membrane
7	gene6	hoge06	netteba	humei
8	gene7	hoge07	tebasaki	nuclear
9	gene8	hoge08	biiru	nuclear
10	gene9	hoge09	nihonshu	nuclear
11	gene10	hoge10	agene1	membrane
12	gene11	hoge11	iyaaaa	endoplasmic

保存。

解答例

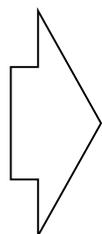
1. 目的のキーワードリストを含むファイルを作成し、例: list.txt
2. 該当箇所を変更し、Rコンソール画面上でコピー

```
list.txt - メモ帳
ファイル(F) 編集(E)
nuclear
membrane
```

```
run1.txt - メモ帳
ファイル(F) 編集(E) 書式(O) 表
in_f1 <- "annotation.txt"
in_f2 <- "genelist1.txt"
out_f <- "hogel.txt"
param <- 1

#ファイルの読み込み
data <- read.table(in_f1,
keywords <- readLines(in_f
dim(data)

#本番
obj <- is.element(as.chara
out <- data[obj,]
dim(out)
write.table(out, out_f, se
```



```
run1.txt - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
in_f1 <- "annotation.txt"
in_f2 <- "list.txt"
out_f <- "hogel.txt"
param <- 4

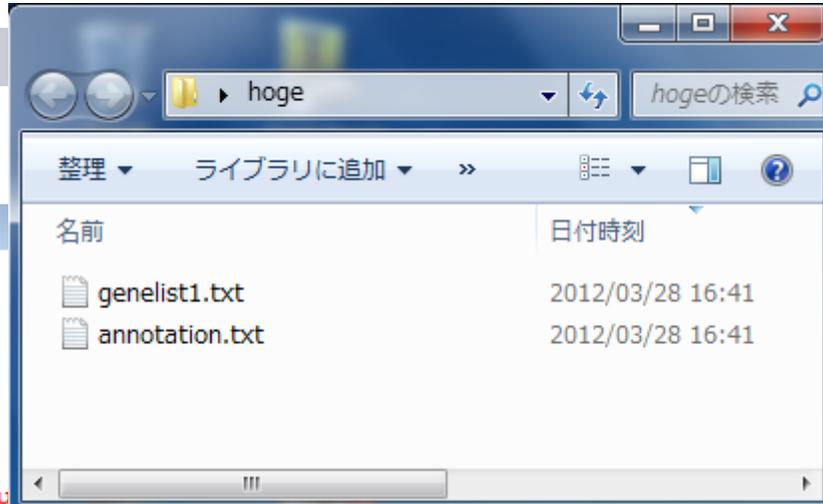
#ファイルの読み込み
data <- read.table(in_f1, header=TRUE, sep="¥t", quote="")
keywords <- readLines(in_f2)
dim(data)

#本番
obj <- is.element(as.character(data[,param]), keywords)
out <- data[obj,]
dim(out)
write.table(out, out_f, sep="¥t", append=F, quote=F, row.names=
```

一連の作業手順を記述した「スクリプト」を一つのファイルとして保存することをお勧め。(list.txtファイル作成時に、“membrane”と打った後に改行を入れた場合と入れない場合を試してみよう)

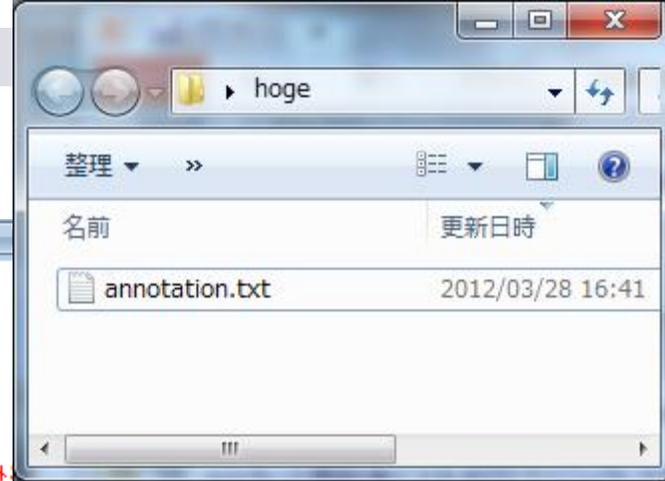
ありがちなミス1

```
R Console
> in_f1 <- "annotation.txt"
> in_f2 <- "genelist1.txt"
> out_f <- "hoge1.txt"
> param <- 1
>
> #ファイルの読み込み
> data <- read.table(in_f1, header=TRUE, sep="\t", qu
以下にエラー file(file, "rt") : コネクションを開くことができません
追加情報: 警告メッセージ:
In file(file, "rt") :
  ファイル 'annotation.txt' を開くことができません: No such file or director$
> keywords <- readLines(in_f2) #入力ファ$
以下にエラー file(con, "r") : コネクションを開くことができません
追加情報: 警告メッセージ:
In file(con, "r") :
  ファイル 'genelist1.txt' を開くことができません: No such file or directory
> dim(data) #オブジ$
NULL
>
> #本番
> obj <- is.element(as.character(data[,param]), keywords) #in_f1で読$
以下にエラー data[, param] :
  'closure' 型のオブジェクトは部分代入可能ではありません
> out <- data[obj,] #行列data$
エラー: オブジェクト 'obj' がありません
> dim(out) #オブジ$
エラー: オブジェクト 'out' がありません
> write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F) #outの中身$
以下にエラー is.data.frame(x) : オブジェクト 'out' がありません
>
> getwd()
[1] "C:/Users/kadota/Documents"
```



作業ディレクトリの変更を忘れている...

ありがちなミス2



```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge"
> in_f1 <- "annotation.txt"
> in_f2 <- "genelist1.txt"
> out_f <- "hogel.txt"
> param <- 1
>
> #ファイルの読み込み
> data <- read.table(in_f1, header=TRUE, sep="\t", quote="")
> keywords <- readLines(in_f2)
以下にエラー file(con, "r") : コネクションを開くことができません
追加情報: 警告メッセージ:
In file(con, "r") :
  ファイル 'genelist1.txt' を開くことができません: No such file or directory
> dim(data)
[1] 11 4
>
> #本番
> obj <- is.element(as.character(data[,param]), keywords)
以下にエラー match(el, set, 0L) : オブジェクト 'keywords' がありません
> out <- data[obj,]
以下にエラー `[.data.frame' (data, obj, ) : オブジェクト 'obj' がありません
> dim(out)
エラー: オブジェクト 'out' がありません
> write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F) #outの中身をo$
以下にエラー is.data.frame(x) : オブジェクト 'out' がありません
>
```

#入力ファイル\$
#入力ファイル\$
#出力ファイル\$
#in_f1で読み\$
#入力ファイル\$
#入力ファイル\$
#オブジェクト\$
#in_f1で読み\$
#行列dataから\$
#オブジェクト\$

必要な入力ファイルが作業ディレクトリ中に存在しない...

ありがちなミス3

```
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
> in_f1 <- "annotation.txt"
> in_f2 <- "list.txt"
> out_f <- "hoge1.txt"
> param <- 4
>
> #ファイルの読み込み
> data <- read.table(in_f1, header=TRUE, sep="\t", quote="")
> keywords <- readLines(in_f2)
> dim(data)
[1] 11 4
>
> #本番
> obj <- is.element(as.character(data[,param]), keywords)
> out <- data[obj,]
> dim(out)
[1] 7 4
> write.table(out, out_f, sep="\t", append=F, quote=F, row.names=
以下にエラー file(file, ifelse(append, "a", "w")) :
  コネクションを開くことができません
追加情報: 警告メッセージ:
In file(file, ifelse(append, "a", "w")) :
  ファイル 'hoge1.txt' を開くことができません: Permission denied
> |
```

#入

hoge1.txt - Microsoft Excel

	A	B	C	D	E	F
1	gene1	hoge01	plasma_mer	nuclear		
2	gene7	hoge07	tebasaki	nuclear		
3	gene9	hoge09	nihonshu	nuclear		
4						
5						
6						

```
run1.txt - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
in_f1 <- "annotation.txt"
in_f2 <- "list.txt"
out_f <- "hoge1.txt"
param <- 4

#ファイルの読み込み
data <- read.table(in_f1, header=TRUE, sep="\t", quote="")
keywords <- readLines(in_f2)
dim(data)

#本番
obj <- is.element(as.character(data[,param]), keywords)
out <- data[obj,]
dim(out)
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=
```

出力予定のファイル名と同じものを別のプログラムで開いているため最後のwrite.table関数のところでエラーが出る

ありがちなミス4

```
run1.txt - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
in_f1 <- "annotation.txt" #入力ファイル名(目的のタブ区切りテキストファイル)を
in_f2 <- "list.txt" #入力ファイル名(キーワードなどのリストファイル)を指定
out_f <- "hogel.txt" #出力ファイル名を指定
param <- 4 #in_f1で読み込む目的のファイルの何列目のデータに対し

#ファイルの読み込み
> data <- read.table(in_f1, header=TRUE, sep="\t", quote="") #入力ファイル(目的のファイル)を読み込んでdataに格納
> keywords <- readLines(in_f2) #入力ファイル(リストファイル)を読み込んでkeywordsに
> dim(data) #オブジェクトdataの行数と列数を表示

#本番
> obj <- is.element(as.character(data[,param]), keywords) #in_f1で読み込んだファイル中の(param)列目の文字列へ
> out <- data[obj,] #行列dataからobjがTRUEとなる行のみを抽出した結果をo
> dim(out) #オブジェクトoutの行数と列数を表示
> write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F) #outの中身をout_fで指定したファイル名で保存。

> keywords <- readLines(in_f2)
> dim(data)
[1] 11 4
> #本番
> obj <- is.element(as.character(data[,param]), keywords) #in_f1で読み込んだファイル中の(param)列目の文字列へ
> out <- data[obj,] #行列dataからobjがTRUEとなる行のみを抽出した結果をo
> dim(out) #オブジェクトoutの行数と列数を表示
[1] 7 4
> write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F) #outの中身をout_fで指定したファイル名で保存。 |
```

実行スクリプトをコピーする際、最後の行のところで改行を含ませずにR Console画面上でペーストしたため、最後のコマンドが実行されない(出力ファイルが生成されない)

• イントロダクション | 一般 | 任意のキーワードを含む行を抽出

例えばタブ区切りテキストファイルの `annotation.txt` が手元があり、この中から `genelist1.txt` のようなリストファイル中の文字列を含を示します。

Linux (UNIX) の `grep` コマンドのようなものです。perl のハッシュのようなものです。

「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ペ。

1. 目的のタブ区切りテキストファイル (`annotation.txt`) 中の第 `l` 列目をキーとして、リストファイル (`genelist1.txt`) 中のものが含ま

```
-----   ここから   -----
in_f1 <- "annotation.txt"
in_f2 <- "genelist1.txt"
out_f  <- "hogel.txt"
param <- 1
```

#入力ファイル名(目的のタブ区切りテキストファイル)を指定
#入力ファイル名(キーワードなどのリストファイル)を指定
#出力ファイル名を指定
#in_f1で読み込む目的のファイルの何列目のデータに対して

#ファイルの読み込み

```
data <- read.table(in_f1, header=TRUE, sep="\t", quote="")
keywords <- readLines(in_f2)
dim(data)
```

#入力ファイル(目的のファイル)を読み込んでdataに格納
#入力ファイル(リストファイル)を読み込んでkeywordsに格納
#オブジェクトdataの行数と列数を表示

#本番

```
obj <- is.element(as.character(data[,param]), keywords)
out <- data[obj,]
dim(out)
```

#in_f1で読み込んだファイル中の(param)列目の文字列ベクト
#行列dataからobjがTRUEとなる行のみを抽出した結果をoutに
#オブジェクトoutの行数と列数を表示

```
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F)
```

#outの中身をout_fで指定したファイル名で保存。

```
-----   ここまで   -----
```

「----- ここまで -----」の一つ上の空行には「スクリプト最終行のコマンドを確実に実行するため」という深い意味があります

読み込み

----- ここから -----

```
in_f1 <- "annotation.txt"
in_f2 <- "genelist1.txt"
out_f <- "hogel.txt"
param <- 1
```

#ファイルの読み込み

```
data <- read.table(in_f1, header=TRUE, sep="\t", quote="")
keywords <- readLines(in_f2)
dim(data)
```

#本番

```
obj <- is.element(as.character(1:10), rownames(data))
out <- data[obj,]
dim(out)
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F)
```

----- ここまで -----

	A	B	C	D
1	gene name	accession	description	subcellular_location
2	gene1	hoge01	plasma_mem	nuclear
3	gene2	hoge02	hohinu	membrane
4	gene3	hoge03	agribio	endoplasmic
5	gene4	hoge04	genesis	endoplasmic
6	gene5	hoge05	kamo	membrane
7	gene6	hoge06	netteba	humei
8	gene7	hoge07	tebasaki	nuclear
9	gene8	hoge08	biiru	nuclear
10	gene9	hoge09	nihonshu	nuclear
11	gene10	hoge10	agene1	membrane
12	gene11	hoge11	iyaaaa	endoplasmic

- ① in_f1で指定したファイルを読み込め
- ② 読み込むファイルの最初の行はヘッダー部分です
- ③ ファイルの区切り文字は「タブ」です

行列data

参考

	A	B	C	D
1	gene name	accession	description	subcellular_location
2	gene1	hoge01	plasma_mem	nuclear
3	gene2	hoge02	hohinu	membrane
4	gene3	hoge03	agribio	endoplasmic
5	gene4	hoge04	genesis	endoplasmic
6	gene5	hoge05	kamo	membrane
7	gene6	hoge06	netteba	humei
8	gene7	hoge07	tebasaki	nuclear
9	gene8	hoge08	biiru	nuclear
10	gene9	hoge09	nihonshu	nuclear
11	gene10	hoge10	agen1	membrane
12	gene11	hoge11	iyaaaa	endoplasmic

```
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
>
> #ファイルの読み込み
> data <- read.table(in_f1, header=TRUE, sep="\t", g
>
> data
  gene name accession description subcellular_location
1     gene1   hoge01   plasma_mem          nuclear
2     gene2   hoge02     hohinu            membrane
3     gene3   hoge03     agribio            endoplasmic
4     gene4   hoge04     genesis            endoplasmic
5     gene5   hoge05       kamo            membrane
6     gene6   hoge06     netteba            humei
7     gene7   hoge07     tebasaki           nuclear
8     gene8   hoge08       biiru            nuclear
9     gene9   hoge09     nihonshu           nuclear
10    gene10   hoge10     agen1             membrane
11    gene11   hoge11     iyaaaa            endoplasmic
> |
```

入力ファイルの中身を正しく読み込めていることがわかる

```

----- ここから -----
in_f1 <- "annotation.txt"
in_f2 <- "genelist1.txt"
out_f <- "hoge1.txt"
param <- 1

```

```

#ファイルの読み込み
data <- read.table(i
keywords <- readLine
dim(data)

```

```

#本番
obj <- is.element(as
out <- data[obj,]
dim(out)
write.table(out, out

```

----- ここまで -----

	A	B	C	D
1	gene name	accession	description	subcellular_location
2	gene1	hoge01	plasma_mem	nuclear
3	gene2	hoge02	hohinu	membrane
4	gene3	hoge03	agribio	endoplasmic
5	gene4	hoge04	genesis	endoplasmic
6	gene5	hoge05	kamo	membrane
7	gene6	hoge06	netteba	humei
8	gene7	hoge07	tebasaki	nuclear
9	gene8	hoge08	biiru	nuclear
10	gene9	hoge09	nihonshu	nuclear
11	gene10	hoge10	agen1	membrane
12	gene11	hoge11	iyaaaa	endoplasmic

```

RGui (64-bit)
ファイル 編集 閲覧
R Console
> data
  gene name accession plasma_mem nuclear
1 gene1 hoge01 plasma_mem nuclear
2 gene2 hoge02 hohinu membrane
3 gene3 hoge03 agribio endoplasmic
4 gene4 hoge04
5 gene5 hoge05
6 gene6 hoge06
7 gene7 hoge07
8 gene8 hoge08
9 gene9 hoge09 nihonshu nuclear
10 gene10 hoge10 agen1 membrane
11 gene11 hoge11 iyaaaa endoplasmic
> dim(data)
[1] 11 4
> |

```

オブジェクトdataの行数と列数は11と4。webpage中の表記が灰色なのは、「特にやらなくてもいいコマンド」だから。

行列の要素へのアクセス

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

```

10  gene10  hoge10  agene1  memk
11  gene11  hoge11  iyaaaa  endopla
> dim(data)
[1] 11 4
> data[6,4]
[1] humei
Levels: endoplasmic humei membrane nuclear
> data[2,]
  gene10  hoge10  agene1  memk
  gene11  hoge11  iyaaaa  endopla
  gene2    hoge2    hohinu   membrane
  gene3    hoge3    agribio  endoplasmic
  gene4    hoge4    genesis  endoplasmic
  gene5    hoge5    kamo     membrane
  gene6    hoge6    netteba  humei
  gene7    hoge7    tebasaki nuclear
  gene8    hoge8    biiru    nuclear
  gene9    hoge9    nihonshu nuclear
  gene10   hoge10   agene1   membrane
  gene11   hoge11   iyaaaa   endoplasmic
> data[,2]
[1] hoge01 hoge02 hoge03 hoge04 hoge05 hoge06 hoge07
[8] hoge08 hoge09 hoge10 hoge11
11 Levels: hoge01 hoge02 hoge03 hoge04 hoge05 ... hoge11
> data[,param]
[1] gene1  gene2  gene3  gene4  gene5  gene6  gene7
[8] gene8  gene9  gene10 gene11
11 Levels: gene1 gene10 gene11 gene2 gene3 ... gene9
> |

```

data[行, 列]

	A	B	C	D
1	gene10	hoge10	agene1	membrane
2	gene11	hoge11	iyaaaa	endoplasmic
3	gene2	hoge02	hohinu	membrane
4	gene3	hoge03	agribio	endoplasmic
5	gene4	hoge04	genesis	endoplasmic
6	gene5	hoge05	kamo	membrane
7	gene6	hoge06	netteba	humei
8	gene7	hoge07	tebasaki	nuclear
9	gene8	hoge08	biiru	nuclear
10	gene9	hoge09	nihonshu	nuclear
11	gene10	hoge10	agene1	membrane
12	gene11	hoge11	iyaaaa	endoplasmic

```

----- ここから -----
in_f1 <- "annotation.txt"
in_f2 <- "genelist1.txt"
out_f  <- "hoge1.txt"
param <- 1

```

**paramには1という数値
が代入されていたから**

	A
1	gene1
2	gene7
3	gene9

やりたかったことをおさらい

```
----- ここから -----
in_f1 <- "annotation.txt"
in_f2 <- "genelist1.txt"
out_f <- "hoge1.txt"
param <- 1
```

```
#ファイルの読み込み
data <- read.table(in_f2)
keywords <- readLines(in_f1)
dim(data)
```

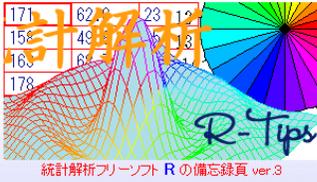
```
#本番
obj <- is.element(keywords, data[,1])
out <- data[obj,]
dim(out)
write.table(out, out_f)
```

```
----- ここまで -----
```

```
R Console
> obj
[1] TRUE FALSE FALSE FALSE FALSE FALSE TRUE FALSE TRUE FALSE FALSE
> data
  gene_name accession description subcellular_location
1   gene1   hoge01  plasma_mem          nuclear
2   gene2   hoge02    hohinu            membrane
3   gene3   hoge03   agribio          endoplasmic
4   gene4   hoge04   genesis          endoplasmic
5   gene5   hoge05     kamo            membrane
6   gene6   hoge06   netteba           humei
7   gene7   hoge07   tebasaki          nuclear
8   gene8   hoge08    hiiru            nuclear
9   gene9   hoge09   nihonshu          nuclear
10  gene10   hoge10
11  gene11   hoge11
> obj
[1] TRUE FALSE FALSE FALSE FALSE FALSE TRUE FALSE TRUE FALSE FALSE
> data[obj,]
  gene_name accession description subcellular_location
1   gene1   hoge01  plasma_mem          nuclear
7   gene7   hoge07   tebasaki          nuclear
9   gene9   hoge09   nihonshu          nuclear
>
```

論理値ベクトルobjを用いてTRUEの要素に対応する行を抽出している

R-Tipsでお勉強



ときどき「R-Tips」に立ち返るといいと思います

ブラウザのアドレスバー: <http://cse.naro.affrc.go.jp/takezawa/r>

フレーム: なしあり

R (有名な統計言語『S』を
まなプラットフォーム(OS)に
にも関わらず、世界中の専門
加えられていますので、機能
りません。とにかく計算が速
です。このドキュメントは Win
へた足跡です。

- リンク [R-Web\(ウェブ版でん\)](#)
- リンク [RjpWiki\(岡田先生\)](#)
- リンク [PDF版 R-Tips\(2008\)](#)
- 第01節 [R のセットアップ+](#)
- 第03節 [簡単な計算](#)
- 第05節 [オブジェクトと代入](#)
- 第07節 [ヘルプを見る](#)
- 第09節 [データの型](#)

ベクトル篇

第12節	ベクトルの作成	第13節	要素へのアクセス
第14節	ベクトルの計算	第15節	ベクトル要素の置換
第16節	種々のベクトル	第17節	文字列を操作する
第18節	NULL, NA, NaN, Infの操作		

行列・配列・リスト

第19節	行列の作成	第20節	
第21節	行列の操作	第22節	
第23節	リスト	第24節	
第25節	データ型とデータ構造	第26節	

関数とプログラミング

第27節	関数事始	第28節	
第29節	条件分岐	第30節	
第31節	関数の定義	第32節	
第33節	引数について	第34節	

ベクトル要素へのアクセス

ベクトルの中の数を「要素」と呼び、各要素には左から順に 1, 2, ... と番号が振られている。以下ではベクトル x について要素にアクセスする方法を一覧表で示している。

コマンド	機能
x[k]	k 番目の要素を取り出す。要素番号として 0 を指定すると、長さ 0 で元のベクトルと同じ型のベクトルが返る。
x[k] <- a	k 番目の要素を a に変更。
x[正整数ベクトル]	いくつかの要素をまとめて取り出す。
x[負整数ベクトル]	対応する要素番号の要素を取り除く。
x[論理値ベクトル]	TRUE の要素に対応した要素を取り出す。
x[条件式]	条件に合致した要素を取り出す。
x[文字型ベクトル]	要素ラベルを指定して要素を取り出す (names 属性が付いている場合)。

以下に簡単な例を示す。

```
x <- c(1, 2, 3, 4, 5)
x[3]
[1] 3
```

3 番目の要素を取り出す

論理値ベクトルを理解

genelist1.txt

	A
1	gene1
2	gene7
3	gene9

```
----- ここから -----
in_f1 <- "annotation.txt"
in_f2 <- "genelist1.txt"
out_f <- "hogel.txt"
param <- 1
```

```
#ファイルの読み込み
data <- read.table(in_f1, header=T)
keywords <- readLines(in_f2)
dim(data)
```

#本番

```
obj <- is.element(as.character(data[,param]), keywords)
out <- data[obj,]
dim(out)
write.table(out, out_f, sep="%t", append=F, quote=F, row.names=F)
```

```
----- ここまで -----
```

```
R Console
>
> keywords #keywordsの中身を表示
[1] "gene1" "gene7" "gene9"
> length(keywords) #keywordsベクトルの要素数を表示
[1] 3
> keywords[3] #3番目の要素を表示
[1] "gene9"
> keywords[4] #4番目の要素は...ない
[1] NA
> keywords == "gene7" #"gene7"という文字の位置情報を表示
[1] FALSE TRUE FALSE
> obj <- keywords == "gene7" #上記結果をobjに格納
> keywords[obj] #objがTRUEとなる要素のみを表示
[1] "gene7"
> |
```

論理値ベクトルを理解

genelist1.txt

	A
1	gene1
2	gene7
3	gene9

R Console

```
> hoge <- c("gene7", "gene9") #二つの要素からなるベクトルhogeを作成
> hoge #hogeの中身を表示させてるだけ
[1] "gene7" "gene9"
> keywords == hoge #FALSE TRUE TRUEになることを期待したが...
[1] FALSE FALSE FALSE
警告メッセージ:
In keywords == hoge :
長いオブジェクトの長さが短いオブジェクトの長さの倍数になっていません
> is.element(keywords, hoge) #keywords中の各要素は集合hogeに含まれるか否か
[1] FALSE TRUE TRUE
> |
```

疑問に思ったら、自分の理解できるところからボトムアップ式に試す

#本番

```
obj <- is.element(as.character(data[,param]), keywords)
out <- data[obj,]
dim(out)
write.table(out, out_f, sep="%t", append=F, quote=F, row.names=F)
```

----- ここまで -----

Contents

■ Rでゲノム解析

- アノテーションファイル(行列形式)からの情報抽出
 - 条件を満たす行のみの抽出。ハッシュとかgrep周辺がRで可能
- multi-fastaファイルからの情報抽出
 - 一定の長さ以上のものを抽出。コンティグごとのGC含量計算。

■ Rで(比較)トランスクリプトーム解析

- 研究目的別留意点(サンプル内比較、サンプル間比較)
- リードマッピング → 数値化(カウントデータの取得)
- 発現変動解析
 - 既存のRパッケージの弱点(unbiased DE: ○、biased DE: ×)
 - TCCパッケージを用いて二群間比較(複製あり、複製なし)



multi-fastaファイルからの情報抽出

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 2012/05/28)
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 2012/05/28)
イントロダクション	一般	任意の長さの連続塩基の出現頻度情報を取得 (last modified 2012/05/28)
イントロダクション	NGS	NGSデータ取得 (last modified 2011/07/19)
イントロダクション	NGS	NGSデータ取得(SRAdb) (last modified 2011/01/14)
イントロダクション	NGS	マッピング basic aligner (基本的なマッピングプログラム) (last modified 2012/09/11)
イントロダクション	NGS	マッピング splice-aware aligner (spliced readもマッピング可能なもの) (last modified 2012/08/01)
イントロダクション	NGS	マッピング (short) readの入力形式について (last modified 2011/08/03)
イントロダクション	NGS	マッピング (short) readの出力形式について (last modified 2010/12/06)
イントロダクション	NGS	マッピング (ESTレベルの長さの)contig (last modified 2010/12/06)
イントロダクション	NGS	ファイル形式の変換 (last modified 2011/10/31)
イントロダクション	NGS	アセンブルプログラムについて (last modified 2011/07/26)
イントロダクション	NGS	アセンブルプログラム(ゲノム) (last modified 2012/08/01)
イントロダクション	NGS	アセンブルプログラム(転写物) (last modified 2012/08/30)
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 (last modified 2012/04/03)
イントロダクション	NGS	可視化ツール(R用) (last modified 2010/7/8)
イントロダクション	NGS	可視化ツール(R以外) (last modified 2010/7/8)
イントロダクション	NGS	NGSとqPCRやmicroarrayなどとの比較 (last modified 2010/12/16)
ファイルの読み込み	マップ前	FASTQ形式 (last modified 2013/01/16)
ファイルの読み込み	マップ前	FASTA形式 (last modified 2010/6/18)
ファイルの読み込み	マップ前	Illuminaの* seq.txt (last modified 2010/6/3)
ファイルの読み込み	マップ前	Illuminaの* gseq.txt (last modified 2012/08/06)
ファイルの読み込み	マップ後	Bowtie形式 (last modified 2013/01/16)

Rでmulti-fastaファイルを読み込んで自由に解析できます

まずはコピペ

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

アセンブルプログラムを実行して得られるmulti-fastaファイルを読み込んで、Total lengthやaverage lengthのやり方を示します。ここでは以下の二つを例題として行います:

- 1 「イントロダクション | 一般 | ランダムな塩基配列を作成」の4を実行して得られたmulti-fastaファイル
 2. 130MB程度のRefSeqのhuman mRNAのmulti-fastaファイル(h_rna.fasta)
- 「ファイル」-「ディレクトリの変更」でファイルを保存したいディレクトリに移動し、以下をコピー

1. hoge4.faファイルの場合:

```
----- ここから -----
in_f <- "hoge4.fa"
out_f <- "hoge1.txt"
```

#必要なパッケージなどをロード

```
library(Biostrings)
```

```
reads <- readDNASTringSet(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計算のところで

```
sorted_length <- rev(sort(width(reads)))
N50 <- sorted_length[cumsum(sorted_length) >= Total_length][1]
```

#GC含量(GC content)計算のところで

```
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))
```

```
hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)

>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCGAGAGACCCAGCACACCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
```

```
RGui (64-bit)
ファイル 編集 閲覧 その他
[Icons] [STOP]

R Console
> #GC含量(GC content)計算のところ
> count <- alphabetFrequency(reads) #A,C,G,$
> CG <- rowSums(count[,2:3]) #C,Gの$
> ACGT <- rowSums(count[,1:4]) #A,C,G,$
> GC_content <- sum(CG)/sum(ACGT) #トータ$
>
> #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$
>
> |
```

「hoge」フォルダにhoge1.txtが作成されているはず

結果ファイルを眺めて動作確認

ID	Length
contig_1	24
contig_2	103
contig_3	65
contig_4	49

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

アセンブルプログラムを実行して得られるmulti-fastaファイルを読み込んで、Total lengthやaverage lengthなどの各種情報取得を行うためのやり方を示します。ここでは以下の二つを例題として行います:

- 1 「イントロダクション | 一般 | ランダムな塩基配列を作成」の4を実行して得られたmulti-fastaファイル(hoge4.fa)
 2. 130MB程度のRefSeqのhuman mRNAのmulti-fastaファイル(h_rna.fasta)
- 「ファイル」-「ディレクトリの変更」でファイルを保存したいディレクトリに移動し以下をコピー。

1. hoge4.faファイルの場合:

```
----- ここから -----
in_f <- "hoge4.fa" #multi-fasta形式の入力ファイルを指定
out_f <- "hoge1.txt" #出力ファイル名を指定
```

```
#必要なパッケージなどをロード
hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACACCCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
```

ファイルをFASTA形式で読み込み
 配列の長さ」の情報を取得
 情報の取得
 長」の情報を取得
 値」の情報を取得
 「最大値」の情報を取得
 「最小値」の情報を取得
 ートした結果をsorted length
 ら足しこんでいってTotal leng
 配列ごとにカウントした結果を
 CGGに格納
 算してACGTに格納
 情報を取得

```
<- 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))
```

	A	B
1	Total length (bp)	241
2	Number of contigs	4
3	Average length	60.25
4	Median length	57
5	Max length	103
6	Min length	24
7	N50	65
8	GC content	0.577

N50

■ アセンブルがどれだけうまくいっているかを表す指標の一つ

- 長いコンティグから足していったTotal_lengthの50%に達したときのコンティグの長さ

ID	Length
contig_1	24
contig_2	103
contig_3	65
contig_4	49



情報抽出手順(の一部)

1. hoge4.faファイルの場合:

```
----- ここから -----
in_f <- "hoge4.fa"
out_f <- "hoge1.txt"
```

```
#必要なパッケージなどをロード
library(Biostrings)
```

```
reads <- readDNASTringSet(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))
```

```
#multi-fasta形式の入力ファイルを指定
#出力ファイル名を指定
```

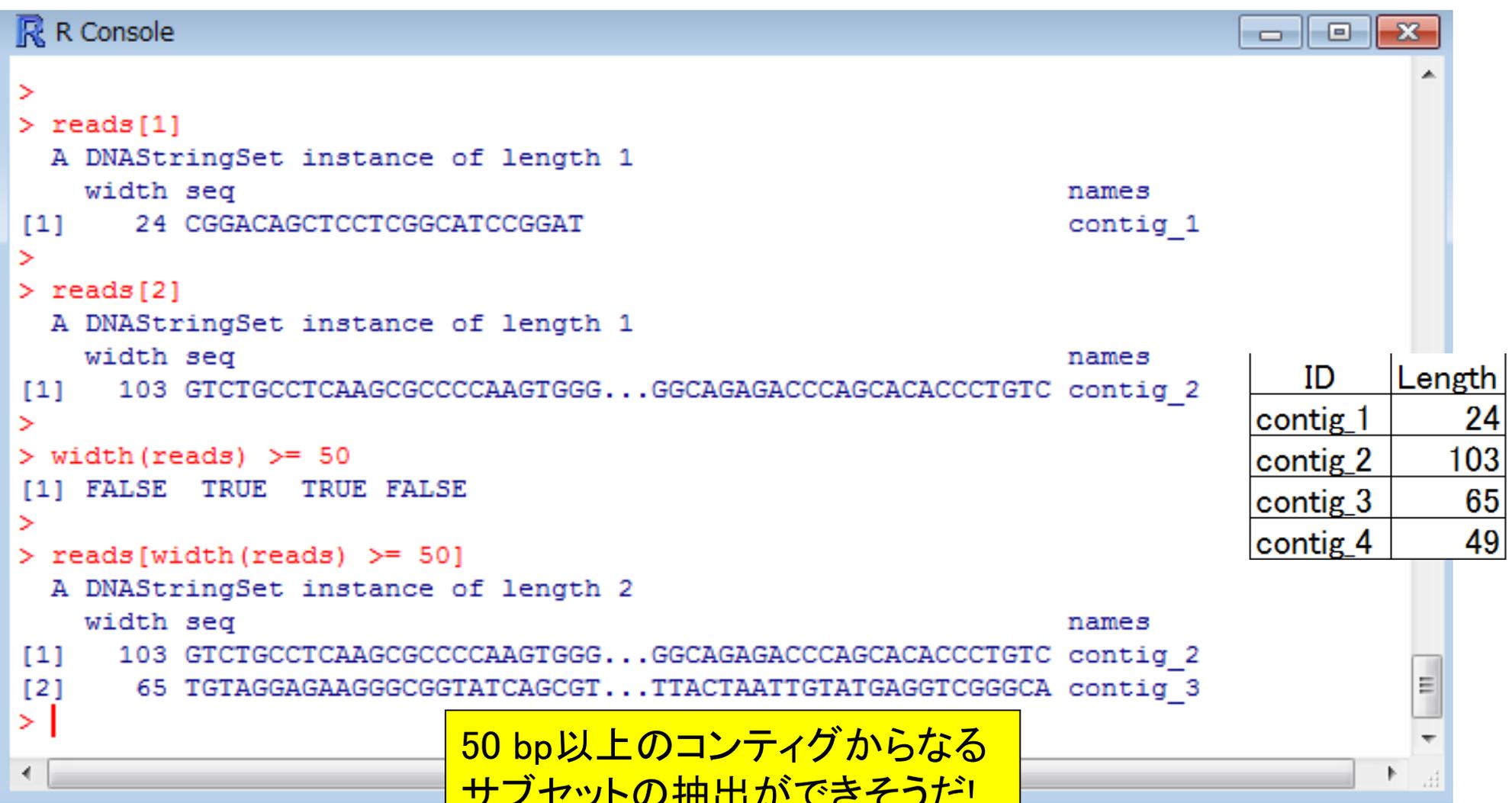
```
#パッケージの読み込み
```

```
#in_fで指定したファイルをFASTA形式で読み込み
```

```
R Console
> reads
A DNASTringSet instance of length 4
  width seq                                     names $
[1]   24 CGGACAGCTCCTCGGCATCCGGAT             contig_1
[2]  103 GTCTGCCTCAAGCGCCCAAG...GAGACCCAGCACACCCTGTC contig_2
[3]   65 TGTAGGAGAAGGGCGGTATCA...TAATTGTATGAGGTCGGGCA contig_3
[4]   49 CGTGCTGATTCCACACAGCAG...GACCTCTACCTATGAACATG contig_4
>
> width(reads)
[1] 24 103 65 49
>
> sum(width(reads))
[1] 241
>
```

width関数を使えば配列長
情報を取り出せるようだ

情報抽出手順(の一部)



```
>
> reads[1]
A DNASTringSet instance of length 1
width seq
[1] 24 CGGACAGCTCCTCGGCATCCGGAT
names
contig_1
>
> reads[2]
A DNASTringSet instance of length 1
width seq
[1] 103 GTCTGCCTCAAGCGCCCAAGTGGG...GGCAGAGACCCAGCACACCCTGTC
names
contig_2
>
> width(reads) >= 50
[1] FALSE TRUE TRUE FALSE
>
> reads[width(reads) >= 50]
A DNASTringSet instance of length 2
width seq
[1] 103 GTCTGCCTCAAGCGCCCAAGTGGG...GGCAGAGACCCAGCACACCCTGTC
[2] 65 TGTAGGAGAAGGGCGGTATCAGCGT...TTACTAATTGTATGAGGTCGGGCA
names
contig_2
contig_3
> |
```

ID	Length
contig_1	24
contig_2	103
contig_3	65
contig_4	49

50 bp以上のコンティグからなるサブセットの抽出ができそうだ!

- ファイルの読み込み マップ後 [Bowtie形式](#) (last modified 2013/01/16)
- ファイルの読み込み マップ後 [BAM形式](#) (last modified 2010/7/13)
- ファイルの読み込み マップ後 [GFF3形式](#) (last modified 2010/6/9)
- ファイルの読み込み マップ後 [SOAP形式](#) (last modified 2011/07/20)
- [クオリティチェック](#) | [NGS\(一般\)について](#) (last modified 2012/08/02)
- [クオリティチェック](#) | [NGS\(一般\) | qrc \(Quick Read Quality Control\)](#) (last modified 2012/02/20)
- [フィルタリング](#) | 一般 [fastqファイルからACGTのみからなる配列\(重複あり\)の抽出](#) (last modified 2012/07/13)
- [フィルタリング](#) | 一般 [fastqファイルからACGTのみからなる配列\(重複なし\)の抽出](#) (last modified 2012/07/13)
- [フィルタリング](#) | 一般 [multi-fastaファイルのdescription行の記述を整形](#) (last modified 2012/07/13)
- [フィルタリング](#) | 一般 [multi-fastaファイルからACGT以外の文字の出現数でフィルタリング](#) (last modified 2012/07/13)
- [フィルタリング](#) | 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 2012/04/03)
- [フィルタリング](#) | 一般 [multi-fastaファイルから任意のサブセットを抽出](#) (last modified 2012/07/17)
- [フィルタリング](#) | [NGS\(miRNA\) | アダプター配列除去0\(基本的なところ\)](#) (last modified 2010/5/27)
- [フィルタリング](#) | 一般 | [multi-fastaファイルから指定した配列長のもののみ抽出](#)

multi-fasta形式のファイルが手元にあったとして、任意の配列長のもののみ抽出するやり方を示します。
 ここでは以下の二つを例題として行います：

1. 「イントロダクション | 一般 | [ランダムな塩基配列を作成](#)」の4.を実行して得られたmulti-fastaファイル([hoge4.fa](#))
2. 130MB程度のRefSeqのhuman mRNAのmulti-fastaファイル([h_rna.fasta](#))

「ファイル」-「ディレクトリの変更」でファイルを保存したいディレクトリに移動し以下をコピー。

1. [hoge4.fa](#)ファイルの場合：

```

-----   ここから   -----
in_f <- "hoge4.fa"
out_f <- "hoge1.fasta"
param <- 50

#必要なパッケージなどをロード
library(Biostrings)

reads <- readDNASTringSet(in_f, format="fasta")
reads
reads <- reads[width(reads) >= param]
reads
writeXStringSet(reads, file=out_f, format="fasta")

-----   ここまで   -----
    
```

#multi-fasta形式の入力ファイルを指定
 #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read
 #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

指定した配列長以下のものを抽出
 したいときは「<=」とすればよい

- インタロダクション | NGS | [アンテーション情報取得\(BioMart and biomaRt\)](#) (last modified 2011/08/26)
- インタロダクション | 一般 | [任意のキーワードを含む行を抽出](#) (last modified 2012/03/29) NEW
- インタロダクション | 一般 | [ランダムな塩基配列を作成](#) (last modified 2012/04/02) NEW
- インタロダクション | 一般 | [配列取得](#) (last modified 2010/7/7)
- インタロダクション | 一般 | [指定した範囲の配列を取得](#) (last modified 2012/04/03) NEW
- インタロダクション | 一般 | [翻訳配列\(translate\)を取得](#) (last modified 2011/07/27)

● インタロダクション | 一般 | 指定した範囲の配列を取得

1.では、12塩基(AGTGACGGTCTT)からなる一つの塩基配列(description行が">kadota")からなるFASTA形式ファイル(sample1.fasta)を入力として、この塩基配列の任意の範囲(始点が3, 終点が9)の配列を抽出し、得られた部分配列をFASTA形式ファイル(tmp1.fasta)に出力するやり方を示します。

2.では、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RefSeq ID(例:NM203348.1)の任意の範囲(例:始点が2, 終点が5)の配列の抽出を行います。

3.では、2.の延長線上で、目的のaccession番号が複数ある場合に対応したものです。start位置, 3列目:end位置」からなるリストファイル(list_sub1.txt)を読み込ませて、目的の方を示します。

ここでは、得られた部分配列をFASTA形式ファイル(ファイル名:tmp3.fasta)で保存する

4.では、3.と同じことを異なるファイル(multi-fastaファイル:hoge4.fa、リストファイル:list

入力ファイル: sample1.fasta

```
>kadota
AGTGACGGTCTT
```

2列目:
るやり

出力ファイル: tmp1.fasta

```
>kadota
TGACGGT
```

```
1. sample1.fastaファイルの場合:
----- ここから -----
in_f <- "sample1.fasta"
out_f <- "tmp1.fasta"
param <- c(3, 9)

#必要なパッケージなどをロード
library(Biostrings)

reads <- readDNASTringSet(in_f, format="fasta")
out <- subseq(reads, param[1], param[2])
writeXStringSet(out, file=out_f, format="fasta", width=80)

----- ここまで -----
```

```
#multi-fasta形式のフ
#出力ファイル名を指定
#抽出したい範囲の始点
#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
#paramで指定した始点と終点の範囲の配列を抽出してoutに格納
#outの中身をout_fで指定したファイル名で保存
```

1. `sample1.fasta` ファイルの場合:

```
----- ここから -----
in_f <- "sample1.fasta"
out_f <- "tmp1.fasta"
param <- c(3, 9)
```

#必要なパッケージなどをロード
library(Biostrings)

```
reads <- readDNASTringSet(in_f, format="fasta")
out <- subseq(reads, param[1], param[2])
writeXStringSet(out, file=out_f, format="fasta", width=80)
```

----- ここまで

```
> reads
  A DNASTringSet instance of length 1
    width seq
[1]    12 AGTGACGGTCTT

> out
  A DNASTringSet instance of length 1
    width seq
[1]     7 TGACGGT

> param
[1] 3 9

> param[1]
[1] 3

> param[2]
[1] 9

> subseq(reads, 3, 9)
  A DNASTringSet instance of length 1
    width seq
[1]     7 TGACGGT

> |
```

入力ファイル: sample1.fasta

```
>kadota
AGTGACGGTCTT
```

出力ファイル: tmp1.fasta

```
>kadota
TGACGGT
```

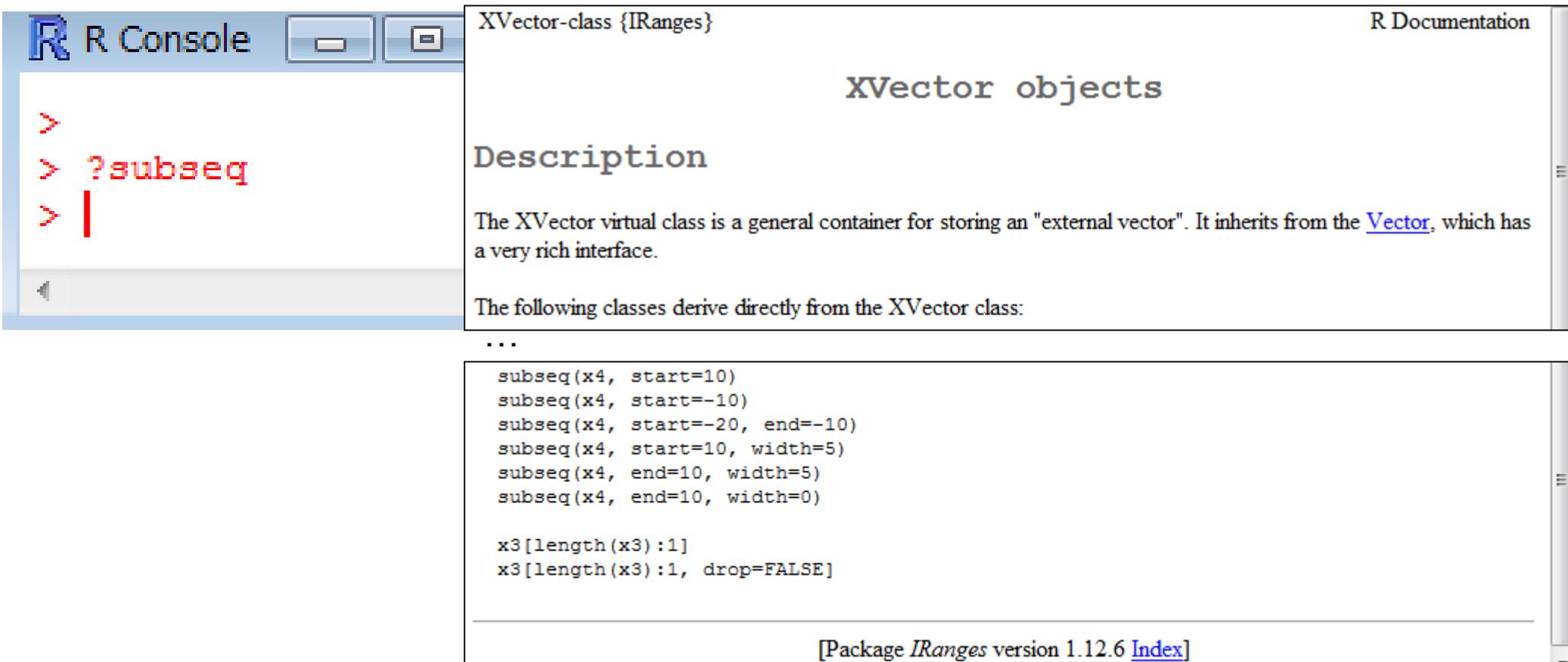
```
names
kadota
```

```
names
kadota
```

```
names
kadota
```

subseq関数は「塩基配列, start, end」という形式で使うようだ

関数の使用法について



The image shows a screenshot of an R environment. On the left is the R Console window with the following input:

```
>  
> ?subseq  
> |
```

On the right is the R Documentation window for the XVector class from the IRanges package. The title is "XVector objects". The description states: "The XVector virtual class is a general container for storing an 'external vector'. It inherits from the [Vector](#), which has a very rich interface." Below the description, it lists classes that derive directly from the XVector class:

```
...  
subseq(x4, start=10)  
subseq(x4, start=-10)  
subseq(x4, start=-20, end=-10)  
subseq(x4, start=10, width=5)  
subseq(x4, end=10, width=5)  
subseq(x4, end=10, width=0)  
  
x3[length(x3):1]  
x3[length(x3):1, drop=FALSE]
```

At the bottom of the documentation window, it says "[Package IRanges version 1.12.6 [Index](#)]"

- ・「?関数名」で使用方法を記したウェブページが開く
- ・ページの下のように、(大抵の場合)使用例が掲載されている
- ・使用方法既知の関数のマニュアルをいくつか読んで、慣れておく

```
R R Console  
  
> reads  
A DNAStringSet instance of length 1  
width seq names  
[1] 12 AGTGACGGTCTT kadota  
  
> subseq(reads, start=3, end=9)  
A DNAStringSet instance of length 1  
width seq names  
[1] 7 TGACGGT kadota  
  
> subseq(reads, start=3, width=9)  
A DNAStringSet instance of length 1  
width seq names  
[1] 9 TGACGGTCT kadota  
  
> subseq(reads, start=3, width=11)  
以下にエラー .Call2("solve_user_SEW", refwidths, start, end  
solving row 1: 'allow.nonnarrowing' is FALSE and the so  
  
> subseq(reads, start=3, width=10)  
A DNAStringSet instance of length 1  
width seq names  
[1] 10 TGACGGTCTT kadota  
  
> |
```

入力ファイル: sample1.fasta

```
>kadota  
AGTGACGGTCTT
```

出力ファイル: tmp1.fasta

```
>kadota  
TGACGGT
```

マニュアルの使用例をいくつか試して、ステップアップ

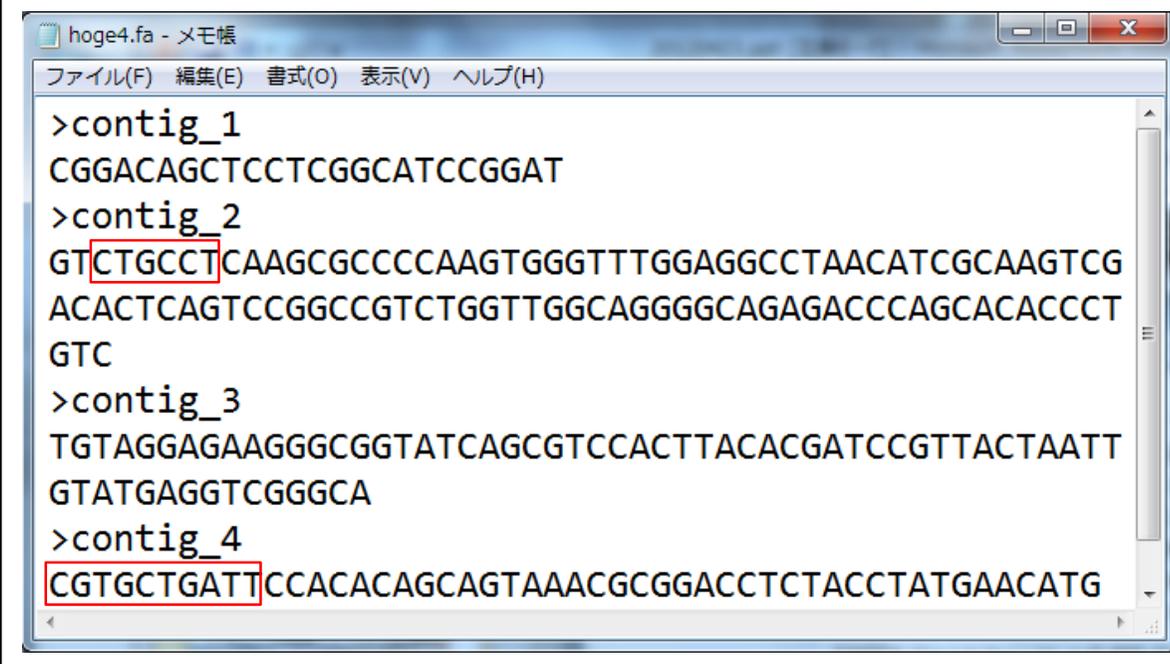
4. hoge4.faファイルで抽出したい情報をリストファイル (list_sub2.txt) で与える場合:

```
----- ここから -----
in_f1 <- "hoge4.fa"
in_f2 <- "list_sub2.txt"
out_f <- "tmp4.fasta"
```

#multi-fasta形式のファイルを指定
#リストファイルを指定
#出力ファイル名を指定

```
library(Biostrings)
reads <- readFasta(in_f1)
list_obj <- readFastaList(in_f2)

out <- NULL
for(i in 1:length(list_obj)) {
  obj <- list_obj[[i]]
  out <- append(out, obj)
}
write.XStringSet(out, out_f)
```



list_sub2.txt

contig_4	1	10
contig_2	3	8

式で読み込み
スペースホル
一致する位置を
てobjがTRUE
名で保存

出力ファイル: tmp4.fasta

```
>contig_4
CGTGCTGATT
>contig_2
CTGCCT
```

Contents

■ Rでゲノム解析

- アノテーションファイル(行列形式)からの情報抽出
 - 条件を満たす行のみの抽出。ハッシュとかgrep周辺がRで可能
- multi-fastaファイルからの情報抽出
 - 一定の長さ以上のものを抽出。コンティグごとのGC含量計算。

■ Rで(比較)トランスクリプトーム解析

- 研究目的別留意点(サンプル内比較、サンプル間比較)
- リードマッピング → 数値化(カウントデータの取得)
- 発現変動解析
 - 既存のRパッケージの弱点(unbiased DE: ○、biased DE: ×)
 - TCCパッケージを用いて二群間比較(複製あり、複製なし)



GC含量

前処理	正規化(サンプル間)	二群間用	複製あり	IDEGES/edgeR正規化(Sun submitted) (last modified
前処理	正規化(サンプル間)	二群間用	複製あり	DEGES/TbT正規化(TbT, Kadota 2012) (last modified
前処理	正規化(サンプル間)	二群間用	複製あり	TMM正規化(Robinson 2010) (last modified 2013/03/
前処理	正規化(サンプル間)	全般		Anders and Huberの(AH)正規化(Anders 2010) (last modified 2013/01
前処理	正規化(サンプル間)	全般	Up	
前処理	正規化(サンプル間)	全般	Qu	
前処理	正規化(サンプル間)	全般	RP	
前処理	正規化(混合)	RPKM正規化(Mc		
前処理	正規化(混合)	RPKM正規化(Mc		
解析	一般	アラインメント(ペアワイス)		
解析	一般	アラインメント(ペアワイス)		
解析	一般	アラインメント(ペアワイス)		
解析	一般	パターンマッチング (last mo		
解析	一般	GC含量 (GC contents) (last		
解析	一般	Sequence logos (Schneider		
解析	一般	上流配列解析 Local Distr		
解析	一般	上流配列解析 Relative Ag		
解析	NGS(RNA-seq)	その他	Technic	
解析	NGS(RNA-seq)	その他	ポアソ	
解析	NGS(RNA-seq)	その他	Biologic	
解析	NGS(RNA-seq)	その他	負の二	
解析	NGS(RNA-seq)	その他	負の二	
解析	NGS(RNA-seq)	その他	負の二	

解析 | 一般 | GC含量 (GC contents)

ここでは、multi-fasta形式ファイルを読み込んでコンティグごとのGC含量 (GC contents)を出力するやり方を示します。

ここでは以下の二つを例題として行います:

1. 「イントロダクション | 一般 | [ランダムな塩基配列を作成](#)」の4.を実行して得られたmulti-fastaファイル([hoge4.fa](#))
 2. 「ファイルの読み込み | マップ前 | [FASTA形式](#)」で読み込んだ250 readsからなる[test1.fasta](#)ファイル
- 「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含量を計算をする数式中の「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。

```
1. hoge4.faファイルの場合:
----- ここから -----
in_f <- "hoge4.fa"
out_f <- "hogel.txt"
```

#読み込みたいFASTA形式のファイル名を指定してin_f
#出力ファイル名を指定

```
#必要なパッケージをロード
library(Biostrings)
```

配列ごとのGC含量を計算したいとき

```
#入力ファイルの読み込み
reads <- readDNASTringSet(in_f, format="fasta")
```

#in_fで指定したファイルの読み込み

```
#GC含量(GC content)計算のところ
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
out <- CG/ACGT*100
```

#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

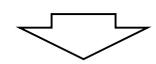
```
#出力用に結果をまとめている
tmp <- cbind(names(reads), CG, ACGT, width(reads), out)
colnames(tmp) <- c("description", "CG", "ACGT", "Length", "%GC_contents")#列名情報を与えている
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F, col.names=T)#tmpの中身をout_fで指定したファイルに出力
```

----- ここまで -----

GC含量

```
R Console
> in_f <- "hoge4.fa"
> out_f <- "hoge1.txt"
>
> #必要なパッケージをロード
> library(Biostrings)
>
> #入力ファイルの読み込み
> reads <- readDNASTringSet(in_f, format="fasta")
>
> #GC含量 (GC content) 計算のところ
> count <- alphabetFrequency(reads)
> CG <- rowSums(count[,2:3])
> ACGT <- rowSums(count[,1:4])
> out <- CG/ACGT*100
>
> #出力用に結果をまとめている
> 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)
> |
```

```
hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
#
>contig_2
GTCTGCCTCAAGCGCCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACACCCT
#
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
#
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
#
```



	A	B	C	D	E
1	description	CG	ACGT	Length	%GC_contents
2	contig_1	16	24	24	66.67
3	contig_2	65	103	103	63.11
4	contig_3	33	65	65	50.77
5	contig_4	25	49	49	51.02

(スライド7の)シロイヌナズナのGC含量計算は、ファイル名部分を変更して実行しただけです

GC含量

```

----- ここから -----
in_f <- "hoge4.fa"
out_f <- "hoge1.txt"

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
reads <- readDNAStringSet(in_f, format="fasta")

#GC含量(GC content)計算のところ
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
out <- CG/ACGT*100

#出力用に結果をまとめている
tmp <- cbind(names(reads), CG, ACGT, width(reads))
colnames(tmp) <- c("description", "CG", "ACGT", "width")
write.table(tmp, out_f, sep="¥t", append=F, quot=FALSE)

----- ここまで -----

```

alphabetFrequency関数を適用すると塩基ごとの出現頻度が得られます
 …が「M, R, W, …」は？

```

hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)

>contig_1
CGGACAGCTCCTCGGCATCCGGAT
#
#
>contig_2
GTCTGCCTCAAGCGCCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACACCCT
#
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
#
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
#
#

```

```

R Console

>
> reads
  A DNASTringSet instance of length 4
  width seq                      names
[1]   24 CGGACAGCTCCTCGGCATCCGGAT contig_1
[2]  103 GTCTGCCTCAAGCGC...CCAGCACACCCTGTC contig_2
[3]   65 TGTAGGAGAAGGGCG...GTATGAGGTCGGGCA contig_3
[4]   49 CGTGCTGATTCCACA...CTACCTATGAACATG contig_4

> count
      A C G T M R W S Y K V H D B N - +
[1,]  4  9  7  4  0  0  0  0  0  0  0  0  0  0  0  0  0  0
[2,] 20 34 31 18  0  0  0  0  0  0  0  0  0  0  0  0  0  0
[3,] 16 13 20 16  0  0  0  0  0  0  0  0  0  0  0  0  0  0
[4,] 14 15 10 10  0  0  0  0  0  0  0  0  0  0  0  0  0  0

> ?alphabetFrequency
> |

```

GC含量

```
R Console
> ?alphabetFrequency
starting httpd help server ... done
> |
```

letterFrequency {Biostrings} R Documentation

Calculate the frequency of letters in a biological sequence, or the consensus matrix of a set of sequences

Description

Given a biological sequence (or a set of biological sequences), the `alphabetFrequency` function computes the frequency of each letter of the relevant [alphabet](#).

Usage

```
alphabetFrequency(x, as.prob=FALSE, ...)
hasOnlyBaseLetters(x)
```

Arguments

`x` An [XString](#), [XStringSet](#), [XStringViews](#) or [MaskedXString](#) object for

Details

`alphabetFrequency`, `letterFrequency`, and `letterFrequencyInSlidingView` are generic functions defined in the `Biostrings` package.

Value

`alphabetFrequency` returns an integer vector when `x` is an [XString](#) or [MaskedXString](#) object. When `x` is an [XStringSet](#) or [XStringViews](#) object, then it returns an integer matrix with `length(x)` rows where the `i`-th row contains the frequencies for `x[[i]]`. If `x` is a DNA or RNA input, then the returned vector is named with the letters in the alphabet. If the `baseOnly` argument is `TRUE`, then the returned vector has only 5 elements: 4 elements corresponding to the 4 nucleotides + the 'other' element.

alphabetFrequency関数が
 どういう出力結果を返すの
 か詳細を知りたい場合は
Value のところを読むとよい

GC含量

Value

alphabetFrequency returns an integer vector when x is an [XString](#) or [MaskedXString](#) object. When x is an [XStringSet](#) or [XStringViews](#) object, then it returns an integer matrix with length(x) rows where the i-th row contains the frequencies for x[[i]]. If x is a DNA or RNA input, then the returned vector is named with the letters in the alphabet. If the baseOnly argument is TRUE, then the returned vector has only 5 elements: 4 elements corresponding to the 4 nucleotides + the 'other' element.

M(A/C), R(A/G), W(A/T), S(C/G), ..., N(A/C/G/T) という事実は常識?!なので書かれていない…。

```
R Console
> alphabetFrequency(reads)
      A C G T M R W S Y K V H D B N - +
[1,]  4  9  7  4  0  0  0  0  0  0  0  0  0  0  0  0  0
[2,] 20 34 31 18  0  0  0  0  0  0  0  0  0  0  0  0  0
[3,] 16 13 20 16  0  0  0  0  0  0  0  0  0  0  0  0  0
[4,] 14 15 10 10  0  0  0  0  0  0  0  0  0  0  0  0  0
> alphabetFrequency(reads, baseOnly = TRUE)
      A C G T other
[1,]  4  9  7  4     0
[2,] 20 34 31 18     0
[3,] 16 13 20 16     0
[4,] 14 15 10 10     0
> |
```

argument (独立変数 or 引数)とは、出力結果の形式などを指定するオプション。関数ごとに異なる

GC含量

Usage

```
alphabetFrequency(x, as.prob=FALSE, ...)
hasOnlyBaseLetters(x)
uniqueLetters(x)
```

Arguments

x An [XString](#), [XStringSet](#), [XStringView](#) object for `alphabetFrequency`, `letterFrequency`, `uniqueLetters`.

DNA or RNA input for `hasOnlyBaseLetters`.

An [XString](#) object for `letterFrequency`.

A character vector, or an [XStringSet](#) or [XStringView](#) object for `consensusMatrix`.

A consensus matrix (as returned by `consensusMatrix`), or an [XStringSet](#) or [XStringView](#) object for `consensusString`.

as.prob If TRUE then probabilities are reported, otherwise counts (the default).

```
R Console
> alphabetFrequency(reads, baseOnly = TRUE)
      A C G T other
[1,]  4 9 7 4     0
[2,] 20 34 31 18    0
[3,] 16 13 20 16    0
[4,] 14 15 10 10    0

> alphabetFrequency(reads, as.prob = TRUE)
      A C G T M R W S Y K V H D B N - +
[1,] 0.1666667 0.3750000 0.2916667 0.1666667 0 0 0 0 0 0 0 0 0 0 0 0 0 0
[2,] 0.1941748 0.3300971 0.3009709 0.1747573 0 0 0 0 0 0 0 0 0 0 0 0 0 0
[3,] 0.2461538 0.2000000 0.3076923 0.2461538 0 0 0 0 0 0 0 0 0 0 0 0 0 0
[4,] 0.2857143 0.3061224 0.2040816 0.2040816 0 0 0 0 0 0 0 0 0 0 0 0 0 0

> alphabetFrequency(reads, as.prob = TRUE, baseOnly = TRUE)
      A C G T other
[1,] 0.1666667 0.3750000 0.2916667 0.1666667 0
[2,] 0.1941748 0.3300971 0.3009709 0.1747573 0
[3,] 0.2461538 0.2000000 0.3076923 0.2461538 0
[4,] 0.2857143 0.3061224 0.2040816 0.2040816 0
> |
```

```
hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACACCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
```

相対出現頻度 (probability) にすることも可能です。

GC含量

```

----- ここから -----
in_f <- "hoge4.fa"
out_f <- "hoge1.txt"

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
reads <- readDNAStringSet(in_f, format="fasta")

#GC含量 (GC content) 計算のところ
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
out <- CG/ACGT*100

#出力用に結果をまとめている
tmp <- cbind(names(reads), CG, ACGT, width(reads), out)
colnames(tmp) <- c("description", "CG", "ACGT", "Length", "GC")
write.table(tmp, out_f, sep="¥t", append=F, quote=F, row.names=F)

----- ここまで -----

```

「コンティグごとのCとGの出現頻度の和」は、「行の和」を計算することで得ることができる

```

R Console
> count
      A  C  G  T  M  R  W  S  Y  K  V  H  D  B  N  -  +
[1,]  4  9  7  4  0  0  0  0  0  0  0  0  0  0  0  0  0
[2,] 20 34 31 18  0  0  0  0  0  0  0  0  0  0  0  0  0
[3,] 16 13 20 16  0  0  0  0  0  0  0  0  0  0  0  0  0
[4,] 14 15 10 10  0  0  0  0  0  0  0  0  0  0  0  0  0
> dim(count)
[1]  4 17
> count[,2:3]
      C  G
[1,]  9  7
[2,] 34 31
[3,] 13 20
[4,] 15 10
> rowSums(count[,2:3])
[1] 16 65 33 25
> |

```

GC含量

```
----- ここから -----
in_f <- "hoge4.fa"
out_f <- "hoge1.txt"
```

```
#必要なパッケージをロード
library(Biostrings)
```

```
#入力ファイルの読み込み
reads <- readDNASTringSet(in_f, format="fasta")
```

```
#GC含量 (GC content) 計算のところ
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
out <- CG/ACGT*100
```

```
#出力用に結果をまとめている
tmp <- cbind(names(reads), CG, ACGT, width(reads), out)
colnames(tmp) <- c("description", "CG", "ACGT", "Length", "%GC_contents") #列名情報を与えている
write.table(tmp, out_f, sep="¥t", append=F, quote=F, row.names=F, col.names=T) #tmpの中身をout_fで指定したファイル名で保存。
```

```
----- ここまで -----
```

```
R Console
> CG
[1] 16 65 33 25
> ACGT
[1] 24 103 65 49
> CG/ACGT
[1] 0.6666667 0.6310680 0.5076923 0.5102041
> CG/ACGT*100
[1] 66.66667 63.10680 50.76923 51.02041
> ATGC
エラー: オブジェクト 'ATGC' がありません
> |
```

出力ファイルの形式やヘッダ行はどのようにして決めているのか？

	A	B	C	D	E
1	description	CG	ACGT	Length	%GC_contents
2	contig_1	16	24	24	66.67
3	contig_2	65	103	103	63.11
4	contig_3	33	65	65	50.77
5	contig_4	25	49	49	51.02

GC含量

```

----- ここから -----
in_f <- "hoge4.fa"
out_f <- "hoge1.txt"

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
reads <- readDNAStringSet(in_f, format="fasta")

#GC含量 (GC content) 計算のところ
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
out <- CG/ACGT*100

#出力用に結果をまとめている
tmp <- cbind(names(reads), CG, ACGT, width(reads), out)
colnames(tmp) <- c("description", "CG", "ACGT", "Length")
write.table(tmp, out_f, sep="¥t", append=F, quote=F, row

----- ここまで -----

```

```

R Console
> reads
A DNAStringSet instance of length 4
  width seq          names
[1]    24 CGGACAGCT...ATCCGGAT contig_1
[2]   103 GTCTGCCTC...ACCCTGTC contig_2
[3]    65 TGTAGGAGA...GTCGGGCA contig_3
[4]    49 CGTGCTGAT...TGAACATG contig_4
> names(reads)
[1] "contig_1" "contig_2" "contig_3" "contig_4"
> CG
[1] 16 65 33 25
> cbind(names(reads), CG)
           CG
[1,] "contig_1" "16"
[2,] "contig_2" "65"
[3,] "contig_3" "33"
[4,] "contig_4" "25"
> |

```

	A	B	C	D	E
1	description	CG	ACGT	Length	%GC_contents
2	contig_1	16	24	24	66.67
3	contig_2	65	103	103	63.11
4	contig_3	33	65	65	50.77
5	contig_4	25	49	49	51.02

同じ要素数からなるベクトル同士を列 (column) 方向で結合 (bind)

GC含量

```
----- ここから -----
in_f <- "hoge4.fa"
out_f <- "hoge1.txt"
```

```
#必要なパッケージをロード
library(Biostrings)
```

```
#入力ファイルの読み込み
reads <- readDNAStringSet(in_f, format="fasta")
```

```
#GC含量 (GC content) 計算のところ
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
out <- CG/ACGT*100
```

```
#出力用に結果をまとめている
tmp <- cbind(names(reads), CG, ACGT, width(reads), out)
colnames(tmp) <- c("description", "CG", "ACGT", "Length", "%GC_contents") #列名情報を与えている
write.table(tmp, out_f, sep="¥t", append=F, quote=F, row.names=F, col.names=T) #mpの中身をout_fで指定したファイル名で保存。
```

```
----- 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
#A,C,G,Tの総数を計算
#GC含量を計算

cbind関数適用時に、出力したい順番に並べているだけです

#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
#mpの中身をout_fで指定したファイル名で保存。

	A	B	C	D	E
1	description	CG	ACGT	Length	%GC_contents
2	contig_1	16	24	24	66.67
3	contig_2	65	103	103	63.11
4	contig_3	33	65	65	50.77
5	contig_4	25	49	49	51.02

GC含量

```
#出力用に結果をまとめている
tmp <- cbind(names(reads), CG, ACGT, width(reads), out) #ファイルに出力したい情報を連結してtmpに格納
colnames(tmp) <- c("description", "CG", "ACGT", "Length", "%GC_contents") #列名情報を与えている
write.table(tmp, out_f, sep="¥t", append=F, quote=F, row.names=F, col.names=T) #tmpの中身をout_fで指定したファイル名で保存。

----- ここまで -----
```

```
R Console
>
> #GC含量 (GC content) 計算のところ
> count <- alphabetFrequency(reads)
> CG <- rowSums(count[,2:3])
> ACGT <- rowSums(count[,1:4])
> out <- CG/ACGT*100
>
> #出力用に結果をまとめている
> tmp <- cbind(names(reads), CG, ACGT, width(reads), out)
> tmp
      CG  ACGT  out
[1,] "contig_1" "16" "24" "24" "66.6666666666667"
[2,] "contig_2" "65" "103" "103" "63.1067961165049"
[3,] "contig_3" "33" "65" "65" "50.7692307692308"
[4,] "contig_4" "25" "49" "49" "51.0204081632653"
> colnames(tmp) <- c("description", "CG", "ACGT", "Length", "%GC_contents")
> tmp
  description CG  ACGT Length %GC_contents
[1,] "contig_1" "16" "24" "24" "66.6666666666667"
[2,] "contig_2" "65" "103" "103" "63.1067961165049"
[3,] "contig_3" "33" "65" "65" "50.7692307692308"
[4,] "contig_4" "25" "49" "49" "51.0204081632653"
> |
```

	A	B	C	D	E
1	description	CG	ACGT	Length	%GC_contents
2	contig_1	16	24	24	66.67
3	contig_2	65	103	103	63.11
4	contig_3	33	65	65	50.77
5	contig_4	25	49	49	51.02

列 (column) の名前 (names) を 適当に変更して出力しています

パッケージ (package) って何？

```
R Console

R version 2.15.2 (2012-10-26) -- "Trick or Treat"
Copyright (C) 2012 The R Foundation for Statistical Comp$
ISBN 3-900051-07-0
Platform: x86_64-w64-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()'$

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$
'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
'q()'と入力すればRを終了します。

[以前にセーブされたワークスペースを復帰します]

> ?alphabetFrequency
No documentation for 'alphabetFrequency' in specified $
you could try '??alphabetFrequency'
> |
```

Rを再起動した状態で、「?関数名」と打ち込んでも使用法を記したウェブページが開かずにエラーが出ることがあります

パッケージ (package) って何？

```

-----   ここから   -----
in_f <- "hoge4.fa"
out_f <- "hoge1.txt"

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
reads <- readDNASTringSet(in_f, format="fasta")

#GC含量(GC content)計算のところ
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
out <- CG/ACGT*100

#出力用に結果をまとめている
tmp <- cbind(names(reads), CG, ACGT,
             colnames(tmp) <- c("description", "CG", "ACGT"))
write.table(tmp, out_f, sep="\t", append=TRUE)

-----   ここまで   -----

```

「Biostrings」というパッケージをlibrary関数を用いて読み込むことによって、Biostringsパッケージが提供する関数群の中の一つであるalphabetFrequency関数を利用できるんです。

```

R Console
> ?alphabetFrequency
No documentation for 'alphabetFrequency' in specified $
you could try '??alphabetFrequency'
> library(Biostrings)
要求されたパッケージ BiocGenerics をロード中です

次のパッケージを付け加えます: 'BiocGenerics'

The following object(s) are masked from 'package:stats':
  xtabs

The following object(s) are masked from 'package:base':
  anyDuplicated, cbind, colnames, duplicated,
  eval, Filter, Find, get, intersect, lapply,
  Map, mapply, mget, order, paste, pmax,
  pmax.int, pmin, pmin.int, Position, rbind,
  Reduce, rep.int, rownames, sapply, setdiff,
  table, tapply, union, unique

要求されたパッケージ IRanges をロード中です
> ?alphabetFrequency
starting httpd help server ... done
> |

```

た結果をcountに格納

に格納

たファイル名で保存。

ここでは、multi-fasta形式ファイルを読み込んでコンテグごとのGC含量 (GC contents)を出力するやり方を示します

ここでは以下の二つを例題として行います:
1. 「イントロダクション | 一般 | ランダムな塩基配列を作成」の4を実行して得られた
2. 「ファイルの読み込み | マップ前 | FASTA形式」で読み込んだ250 readsからなる
「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
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、
「長」を分子にしたい場合にはGC含量を計算する数式中の「CG/ACGT*100」を「CG
「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

```
1. hoge4.faファイルの場合:
-----
in_f <- "hoge4.fa"
out_f <- "hoge1.txt"

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
reads <- readDNASTringSet(in_f, format="fasta")

#GC含量 (GC content)計算のところで
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
out <- CG/ACGT*100

#出力用に結果をまとめている
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)

-----
2. test1.fastaファイルの場合:
-----
in_f <- "test1.fasta"
out_f <- "hoge2.txt"

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
reads <- readDNASTringSet(in_f, format="fasta")

#GC含量 (GC content)計算のところで
count <- alphabetFrequency(reads)
CG <- rowSums(count[,2:3])
ACGT <- rowSums(count[,1:4])
out <- CG/ACGT*100

#出力用に結果をまとめている
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)

-----
```

[BioconductorのBiostringsのwebページ](#)

Bioconductor logo: OPEN SOURCE SOFTWARE FOR BIOINFORMATICS
Home Install Help Devel

Home » Bioconductor 2.12 » Software Packages » Biostrings

Biostrings

String objects representing biological sequences, and matching algorithms

Bioconductor version: Release (2.12)

Memory efficient string containers, string matching algorithms, and other utilities, for fast manipulation of large biological sequences or sets of sequences.

Author: H. Pages, P. Aboyoun, R. Gentleman, and S. DebRoy

Maintainer: H. Pages <hpages at fhrc.org>

To install this package, start R and enter:

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

パッケージを個別にインストールする場合

To cite this package in a publication, start R and enter:

```
citation("Biostrings")
```

Documentation

PDF	R Script	A short presentation of the basic classes defined in Biostrings
PDF	R Script	Biostrings Quick Overview
PDF	R Script	Handling probe sequence information
PDF	R Script	Multiple Alignments
PDF	R Script	Pairwise Sequence Alignments
PDF		Reference Manual
Text		NEWS

使い方の解説記事はPDFのところをクリック。例えば...

Details

biocViews [DataImport](#), [DataRepresentation](#), [Genetics](#), [Infrastructure](#), [SequenceMatching](#), [Sequencing](#), [Software](#)

Hervé Pagès
Fred Hutchinson Cancer Research Center
Seattle, WA

April 3, 2013

Table 2: Basic transformations of sequences.

Please note that *most* but *not all* the functionalities provided by the Biostrings package are listed in this document.

Function	Description
<code>length</code>	Return the number of sequences in an object.
<code>names</code>	Return the names of the sequences in an object.
<code>[]</code>	Extract sequences from an object.
<code>head, tail</code>	Extract the first or last sequences from an object.
<code>rev</code>	Reverse the order of the sequences in an object.
<code>c</code>	Put in a single object the sequences from 2 or more objects.
<code>width, nchar</code>	Return the sizes (i.e. number of letters) of all the sequences in an object.
<code>==, !=</code>	Element-wise comparison of the sequences in 2 objects.
<code>match, %in%</code>	Analog to <code>match</code> and <code>%in%</code> on character vectors.
<code>duplicated, unique</code>	Analog to <code>duplicated</code> and <code>unique</code> on character vectors.
<code>sort, order</code>	Analog to <code>sort</code> and <code>order</code> on character vectors, except that the ordering of DNA or Amino Acid sequences doesn't depend on the locale.
<code>split, relist</code>	Analog to <code>split</code> and <code>relist</code> on character vectors, except that the result is a <code>DNAStringSetList</code> or <code>AAStringSetList</code> object.

Table 1: Low-level manipulation of `DNAStringSet` or `AAStringSet` objects.

Function	Description
<code>subseq, subseq<-</code>	Extract or replace subsequences in a set of sequences.
<code>reverse</code>	Compute the reverse, complement, or reverse-complement, of a set of DNA sequences.
<code>complement</code>	
<code>reverseComplement</code>	
<code>translate</code>	Translate a set of DNA sequences into a set of Amino Acid sequences.
<code>chartr</code>	
<code>replaceLetterAt</code>	

Function	Description
<code>alphabetFrequency</code> <code>letterFrequency</code>	Tabulate the letters (all the letters in the alphabet for <code>alphabetFrequency</code> , only the specified letters for <code>letterFrequency</code>) of a sequence or set of sequences.
<code>letterFrequencyInSlidingView</code>	Specialized version of <code>letterFrequency</code> that tallies the requested letter frequencies for a fixed-width view that is conceptually slid along the input sequence.
<code>consensusMatrix</code>	Computes the consensus matrix of a set of sequences.
<code>dinucleotideFrequency</code> <code>trinucleotideFrequency</code> <code>oligonucleotideFrequency</code>	Fast 2-mer, 3-mer, and k-mer counting for DNA or RNA.
<code>nucleotideFrequencyAt</code>	Tallies the short sequences formed by extracting the nucleotides found at a set of fixed positions from each sequence of a set of DNA or RNA sequences.

Table 3: Counting / tabulating.

Function	Description
<code>matchPattern</code> <code>countPattern</code>	Find/count all the occurrences of a given pattern (typically short) in a reference sequence (typically long). Support mismatches and indels.
<code>vmatchPattern</code> <code>vcountPattern</code>	Find/count all the occurrences of a given pattern (typically short) in a set of reference sequences. Support mismatches and indels.
<code>matchPDict</code> <code>countPDict</code> <code>whichPDict</code>	Find/count all the occurrences of a set of patterns in a reference sequence. (<code>whichPDict</code> only identifies which patterns in the set have at least one match.) Support a small number of mismatches.
<code>vmatchPDict</code> <code>vcountPDict</code> <code>vwhichPDict</code>	[Note: <code>vmatchPDict</code> not implemented yet.] Find/count all the occurrences of a set of patterns in a set of reference sequences. (<code>whichPDict</code> only identifies for each reference sequence which patterns in the set have at least one match.) Support a small number of mismatches.
<code>pairwiseAlignment</code>	Solve (Needleman-Wunsch) global alignment, (Smith-Waterman) local alignment, (Sellers) overlap alignment problems.
<code>matchProbePair</code>	Find all the amplicons that match a pair of probes in a reference sequence.

「Biostrings」パッケージ中の関数を使いこなせると、他の自然言語処理系プログラミング言語(perlやruby)を改めて勉強しなくても必要な解析の大部分が可能です

Contents

■ Rでゲノム解析

- アノテーションファイル(行列形式)からの情報抽出
 - 条件を満たす行のみの抽出。ハッシュとかgrep周辺がRで可能
- multi-fastaファイルからの情報抽出
 - 一定の長さ以上のものを抽出。コンティグごとのGC含量計算。

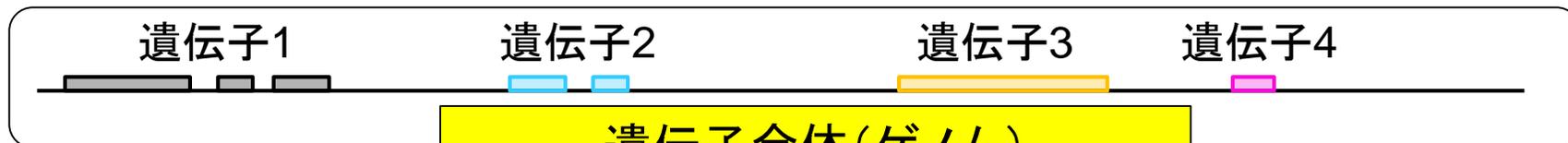
■ Rで(比較)トランスクリプトーム解析

- 研究目的別留意点(サンプル内比較、サンプル間比較)
- リードマッピング → 数値化(カウントデータの取得)
- 発現変動解析
 - 既存のRパッケージの弱点(unbiased DE: ○、biased DE: ×)
 - TCCパッケージを用いて二群間比較(複製あり、複製なし)

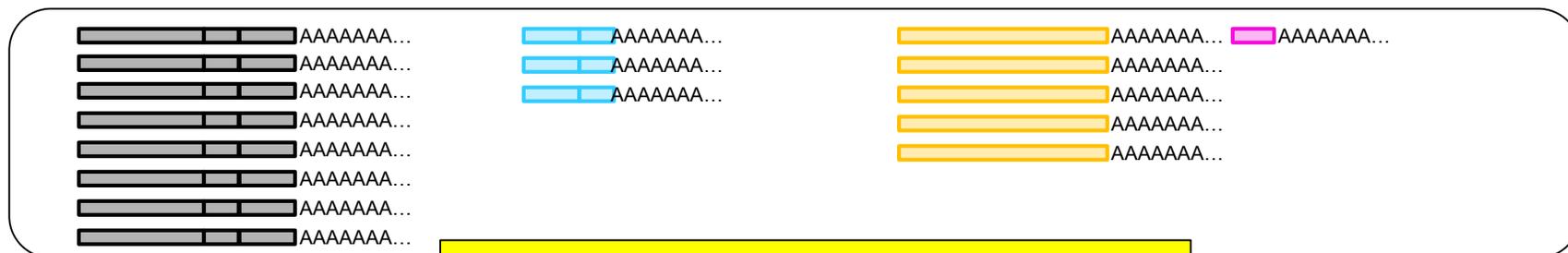
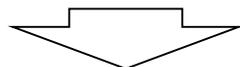


トランスクリプトームとは

- ある状態のあるサンプル(例:目)のあるゲノムの領域



・どの染色体上のどの領域にどの遺伝子があるかは調べる個体(例:ヒト)が同じなら不変(目だろうが心臓だろうが...)



転写物全体(トランスクリプトーム)

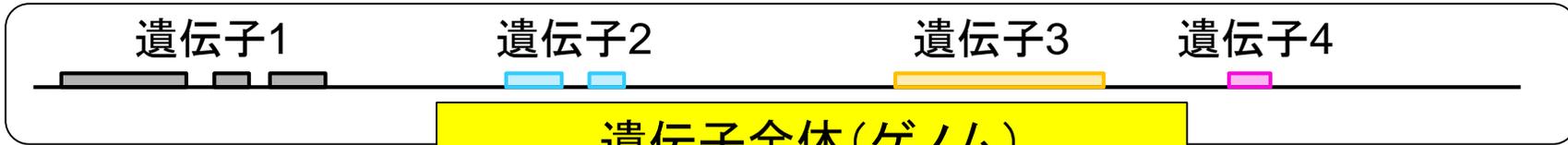
- ・遺伝子1は沢山転写されている(発現している)
- ・遺伝子4はごくわずかしか転写されていない
- ・...

トランスクリプトームとは

- ある状態のあるサンプル(例:目)のあるゲノムの領域

光刺激

ヒト



・どの染色体上のどの領域にどの遺伝子があるかは調べる個体(例:ヒト)が同じなら不変(目だろうが心臓だろうが...)

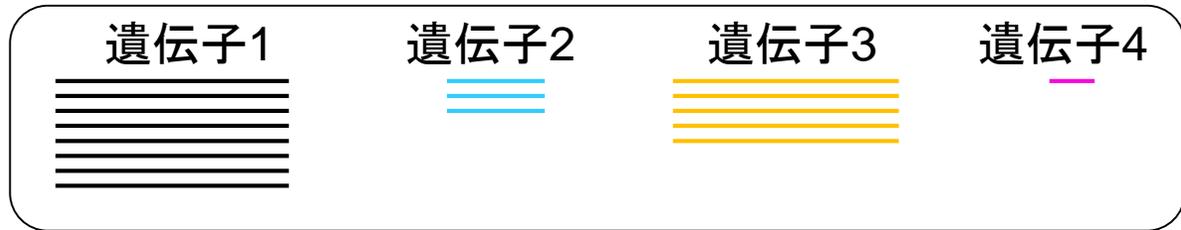


転写物全体(トランスクリプトーム)

- ・遺伝子2は光刺激に应答して発現亢進
- ・遺伝子4も光刺激に应答して発現亢進

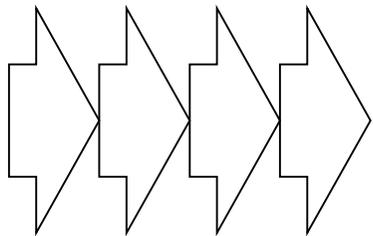
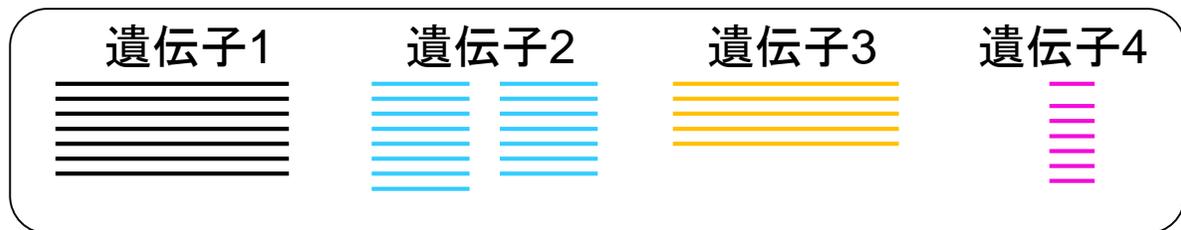
トランスクリプトーム情報を得る手段

■ 光刺激前 (T1) の目のトランスクリプトーム



これがいわゆる
「遺伝子発現行列」

■ 光刺激後 (T2) の目のトランスクリプトーム



	T1	T2
遺伝子1	8	7
遺伝子2	3	15
遺伝子3	5	5
遺伝子4	1	7
...

- ・マイクロアレイ
- ・電気泳動に基づく方法
- ・配列決定に基づく方法

比較トランスクリプトーム解析の流れ

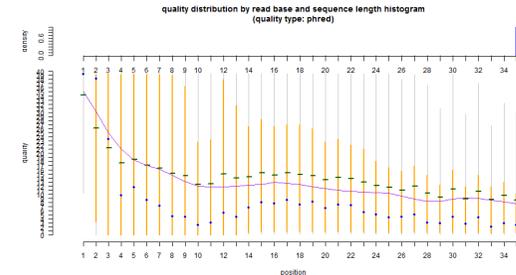
複数のFASTQファイル

リファレンス配列の作成

マッピング

データ解析

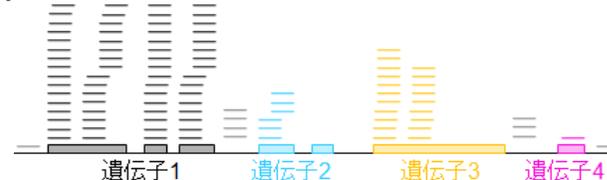
クオリティチェック



アセンブル結果 (multi-fasta) ファイルから平均長やトータルの長さなどの基本情報を抽出

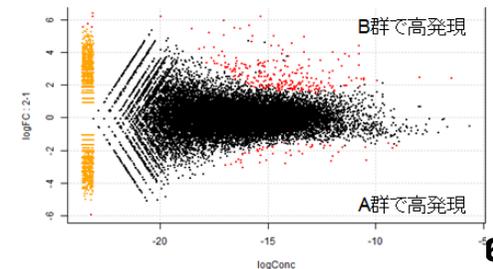
Total length (bp)	2993
Number of contigs	4
Average length	748.3
Median length	784
Max length	888
Min length	537
N50	886
GC content	0.524

マッピング結果 (BED形式) ファイルを入力として、転写物ごとのマップされたリード数をカウント



遺伝子1	40
遺伝子2	6
遺伝子3	20
遺伝子4	1

発現変動遺伝子のリストアップや、作図など



雑感

- 一口にトランスクリプトーム解析といっても目的や手段は多様。
 - トランスクリプトーム配列取得
 - ゲノム配列既知の場合、Cufflinksなどを用いて遺伝子構造推定(アノテーション)
 - ゲノム配列未知の場合、Trinityなどのトランスクリプトーム用アセンブラを実行
 - 遺伝子(or isoform)ごとの発現量推定
 - RSEMなどを利用して発現量情報を得る
 - ある特定のサンプル内での遺伝子間の発現量の大小関係を知りたい
 - Length biasやGC含量biasなどの各種補正がポイント
 - 遺伝子レベルの発現量 \Leftrightarrow isoformレベルの発現量の正確な推定
 - 比較するサンプル間で発現変動している遺伝子(or isoform)の同定
 - TCCパッケージ(など)を利用して発現変動遺伝子(DEG)を得る
 - Sequence depthやサンプル間で発現している遺伝子のcomposition biasの補正がポイント
 - (GO解析など)DEG結果を用いる多くの下流解析結果に影響を及ぼす

研究目的別留意点

■ ある特定のサンプル内での遺伝子間の発現量の大小関係を知りたい場合

□ 「配列長」由来bias: 長いほど沢山sequenceされる

□ 「GC含量」由来bias: カウント数の分布がGC含量依存的である

■ サンプル間比較 (sample A vs. Bなど) で、発現変動遺伝子 (DEG) を調べたい場合

□ 「sequence depthの違い」: 総リード数が x 倍違うと全体的に x 倍変動...

□ 「組成の違い」: サンプル特異的高発現遺伝子の存在で比較困難に...

配列長を考慮した発現量推定のイメージ

- gene1: 3 exons (middle length), 14 reads mapped (**low** coverage)
- gene2: 3 exons (middle length), 56 reads mapped (**high** coverage)
- gene3: 2 exons (**short** length), 12 reads mapped (middle coverage)
- gene4: 2 exons (**long** length), 31 reads mapped (middle coverage)

マップされたリード分布

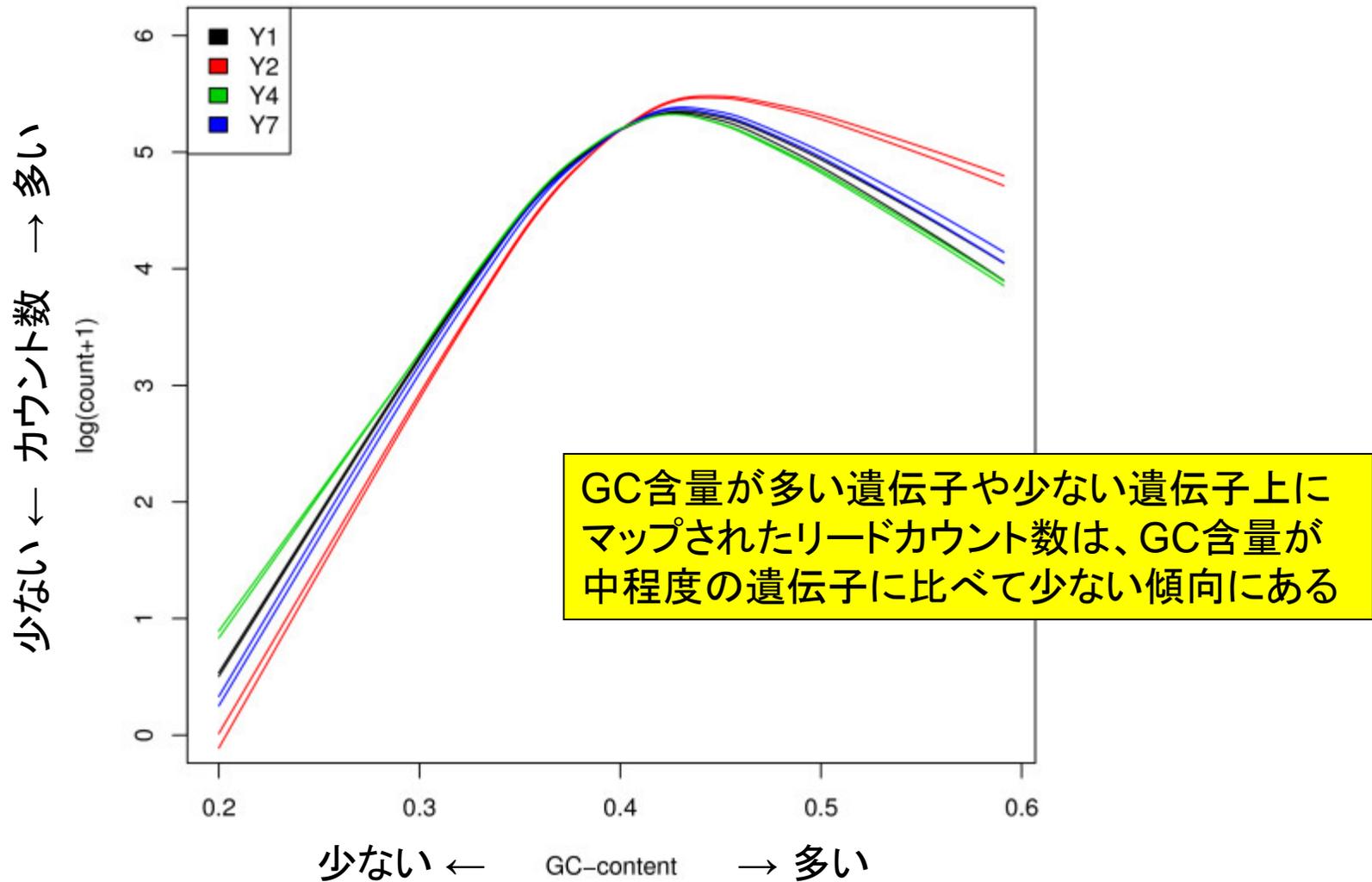
生リードカウント結果

補正度の発現量

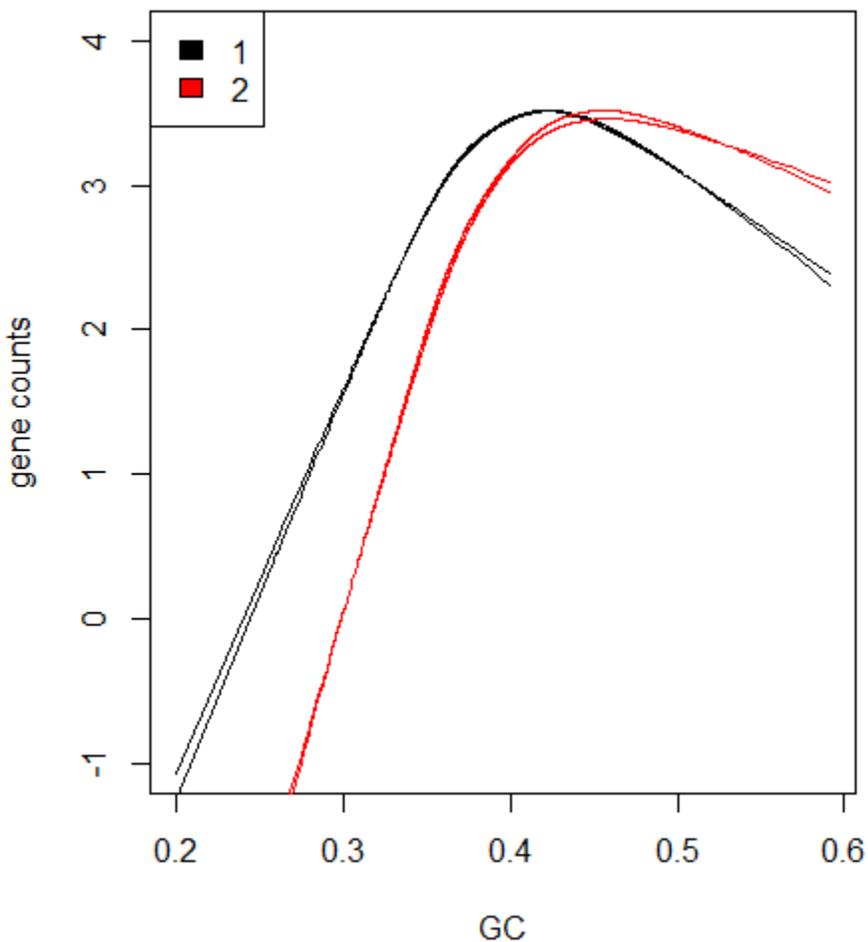
- ・長さが同じならリード数の多い方が発現量高い (gene 1 vs. 2)
- ・長いほどマップされるリード数が多くなる効果を補正する必要がある (gene 3 vs. 4)

一つのサンプル内で転写物(遺伝子)間の発現レベルの大きさを比較したい場合には
配列長を考慮すべきである

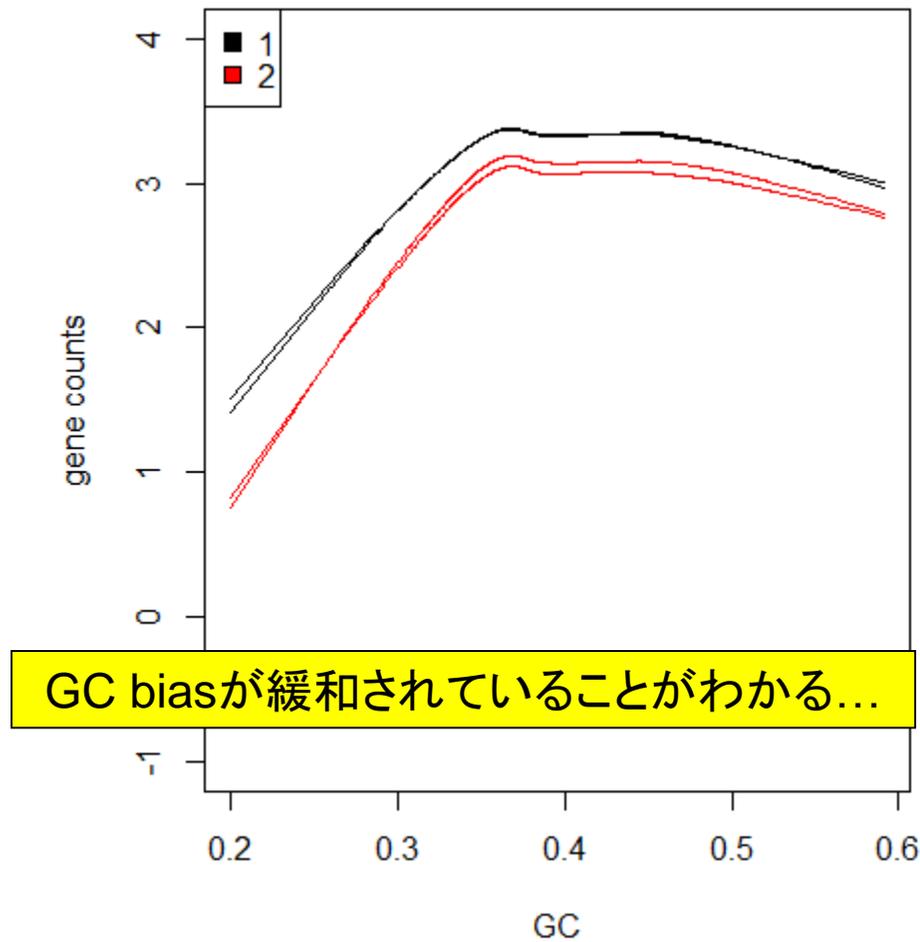
GC biasの実例



GC bias補正 (EDASeqパッケージ)



Quantile
正規化

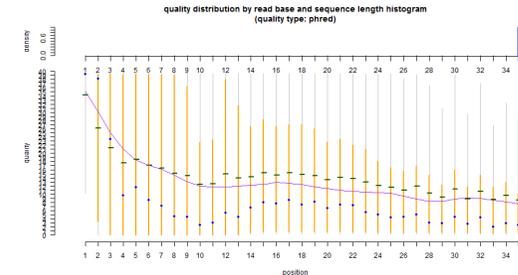


GC biasが緩和されていることがわかる...

比較トランスクリプトーム解析の流れ

複数のFASTQファイル

クオリティチェック



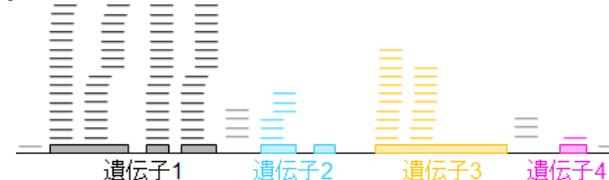
リファレンス配列の作成

アセンブル結果 (multi-fasta) ファイルから平均長やトータルの長さなどの基本情報を抽出

Total length (bp)	2993
Number of contigs	4
Average length	748.3
Median length	784
Max length	888
Min length	537
N50	886
GC content	0.524

マッピング

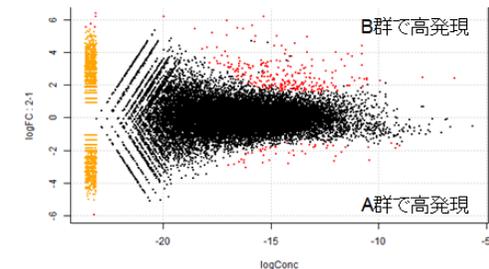
マッピング結果 (BED形式) ファイルを入力として、転写物ごとのマップされたリード数をカウント



遺伝子1	40
遺伝子2	6
遺伝子3	20
遺伝子4	1

データ解析

発現変動遺伝子のリストアップや、作図など



Contents

■ Rでゲノム解析

- アノテーションファイル(行列形式)からの情報抽出
 - 条件を満たす行のみの抽出。ハッシュとかgrep周辺がRで可能
- multi-fastaファイルからの情報抽出
 - 一定の長さ以上のものを抽出。コンティグごとのGC含量計算。

■ Rで(比較)トランスクリプトーム解析

- 研究目的別留意点(サンプル内比較、サンプル間比較)
- リードマッピング → 数値化(カウントデータの取得)
- 発現変動解析
 - 既存のRパッケージの弱点(unbiased DE: ○、biased DE: ×)
 - TCCパッケージを用いて二群間比較(複製あり、複製なし)



マッピング = (大量高速)文字列検索

- マップされる側の配列: 4コンティグ (or 4遺伝子 or 4染色体)
- マップする側のNGS由来塩基配列データ: "AGG"

```
hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACACCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
```

出力ファイル: hoge2.txt

	start	end
contig_2	31	33
contig_2	77	79
contig_3	4	6
contig_3	10	12
contig_3	56	58

Rでやってみよう

パターンマッチング

- 前処理 | [75percentile正規化\(第3四分位数を揃える\)](#) (last modified 2011/08/09)
- 前処理 | [GAM正規化\(Zheng 2011\)](#) (last modified 2011/07/25)
- 解析 | 一般 | [アラインメント\(ペアワイズ;基本編1\)](#) (last modified 2010/6/8)
- 解析 | 一般 | [アラインメント\(ペアワイズ;基本編2\)](#) (last modified 2010/6/8)
- 解析 | 一般 | [アラインメント\(ペアワイズ;応用編\)](#) (last modified 2010/6/8)
- 解析 | 一般 | [パターンマッチング](#) (last modified 2012/04/13) **NEW**
- 解析 | 一般 | [GC含量\(GC content\)](#) (last modified 2010/7/1)
- 解析 | 一般 | [Sequence logos \(Schneider 1990\)](#) (last modified 2012/04/04) **NEW**



• 解析 | 一般 | パターンマッチング

読み込んだリファレンス配列(reference sequence or subject sequence)から(短い)配列パターンを探す場合に利用します。

ここでは、4. multi-fastaファイル [hoge4.fa](#) を入力として、“AGG”でキーワード探索を行う場合：

1. Dihydrof ----- ここから -----

```
Finger Nucl in_f <- "hoge4.fa"
2. 1.と同じ out_f <- "hoge2.txt"
3. 2.と同じ param <- "AGG"
```

```
4. multi-fas
5. multi-fas #必要なパッケージをロード
6. multi-fas library(Biostrings)
7. multi-fas
```

```
(記述の仕 #入力ファイルの読み込み
seq <- readDNASTringSet(in_f, format="fasta")
```

```
#本番
out <- vmatchPattern(pattern=param, subject=seq)
hoge <- cbind(start(unlist(out)), end(unlist(out)))
colnames(hoge) <- c("start", "end")
rownames(hoge) <- names(unlist(out))
tmp <- cbind(rownames(hoge), hoge)
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F)
```

----- ここまで -----

#読み込みたいFASTA形式のファイル
#出力ファイル名を指定してout_fに
#調べたい配列パターンを指定して

#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

#paramで指定した配列と100%マッ
#一致領域の(start, end)の位置情
#行列hogeに列名を付加
#行列hogeに行名を付加
#ファイルに出力したい情報を連結
#tmpの中身をout_fで指定したファ

4. multi-fastaファイルhoge4.faを入力として、“AGG”でキーワード探索を行う場合：

```
----- ここから -----
in_f <- "hoge4.fa"
out_f <- "hoge2.txt"
param <- "AGG"
```

読み込みたいFASTA形式のファイル
出力ファイル名を指定してout_fに
調べたい配列パターンを指定して

• 解析 | 一般 | パターンマッチング

#必要なパッケージをロード
library(Biostrings)

#パッケージの読み込み

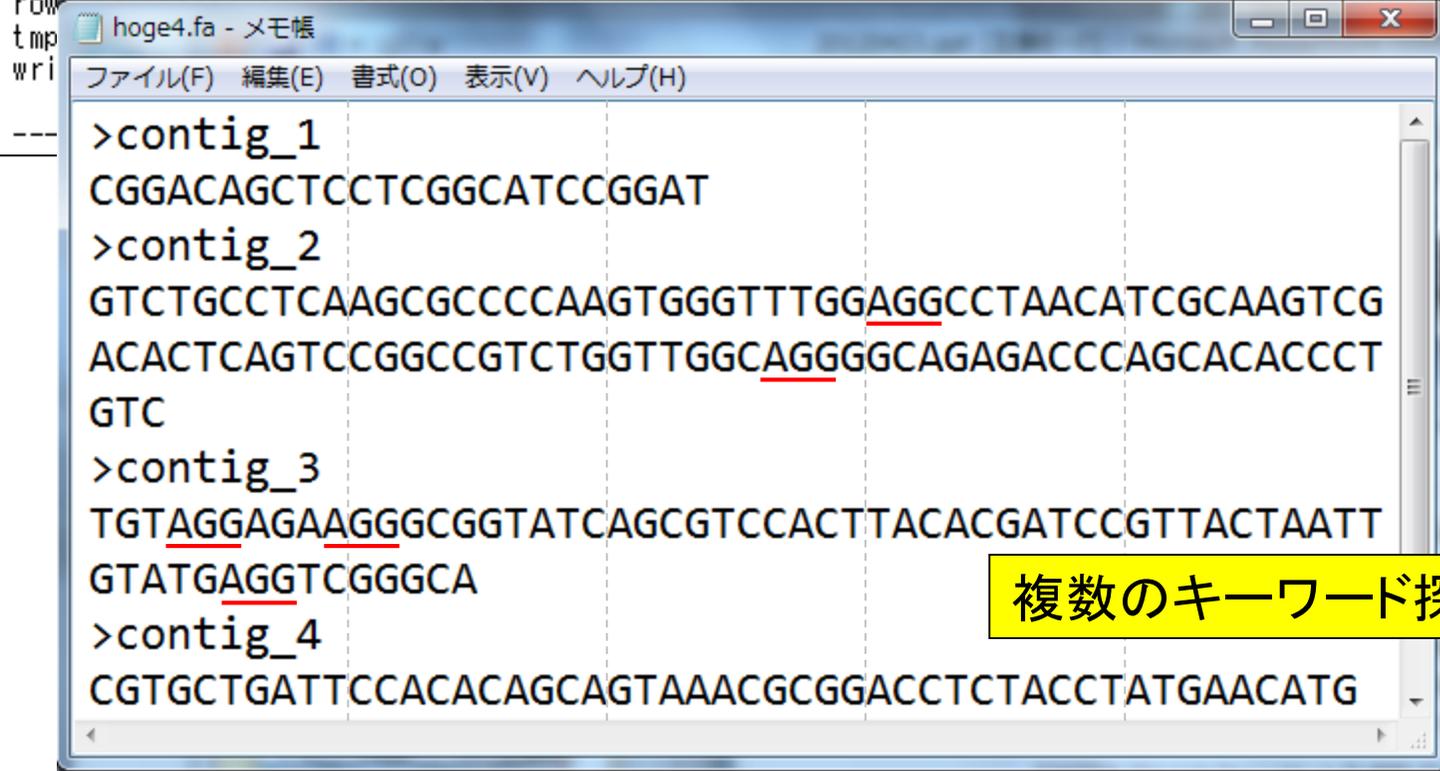
#入力ファイルの読み込み
seq <- readDNAStringSet(in_f, format="fasta")

#in_fで指定したファイルの読み込み

出力ファイル:hoge2.txt

#本番
out <- vmatchPattern(pattern=param, subject=seq)
hoge <- cbind(start(unlist(out)), end(unlist(out)))
colnames(hoge) <- c("start", "end")

#paramで指定した配列と100%一致領域の(start, end)の位置
#行列hogeに列名を付加
#行列hogeに列名を付加



	start	end
contig_2	31	33
contig_2	77	79
contig_3	4	6
contig_3	10	12
contig_3	56	58

複数のキーワード探索も可能です

6. multi-fastaファイル hoge4.fa をリファレンス配列 (マップされる側) として、4リードからなる data_reads.

```
----- ここから -----
in_f1 <- "hoge4.fa"
in_f2 <- "data_reads.txt"
out_f <- "hoge4.txt"
```

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み

```
seq <- readDNAStringSet(in_f1, format="fasta")
reads <- readDNAStringSet(in_f2, format="fasta")
```

#本番

```
out <- c("in_f2", "in_f1", "start", "end")
for(i in 1:length(reads)){
  tmp <- vmatchPattern(pattern=as.character(reads[i]), subject=seq)
```

```
data_reads.txt
>seq1
TTT
>seq2
GGG
>seq3
ACT
>seq4
ACA
```

出力ファイル: hoge4.txt

in_f2	in_f1	start	end
TTT	contig_2	26	28
GGG	contig_2	23	25
GGG	contig_2	78	80
GGG	contig_2	79	81
GGG	contig_3	11	13
GGG	contig_3	61	63
ACT	contig_2	53	55
ACT	contig_3	28	30
ACT	contig_3	44	46
ACA	contig_1	4	6
ACA	contig_2	38	40
ACA	contig_2	51	53
ACA	contig_2	94	96
ACA	contig_3	32	34
ACA	contig_4	13	15
ACA	contig_4	15	17
ACA	contig_4	45	47

```
hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACACCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
```

出力結果ファイルと発現量の関係

出力ファイル:hoge4.txt

```

hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGTTGGCAGGGGCAGAGACCCAGCACACCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
    
```

data_reads.txt

```

>seq1
TTT -
>seq2
GGG -
>seq3
ACT -
>seq4
ACA -
    
```

in_f2	in_f1	start	end
TTT	contig_2	26	28
GGG	contig_2	23	25
GGG	contig_2	78	80
GGG	contig_2	79	81
GGG	contig_3	11	13
GGG	contig_3	61	63
ACT	contig_2	53	55
ACT	contig_3	28	30
ACT	contig_3	44	46
ACA	contig_1	4	6
ACA	contig_2	38	40
ACA	contig_2	51	53
ACA	contig_2	94	96
ACA	contig_3	32	34
ACA	contig_4	13	15
ACA	contig_4	15	17
ACA	contig_4	45	47



multi mapper (複数個所にマップされるリード) の取り扱いは？

よく見かけるカウントデータ取得条件

- basic alignerの一つであるBowtieプログラムを利用して、リファレンス配列(ゲノム or トランスクリプトーム)の1カ所とのみ(最大2塩基ミスマッチまで許容して)一致するリード(uniquely mapped reads or unique mapper)数をカウント
 - Marioni et al., *Genome Res.*, **18**:1509-1517, 2008
 - Bullard et al., *BMC Bioinformatics*, **11**:94, 2010
 - Risso et al., *BMC Bioinformatics*, **12**:480, 2011
 - ReCount (Frazee et al., *BMC Bioinformatics*, **12**:449, 2011)
 - ...

Unique mapperのみにすると...

出力ファイル:hoge4.txt

```

hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGTTGGCAGGGGCAGAGACCCAGCACACCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
    
```

data_reads.txt

```

>seq1
TTT -
>seq2
GGG -
>seq3
ACT -
>seq4
ACA -
    
```

in_f2	in_f1	start	end
TTT	contig_2	26	28
GGG	contig_2	23	25
GGG	contig_2	78	80
GGG	contig_2	79	81
GGG	contig_3	11	13
GGG	contig_3	61	63
ACT	contig_2	53	55
ACT	contig_3	28	30
ACT	contig_3	44	46
ACA	contig_1	4	6
ACA	contig_2	38	40
ACA	contig_2	51	53
ACA	contig_2	94	96
ACA	contig_3	32	34
ACA	contig_4	13	15
ACA	contig_4	15	17
ACA	contig_4	45	47



contig_1	1	➔	contig_1	0
contig_2	8		contig_2	1
contig_3	5		contig_3	0
contig_4	3		contig_4	0

- ・フィルタリング | [NGS | 末端部分のトリミング](#) (last modified 2012/09/12)
- ・前処理 | [について](#) (last modified 2010/12/16)
- ・前処理 | [Trinity出力ファイルからFPKM値を取得](#) (last modified 2011/10/20)
- ・前処理 | [トランスクリプトーム配列へのマップ後のファイルからマップされたリード数をカウント\(BED形式ファイル\)](#)
- ・前処理 | [Ensembl Geneの長さを計算する](#) (last modified 2011/08/15)
- ・前処理 | [転写物の配列長が長いほどマップされたリード数が増えることを確認](#) (last modified 2012/09/11)
- ・前処理 | [転写物の配列長が長いほど発現変動遺伝子とされる確率が上昇することを確認](#) (last modified 2011/09/11)
- ・前処理 | [発現レベルの定量化について](#) (last modified 2013/01/22) **NEW**
- ・前処理 | [正](#)
- ・前処理 | [正](#)



・前処理 | [トランスクリプトーム配列へのマップ後のファイルからマップされたリード数をカウント\(BED形式ファイル\)](#)

手元に「**multi-fasta形式のマップされる側の塩基配列ファイル**(例: [h_rna.fasta](#) や [hoge4.fa](#))」と Bowtie (Langmead et al., 2009)などのマッピングプログラムを用いて得られた「どこにマップされたかを表す**BED形式ファイル**(例: [SRR002324_t.bed](#), [sample_1.bed](#), [sample_2.bed](#))」があるとします。

発現変動遺伝子(DEG)同定の目的で用いられる [TCO](#) や [edgeR](#) などのRパッケージは遺伝子(あるいは転写物)ごとに何個のリードがマップされたかという「**カウント行列(count matrix)**」を入力とすることを前提としています。(ざっくり言えば)BED形式ファイルは、トランスクリプトーム配列にマップした場合、「どのtranscript(or 遺伝子)の(1列目)、どの位置から(2列目; start)、どの位置まで(3列目; end)にリードがマップされた」という情報からなります。それゆえ「**行数がマップされたリード数**」に相当し、「**遺伝子ごとの出現数がカウント数**」に相当するわけです。

したがって、基本的にはBED形式ファイルを読み込んで「**遺伝子名に相当する列**」の部分の情報のみを用いて「**どの遺伝子名が何回出現したか**」をカウントすればいいだけなのですが... 実際には、二つのサンプル(or 二つのBED形式ファイル)を比較する場合などには「**片方のBEDファイル中に一度も出現しない遺伝子**」が存在しうるため、「**multi-fasta形式のマップされる側の塩基配列ファイル**」も同時に読み込んでその遺伝子リストの並びで出力してやったほうが都合がいいのでここではそうしています。

「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。

1. マップされる側の配列 ([hoge4.fa](#)) 中のIDの並びでBED形式ファイル([sample_1.bed](#))中の出現回数をカウントする場合:

```
----- ここから -----
in_f1 <- "sample_1.bed"
in_f2 <- "hoge4.fa"
out_f <- "output1.txt"
```

#読み込みたいBED形式フ
#マップに用いたリファレンス
#出力ファイル名を指定してout_fに格納

```
#必要なパッケージなどをロード
library(ShortRead)
```

#パッケージの読み込み

```
#入力ファイルの読み込みと必要な情報の抽出など
```

1.をやってみましょう。

入力ファイルと目的のおさらい

in_f1	start	end
contig_2	26	28
contig_2	23	25
contig_2	78	80
contig_2	79	81
contig_3	11	13
contig_3	61	63
contig_2	53	55
contig_3	28	30
contig_3	44	46
contig_1	4	6
contig_2	38	40
contig_2	51	53
contig_2	94	96
contig_3	32	34
contig_4	13	15
contig_4	15	17
contig_4	45	47

■ 入力ファイル1: sample_1.bed

- BED形式ファイル。**1列目の情報**のみを用いてコンティグ(遺伝子ID)ごとのカウント(出現回数)情報取得のために利用。

■ 入力ファイル2: hoge4.fa

- マップに用いたリファレンス配列。multi-fasta形式ファイル。**Description行**のコンティグ名(ID)の並びで結果を出力させるために利用。

```

hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACAC
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTA
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACA
    
```

出力ファイル: output1.txt

contig_1	1
contig_2	8
contig_3	5
contig_4	3

Contents

■ Rでゲノム解析

- アノテーションファイル(行列形式)からの情報抽出
 - 条件を満たす行のみの抽出。ハッシュとかgrep周辺がRで可能
- multi-fastaファイルからの情報抽出
 - 一定の長さ以上のものを抽出。コンティグごとのGC含量計算。

■ Rで(比較)トランスクリプトーム解析

- 研究目的別留意点(サンプル内比較、サンプル間比較)
- リードマッピング → 数値化(カウントデータの取得)
- 発現変動解析
 - 既存のRパッケージの弱点(unbiased DE: ○、biased DE: ×)
 - TCCパッケージを用いて二群間比較(複製あり、複製なし)



比較解析用Rパッケージ(or スクリプト)

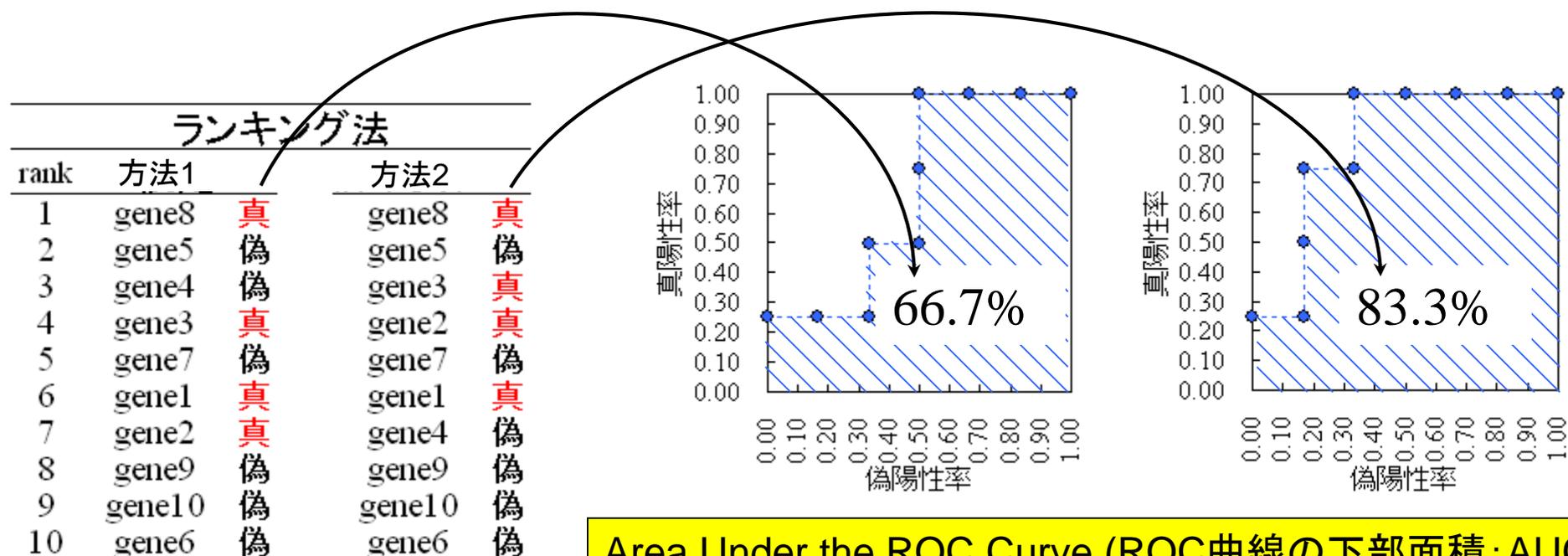
- *DEGSeq* (Wang *et al.*, *Bioinformatics*, **26**: 136-138, 2010)
- *edgeR* (Robinson *et al.*, *Bioinformatics*, **26**: 139-140, 2010)
- *GPseq* (Srivastava and Chen, *Nucleic Acids Res.*, **38**: e170, 2010)
- *baySeq* (Hardcastle and Kelly, *BMC Bioinformatics*, **11**: 422, 2010)
- *DESeq* (Anders and Huber, *Genome Biol.*, **11**: R106, 2010)
- *NBPSeq* (Di *et al.*, *SAGMB*, **10**: article24, 2011)
- *TPSM* (Auer and Doerge, *SAGMB*, **10**: article26, 2011)
- *BBSeq* (Zhou *et al.*, *Bioinformatics*, **27**: 2672-2678, 2011)
- *NOISeq* (Tarazona *et al.*, *Genome Res.*, **21**: 2213-2223, 2011)
- *PoissonSeq* (Li *et al.*, *Biostatistics*, **13**: 523-538, 2012)
- *SAMseq* (Li and Tibshirani, *Stat Methods Med Res.*, 2011 Nov 28)
- *BitSeq* (Glaus *et al.*, *Bioinformatics*, **28**: 1721-1728, 2012)
- *DEXSeq* (Anders *et al.*, *Genome Res.*, **22**: 2008-2017, 2012)
- *ShrinkBayes* (Van DE Wiel *et al.*, *Biostatistics*, **14**: 113-128, 2013)
- *sSeq* (Yu *et al.*, *Bioinformatics*, **29**: 1275-1282, 2013)
- *TCC* (Sun *et al.*, submitted)

黒字のものたち (+ α) の比較結果は...

- *DEGSeq* (Wang et al., *Bioinformatics*, 26: 136-138, 2010)
- *edgeR* (Robinson et al., *Bioinformatics*, 26: 139-140, 2010)
- *GPseq* (Srivastava and Chen, *Nucleic Acids Res.*, 38: e170, 2010)
- *baySeq* (Hardcastle and Kelly, *BMC Bioinformatics*, 11: 422, 2010)
- *DESeq* (Anders and Huber, *Genome Biol.*, 11: R106, 2010)
- *NBPSeq* (Di et al., *SAGMB*, 10: article24, 2011)
- *TPSM* (Auer and Doerge, *SAGMB*, 10: article26, 2011)
- *BBSeq* (Zhou et al., *Bioinformatics*, 27: 2672-2678, 2011)
- *NOISeq* (Tarazona et al., *Genome Res.*, 21: 2213-2223, 2011)
- *PoissonSeq* (Li et al., *Biostatistics*, 13: 523-538, 2012)
- *SAMseq* (Li and Tibshirani, *Stat Methods Med Res.*, 2011 Nov 28)
- *BitSeq* (Glaus et al., *Bioinformatics*, 28: 1721-1728, 2012)
- *DEXSeq* (Anders et al., *Genome Res.*, 22: 2008-2017, 2012)
- *ShrinkBayes* (*ShrinkSeq*; Van DE Wiel et al., *Biostatistics*, 14: 113-128, 2013)
- *sSeq* (Yu et al., *Bioinformatics*, 29: 1275-1282, 2013)
- *TCC* (Sun et al., submitted)

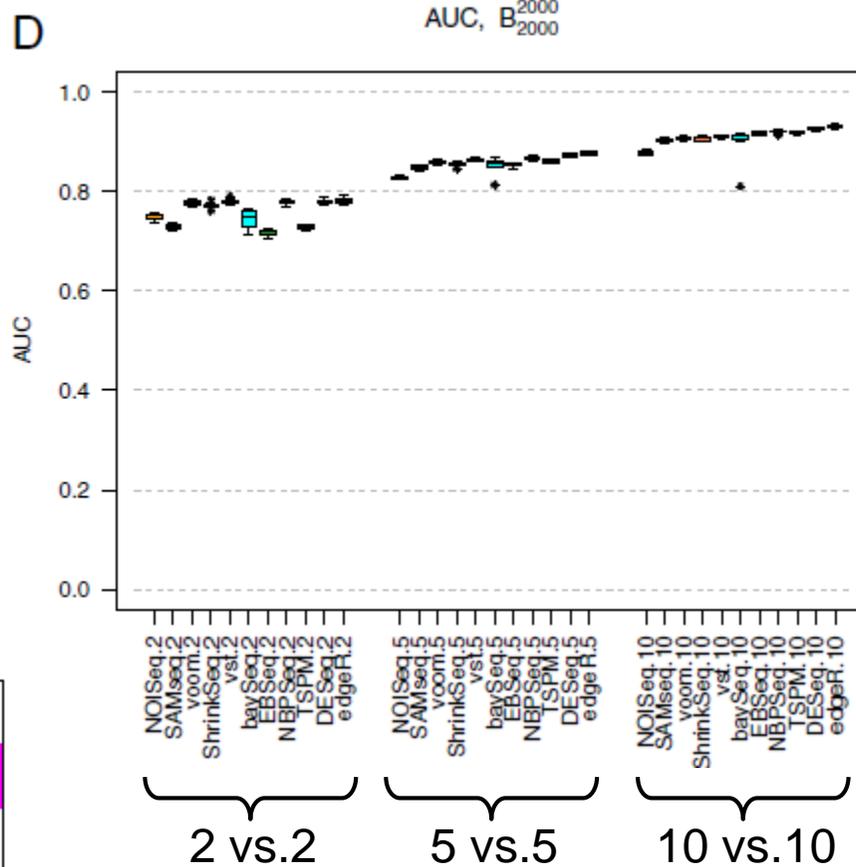
よりよい方法とは？

- その方法を用いて発現変動の度合いでランキングしたときに、**真の発現変動遺伝子 (DEG)** がより上位にランキングされる (感度・特異度高い)

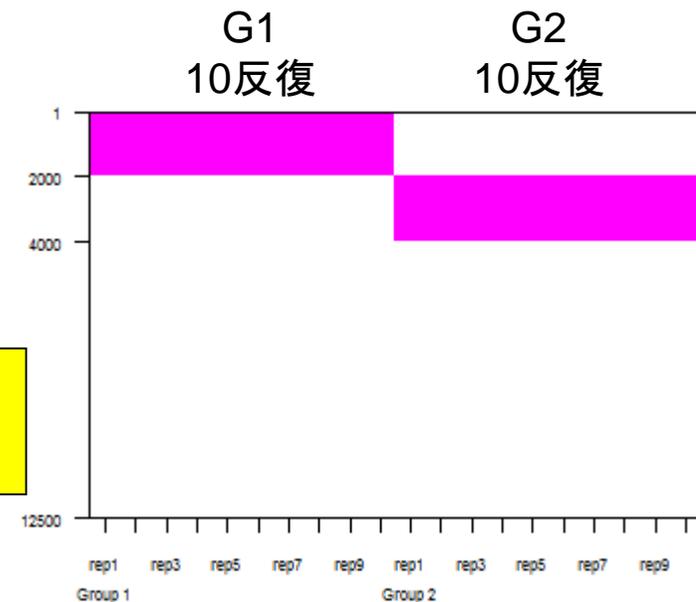
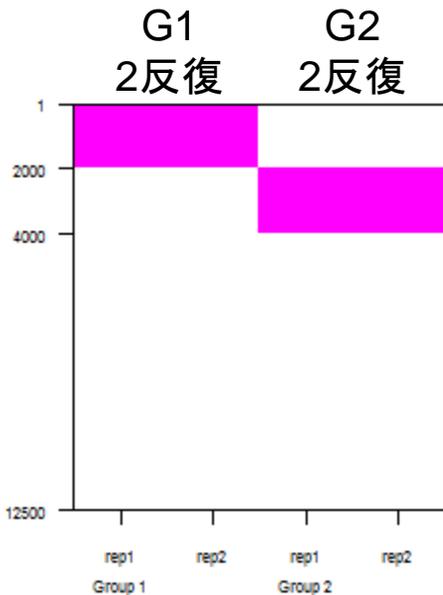


Area Under the ROC Curve (ROC曲線の下部面積: AUC)
 バイオインフォマティクス分野でよく用いられる評価基準です

AUC値の比較結果

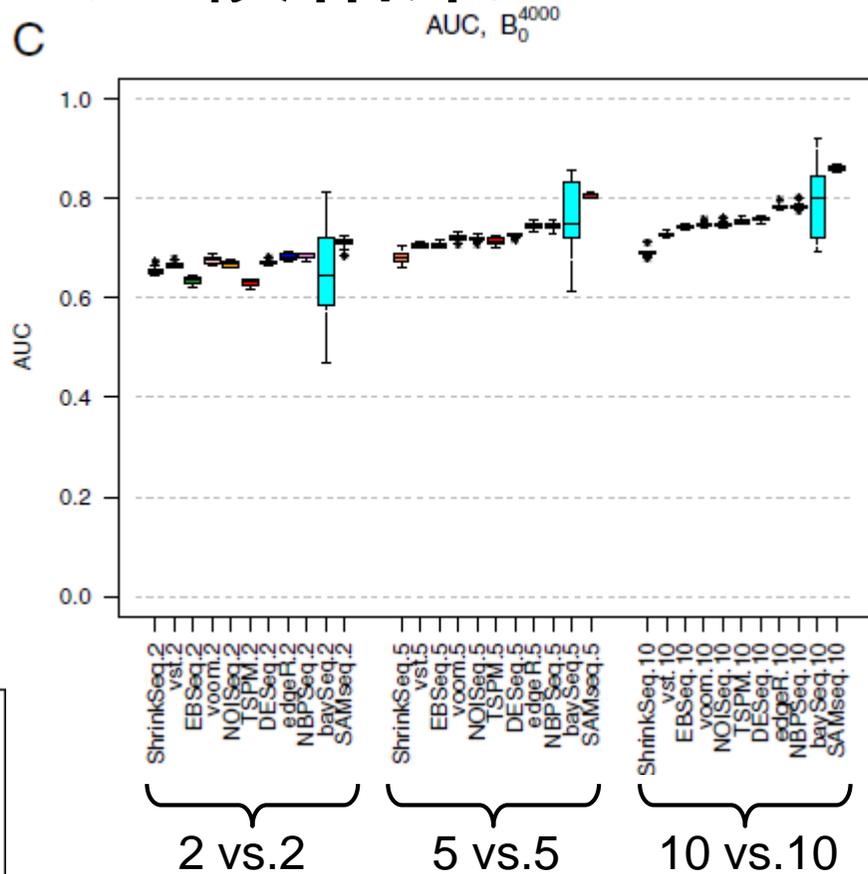


シミュレーション条件: G1 vs. G2
 全遺伝子数: 12500
 発現変動遺伝子(DEG)数: 4000
 G1で高発現: 2000
 G2で高発現: 2000
unbiased DE situation

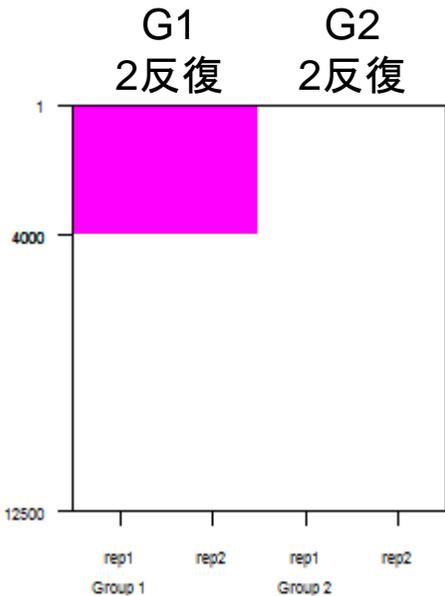


反復実験数を増やすほど精度は上がる
 (これが言いたいわけではない...)

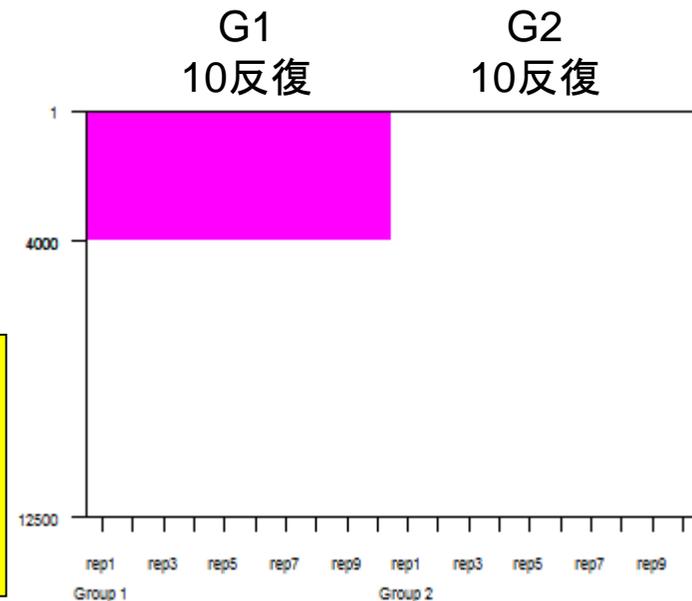
AUC値の比較結果



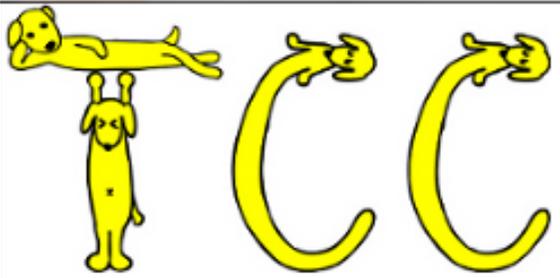
シミュレーション条件: G1 vs. G2
 全遺伝子数: 12500
 発現変動遺伝子(DEG)数: 4000
 G1で高発現: 4000
 G2で高発現: 0
biased DE situation



グループ(群)間でDEG数の組成に偏りがあると精度が大幅に低下する
 理由: データ正規化法がDEG数の組成に偏りが無いことを想定しているため



比較



TCC: an R package for comparing tag count data with robust normalization strategies

The R package, **TCC** provides users with a robust and accurate framework to perform differential expression analysis of tag count data. Differential expression analysis of tag count data (such as RNA-seq) from high-throughput sequencing technologies is a fundamental means of studying gene expression. We recently developed a multi-step normalization method ([TbT](#); Kadota et al., 2012) for two-group RNA-seq data with replicates. The strategy is to remove data that are potential differentially expressed genes (DEGs) before performing the data normalization. We demonstrated that the DEG elimination strategy (called DEGES) for data normalization is essential for obtaining a well-ranked gene list in which true DEGs are top-ranked and non-DEGs are bottom ranked. **TCC** provides integrated analysis pipelines with improved data normalization steps, compared with other packages such as [edgeR](#), [DESeq](#), and [baySeq](#), by appropriately combining their functionalities.

Important note! (last modified: May 20, 2013)

While the older version ([ver. 1.1.3](#)) of this package is currently available at the CRAN repository, we are now moving it from [CRAN](#) to [Bioconductor](#). This webpage is temporal until the next release (perhaps, ver. 1.2.0) of TCC is available upon Bioconductor. The latest version available on this webpage is [ver. 1.1.99](#).

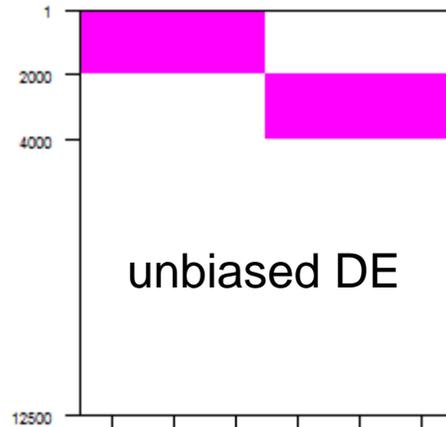
sSeq (Yu et al., *Bioinformatics*, 29: 1275-1282, 2013)

- *DEGSeq* (N)
- *edgeR* (Ro)
- *GPseq* (Sri)
- *baySeq* (H)
- *DESeq* (Ar)
- *NBPSeg* (D)
- *TPSM* (Au)
- *BBSeq* (Zh)
- *NOISeg* (T)
- *PoissonSe*
- *SAMseq* (L)
- *BitSeq* (Gla)
- *DEXSeq* (A)
- *ShrinkBaye*
- *sSeq* (Yu et al., *Bioinformatics*, 29: 1275-1282, 2013)

■ **TCC (Sun et al., submitted)**

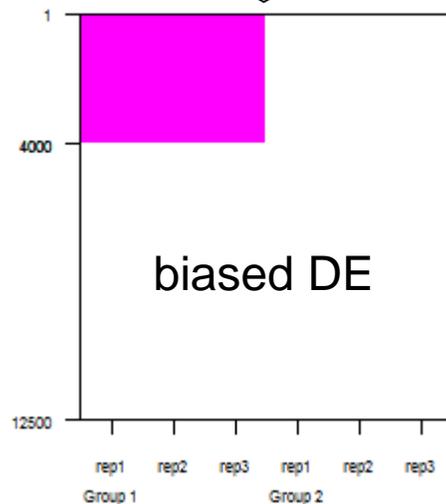
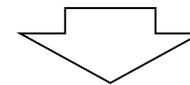
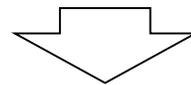
AUC値の比較結果 (*edgeR* vs. *TCC*)

3反復 vs. 3反復



edgeR
90.84%

TCC
90.83%



edgeR
82.96%

TCC
89.86%

*TCC*パッケージ中のシミュレーション解析でも*edgeR*の性能低下、および*TCC*の頑健性の高さが確認できます

TCCでの解析手順 (複製あり; 二群間比較)

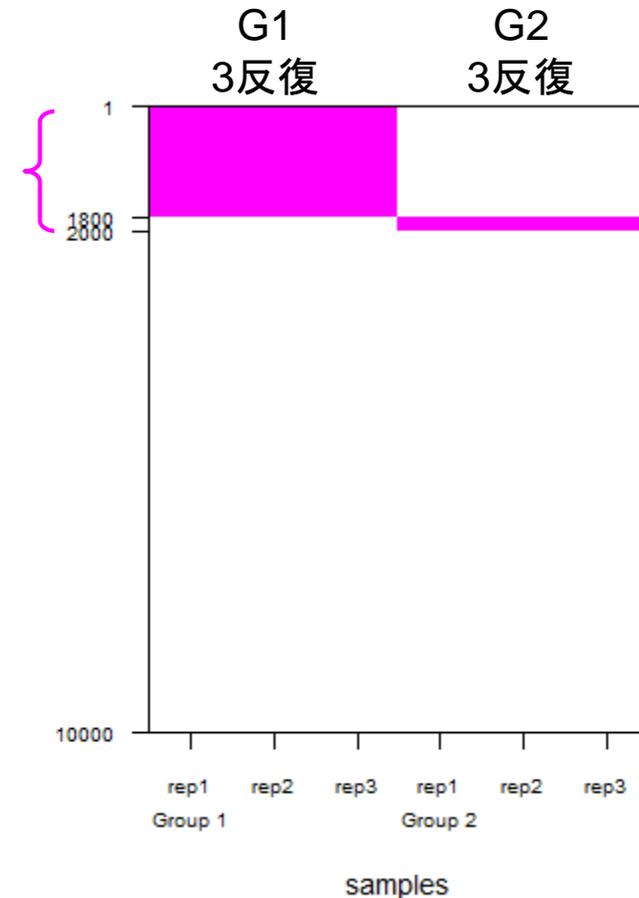
■ 仮想データ (10,000 genes × 6 samples)

□ 2,000 DEGs (20%がDEG)

- Group1 (G1)で高発現: gene1~1800 (90%)
- Group2 (G2)で高発現: gene1801~2000 (10%)

data_hypodata_3vs3.txt

	G1_rep1	G1_rep2	G1_rep3	G2_rep1	G2_rep2	G2_rep3
gene_1	36	56	144	2	1	0
gene_2	84	152	124	52	37	28
gene_3	592	840	800	151	257	200
gene_4	0	8	4	1	1	3
...						



What's new?

RパッケージTCCの最新版は1.1.99です。6月6日に開催されるNAIST植物グローバル教育プロジェクト・平成25年度ワークショップでは、R(ver. 3.0.1)とTCC(ver. 1.1.99)をベースに話を予定。Rのインストールと起動を実行したあとにTCC(ver. 1.1.99)のインストールも忘れずに行っておいてください。(2013/05/23)NEW

・平成25年度外部の方	・R3.0.1が	・解析	NGS(RNA-seq)	その他	負の二項分布に従うシミュレーションデータを作成する(fixed dispersion) (last modified 2011/10/31)
		・解析	NGS(RNA-seq)	その他	負の二項分布に従うシミュレーションデータを作成する(random dispersion) (last modified 2011/10/31)
		・解析	NGS(RNA-seq)	その他	負の二項分布に従うシミュレーションデータを作成する(tagwise dispersion) (last modified 2012/10/18)
		・解析	NGS(RNA-seq)	発現変動遺伝子	二群間について (last modified 2013/05/21) NEW
		・解析	NGS(RNA-seq)	発現変動遺伝子	二群間 複製なし DESeq with iDEGES/DESeq (Sun submitted) (last modified 2013/01/28) recommended
		・解析	NGS(RNA-seq)	発現変動遺伝子	二群間 複製なし DESeq (Anders 2010) (last modified 2013/01/28)
		・解析	NGS(RNA-seq)	発現変動遺伝子	二群間 複製あり edgeR with iDEGES/edgeR (Sun submitted) (last modified 2013/05/22) recommended NEW
		・解析	NGS(RNA-seq)	発現変動遺伝子	二群間 複製あり edgeR with DEGES/TbT (Kadota 2012) (last modified 2013/01/28)
		・解析	NGS(RNA-seq)	発現変動遺伝子	二群間 複製あり edgeR coupled with TbT normalization (上と同じで昔の記述) (last modified 2012/09/11)
		・解析	NGS(RNA-seq)	発現変動遺伝子	二群間 複製あり NBPSeg coupled with TbT normalization (Kadota 2012) (last modified 2012/08/24)
		・解析	NGS(RNA-seq)	発現変動遺伝子	二群間 複製あり DESeq (Anders 2010) (last modified 2013/01/28)



解析 | NGS(RNA-seq) | 発現変動遺伝子 | 二群間 | 複製あり | edgeR with iDEGES/edgeR (Sun submitted)

iDEGES/edgeR正規化(Sun submitted)を実行したのち、edgeRパッケージ中のexact testで発現変動遺伝子(Differentially expressed Genes; DEGs)検出を行うやり方を示します。全てTCC(≥ ver. 1.1.99)パッケージ内で完結します。

ちなみに、param_DEmethodのところでの"edgeR"は「DEG検出の手段としてedgeRパッケージ中の方法を利用する」ことに相当しますが、以下に示す他の二つの選択肢もあります:

- ・DEG検出はDESeqパッケージ中のnegative binomial test(したい場合: "deseq")
- ・DEG検出はbaySeqパッケージ中のempirical Bayes(したい場合: "bayseq")

また、param_FDRでは0.05を指定していますが、これは「DEGと判定したものの中でどの程度の数の偽物(本当はDEGではないのにDEGと判定されたもの)を許容するか?」というfalse discovery rate (FDRと略す)の閾値を指定するところです。0.05だと5%偽物が含まれることを許容することに相当し、一般にこの数を低く設定するほどDEGと判定する数が少なくなります。

ちなみによくp値の閾値(有意水準; significance level)と混同しがちですが、例えばp-value < 0.05は「non-DEGの中で(本当はDEGではないにもかかわらず)DEGと判定してしまうのが5%程度になるような閾値」に相当します。つまり、「分子は同じですが、分母が異なる」ということです。なぜ、p値を使わずにFDR(p値に対応する正しい用語は本当はq-value)がよく用いられるかというと、例えばある10000遺伝子からなるランダムデータの二群間比較を行って、p < 0.05の閾値を満たすものをDEGと判定するとしましょう。p値の定義からいって、大体5%(10000 × 0.05 = 500)程度がDEGと判定されます。この得られた500個が本当のDEG...ではないですよ? つまり、実データの場合で、「有意水準を満たすDEG数」が「搭載遺伝子数 × 設定した有意水準」以下であればそのデータセット中にDEGはないという判定がくだされるべき、という常識的な判断を下せなければなりません。逆に、10000遺伝子からなる実データの二群間比較を行って、p < 0.05の閾値を満たすものが1200個あったとしましょう。この中には偶然にDEGと判定されたものが500個程度含まれると考えるべきですが、残りの700個程度は本物と考えていいですよ。言い換えると1200個中(500/1200) × 100 = 41.6667%がfalse positiveということであり、「FDR p < 0.416667を満たす遺伝子数は1200個あった」という言い方ができます。実用上は、上記事柄が理解できていれば十分でしょう。

つまり
「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ペ。

```
1. サンプルデータ13の10,000 genes × 6 samplesの「複製あり」タグカウントデータ(data_hypodata_3vs3.txt)
Biological replicatesを模倣したシミュレーションデータ(G1群3サンプル vs. G2群3サンプル)です。
gene_1~gene_2000までがDEG (最初の1800個がG1群で高発現、残りの200個がG2群で高発現)
gene_2001~gene_10000までがnon-DEGであることが既知です。
```

```
----- ここから -----
in_f <- "data_hypodata_3vs3.txt"
out_f <- "data_hypodata_3vs3_iDEGESedgeR_DE.txt"
param_A <- 3
param_B <- 3
param1 <- 3
param_DEmethod <- "edgeR"
param_FDR <- 0.05
```

- #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
- #A群(G1群)のサンプル数を指定
- #B群(G2群)のサンプル数を指定
- #正規化におけるTMM-(edgeR-TMM)_nパイプラインのnの値(iterationの回数)を指定(デフォルト)
- #(正規化後に実行するDEG検出法を指定
- #DEG検出時のfalse discovery rate (FDR)閾値を指定

1. サンプルデータ13の10,000 genes×6 samplesの「複製あり」タグカウントデータ(data_hypodata_3vs3.txt) Biological replicatesを模倣したシミュレーションデータ(G1群3サンプル vs. G2群3サンプル)です。 gene_1~gene_2000までがDEG (最初の1800個がG1群で高発現、残りの200個がG2群で高発現) gene_2001~gene_10000までがnon-DEGであることが既知です。

```
----- ここから -----
in_f <- "data_hypodata_3vs3.txt"
out_f <- "data_hypodata_3vs3_iDEGESedgeR_DE.txt"
param_A <- 3
param_B <- 3
param_C <- 3
param_DEmethod <- "edgeR"
param_FDR <- 0.05
```

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
G1群(G1群)のサンプル数を指定

```
#必要なパッケージなどをインストール
library(TCC)

#発現データの読み込みとTCCオブジェクトの作成
data <- read.table(in_f, as.is=T)
data.cl <- c(rep(1, param_A), rep(2, param_B))
tcc <- new("TCC", data, data.cl)

#iDEGES/edgeR正規化の実行
tcc <- calcNormFactors(tcc)

#DEG検出の実行と結果の抽出
tcc <- estimateDE(tcc, test.method="tmm")
result <- getResult(tcc, sort.by="FDR")

#結果をまとめたものをファイルに出力
tmp <- cbind(rownames(tcc$count), tcc$count)
write.table(tmp, out_f, sep="\t", append=T)

#以下は (こんなこともできますという) おまけ
#正規化後のデータでM-A plotを描画
plot(tcc, FDR=param_FDR)

----- ここまで -----
```

- 切り取り(T)
- コピー(C)
- 貼り付け
- すべて選択(S)
- 印刷(I)...
- 印刷プレビュー
- Bingでマップ
- Bingで翻訳
- Googleで検索
- 電子メール
- すべてのアプリケーション
- OneNoteに追加

RGui (64-bit)
_ □ ×

ファイル 履歴 サイズ変更 ウインドウ

📄 🖨️ 🖼️

R Console

```
>
> #iDEGES/edgeR正規化の実行
> tcc <- calcNormFactors(tcc, norm.method="tmm")
TCC::INFO: Calculating normalization factors using the TMM method.
TCC::INFO: (iDEGES pipeline : tmm)
TCC::INFO: Done.
> #DEG検出の実行と結果の抽出
> tcc <- estimateDE(tcc, test.method="tmm")
TCC::INFO: Identifying DE genes using the edgeR pipeline.
TCC::INFO: Done.
> result <- getResult(tcc, sort.by="FDR")
> #結果をまとめたものをファイルに出力
> tmp <- cbind(rownames(tcc$count), tcc$count)
> write.table(tmp, out_f, sep="\t", append=T)
> #以下は (こんなこともできますという) おまけ
> #正規化後のデータでM-A plotを描画
> plot(tcc, FDR=param_FDR)
> |
```

R Graphics: Device 2 (ACTIVE)
_ □ ×

MA plot

コピー

TCCでの解析手順 (複製なし; 二群間比較)

■ 仮想データ (10,000 genes × 2 samples)

□ 2,000 DEGs (20%がDEG)

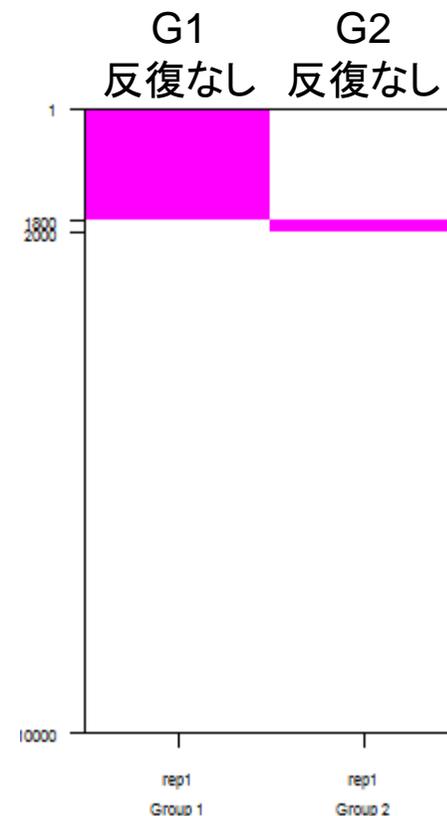
- Group1 (G1)で高発現: gene1~1800 (90%)
- Group2 (G2)で高発現: gene1801~2000 (10%)

data_hypodata_1vs1.txt

	G1_rep1	G2_rep1
gene_1	36	2
gene_2	84	52
gene_3	592	151
gene_4	0	1
...		

data_hypodata_3vs3.txt

	G1_rep1	G1_rep2	G1_rep3	G2_rep1	G2_rep2	G2_rep3
gene_1	36	56	144	2	1	0
gene_2	84	152	124	52	37	28
gene_3	592	840	800	151	257	200
gene_4	0	8	4	1	1	3
...						



What's new?

RパッケージTCCの最新版は1.1.99です。6月6日に開催されるNAIST植物グローバル教育プロジェクト・平成25年度ワークショップでは、R(ver. 3.0.1)とTCC(ver. 1.1.99)をベースに話を予定。Rのインストールと起動を実行したあとにTCC(ver. 1.1.99)のインストールも忘れずに行っておいてください。(2013/05/23)NEW

・平成25年度外部の方	・解析	NGS(RNA-seq)	その他	負の二項分布に従うシミュレーションデータを作成する(fixed dispersion)	(last modified 2011/10/31)
・R3.0.1が	・解析	NGS(RNA-seq)	その他	負の二項分布に従うシミュレーションデータを作成する(random dispersion)	(last modified 2011/10/31)
	・解析	NGS(RNA-seq)	その他	負の二項分布に従うシミュレーションデータを作成する(tagwise dispersion)	(last modified 2012/10/18)
	・解析	NGS(RNA-seq)	発現変動遺伝子	二群間	複製なし
	・解析	NGS(RNA-seq)	発現変動遺伝子	二群間	複製あり
	・解析	NGS(RNA-seq)	発現変動遺伝子	二群間	複製あり

解析 | NGS(RNA-seq) | 発現変動遺伝子 | 二群間 | 複製なし | DESeq with iDEGES/DESeq (Sun_submitted)

iDEGES/DESeq正規化(Sun submitted)を実行したのち、DESeqパッケージ中のnegative binomial testで発現変動遺伝子(Differentially expressed Genes; DEGs)検出を行うやり方を示します。全てTCC(≥ ver. 1.1.99)パッケージ内で完結します。ちなみに、param_DEmethodのところでの"deseq"は「(正規化後に実行する)DEG検出の手段としてDESeqパッケージ中の方法を利用すること」に相当しますが、以下に示す別の選択肢もあります:

- ・DEG検出はbaySeqパッケージ中のempirical Bayesにした場合:"bayseq"

「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。

1. サンプルデータ14の10,000 genes×2 samplesの「複製なし」タグカウントデータ(data_hypodata_lvs1.txt)シミュレーションデータ(G1群1サンプル vs. G2群1サンプル)です。
 gene_1~gene_2000までがDEG(最初の1800個がG1群で高発現、残りの200個がG2群で高発現)
 gene_2001~gene_10000までがnon-DEGであることが既知です。

```

----- ここから -----
in_f <- "data_hypodata_lvs1.txt"
out_f <- "data_hypodata_lvs1_iDEGESDESeq_DE.txt"
param_A <- 1
param_B <- 1
param_I <- 3
param_DEmethod <- "deseq"
param_FDR <- 0.05
  
```

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
 #A群(G1群)のサンプル数を指定
 #B群(G2群)のサンプル数を指定
 #DESeq-(DESeq-DESeq)_nパイプラインのnの値(iterationの回数)を指定
 # (正規化後に実行する)DEG検出法を指定
 #DEG検出時のfalse discovery rate (FDR)閾値を指定

#必要なパッケージなどをロード

```
library(TCC)
```

#パッケージの読み込み

#発現データの読み込みとTCCクラスオブジェクトの作成

```

data <- read.table(in_f, header=TRUE, row.names=1, sep="t", quote="")
data.cl <- c(rep(1, param_A), rep(2, param_B))
tcc <- new("TCC", data, data.cl)
  
```

#発現データファイルの読み込み
 #A群を1、B群を2としたベクトルdata.clを作成
 #TCCクラスオブジェクトtccを作成

#iDEGES/DESeq正規化の実行

```

tcc <- calcNormFactors(tcc, norm.method="deseq", test.method="deseq",
iteration=param_I, FDR=0.1, floorPDEG=0.05)
  
```

#正規化を実行した結果をtccに格納
 #正規化を実行した結果をtccに格納

#DEG検出の実行と結果の抽出

```
tcc <- estimateDE(tcc, test.method=param_DEmethod, FDR=param_FDR)
```

#DEG検出を実行した結果をtccに格納



TCC: an R package for comparing tag count data with robust strategies

The R package, **TCC** provides users with a robust and accurate framework to perform differential expression analysis of tag count data (such as RNA-seq) from various technologies is a fundamental means of studying gene expression. We recently developed a new method ([TbT](#); Kadota et al., 2012) for two-group RNA-seq data with replicates. The strategy to identify potential differentially expressed genes (DEGs) before performing the data normalization and elimination strategy (called DEGES) for data normalization is essential for obtaining a reliable list of DEGs. **TCC** provides integrated analysis of DEGs are top-ranked and non-DEGs are bottom ranked. **TCC** provides integrated analysis of normalization steps, compared with other packages such as [edgeR](#), [DESeq](#), and [baySeq](#) functionalities.

Important note! (last modified: May 22, 2013)

While the older version ([ver. 1.1.3](#)) of this package is currently available at the CRAN repository from [CRAN](#) to [Bioconductor](#). This webpage is temporal until the next release (perhaps, in the future) to [Bioconductor](#). The latest version available on this webpage is [ver. 1.1.99](#).

Installation

To install the latest version ([ver. 1.1.99](#)) of this package, download the [source file](#) and extract it, then start R:

```
install.packages("TCC_1.1.99.tar.gz", repos = NULL, type = "source")
```

Note that you need to enter the following commands if those packages have not been installed in your R environment:

```
source("http://bioconductor.org/biocLite.R")
biocLite(c("edgeR", "baySeq", "DESeq", "ROC"))
```



- [User's Guide \(vignette\)](#) [R script](#) [Manual](#)

Contents		Page
1	Introduction	3
1.1	Installation	3
1.2	Citations	4
1.3	Quick start	5
2	Preparations	7
2.1	Reading the count data	7
2.2	Constructing TCC class object	7
2.3	Filtering low-count genes (optional)	9
3	Normalization	10
3.1	Normalization of two-group count data with replicates	10
3.1.1	DEGES/TbT	10
3.1.2	DEGES/edgeR	11
3.1.3	iDEGES/edgeR	13
3.1.4	DEGES/DESeq	13
3.2	Normalization of two-group count data without replicates	14
3.3	Normalization of multi-group count data with replicates	16
3.3.1	DEGES/TbT	17
3.3.2	DEGES/edgeR	17
3.3.3	DEGES/DESeq	19
3.4	Retrieving normalized data	20
3.4.1	Retrieving two-group DEGES/edgeR-normalized data with replicates	22
3.4.2	Retrieving two-group DEGES/DESeq-normalized data with replicates	23
3.4.3	Retrieving two-group DEGES/DESeq-normalized data without replicates	24
3.4.4	Retrieving multi-group iDEGES/edgeR-normalized data with replicates	26
4	Differential expression (DE)	29
4.1	DE analysis for two-group data with replicates	29
4.1.1	edgeR coupled with iDEGES/edgeR normalization	29
4.1.2	baySeq coupled with iDEGES/edgeR normalization	30
4.2	DE analysis for two-group data without replicates	32
4.3	DE analysis for multi-group data with replicates	33
4.3.1	baySeq coupled with DEGES/edgeR normalization	33
4.3.2	edgeR coupled with DEGES/edgeR normalization	35
4.3.3	DESeq coupled with DEGES/edgeR normalization	37
5	Generation of simulation data	43
5.1	Introduction and basic usage	43
5.2	Two-group data without replicates	44
5.3	Multi-group data with and without replicates	44
5.4	Other utilities	48
6	Session info	52
7	References	53

二群間比較以外にも対応しています



まとめ

Rでいろいろできます

■ Rでゲノム解析

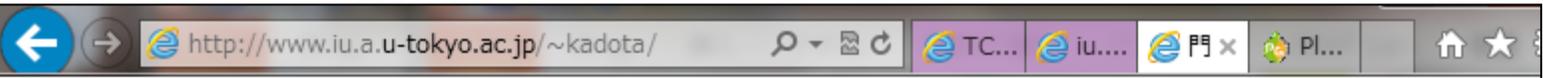
- アノテーションファイル(行列形式)からの情報抽出
 - 条件を満たす行のみの抽出。ハッシュとかgrep周辺がRで可能
- multi-fastaファイルからの情報抽出
 - 一定の長さ以上のものを抽出。コンティグごとのGC含量計算。

■ Rで(比較)トランスクリプトーム解析

- 研究目的別留意点(サンプル内比較、サンプル間比較)
- リードマッピング → 数値化(カウントデータの取得)
- 発現変動解析
 - 既存のRパッケージの弱点(unbiased DE: ○、biased DE: ×)
 - TCCパッケージを用いて二群間比較(複製あり、複製なし)



スライドPDFはウェブから取得可能です



門田 幸二のホームページ

- 名前
門田 幸二(かどた こと)
- 所属
東京大学 大学院農学
- 身分
特任准教授
- 研究分野
バイオインフォマティクス
- 所属学会
日本バイオインフォマティクス学会
日本分子生物学会
- 研究テーマ (last modified: 2013.04.18)
トランスクリプトーム解析などによって得られた生命科学への応用を目的として、トランスクリプトーム解析でひとまじりマイクロアレイデータ解析を併用して述べています。
・次世代シーケンサー
1. 二群間比較データをもとに、トランスクリプトーム解析とマイクロアレイデータを組み合わせ、トランスクリプトーム解析を実装したR

講演など (上記講義以外) (last modified: 2013.04.18) NEW

24. 題目:「ゲノム・トランスクリプトームの各種解析をRで行う」, NAIST植物グローバル教育プロジェクト・平成25年度ワークショップ, 奈良先端科学技術大学院大学(奈良), 2013.06.06
23. 題目:「食品機能解析研究とバイオインフォマティクス」, 日本農芸化学会2013年度大会・シンポジウム4SY08, 東北大学(宮城), 2013.03.27
22. 題目:「Rでトランスクリプトーム解析」, HPCI チュートリアルセミナー, 生命情報工学研究センター(東京), 2013.03.07
21. 題目:「Rでトランスクリプトーム解析」, HPCI チュートリアルセミナー, 生命情報工学研究センター(東京), 2012.03.09
20. 題目:「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～」, Rでつなぐ次世代オミックス
19. 題目:「トランスクリプトーム解析の今昔:なぜマイクロアレイ?なぜRNAシーケンス?」, Illumina Webinar Series・RNAシーケンスを始めよう・セッション3:データ解析
18. 題目:「トランスクリプトーム解析の今昔:なぜマイクロアレイ?なぜRNAシーケンス?」, Illumina Webinar Series・RNAシーケンスを始めよう・セッション1:実験手法, イルミナ株式会社(東京), 2011.9.8
17. 題目:「RNA-Seqデータ解析における正規化法の選択:RPKM値でサンプル間比較は危険?」, イルミナ株式会社バイオインフォマティクス講習会【中級】, 富士ソフトウェア アキバプラザ(東京), 2011.9.29
16. 題目:「トランスクリプトーム解析の今昔:なぜマイクロアレイ?なぜRNAシーケンス?」, Illumina Webinar Series・RNAシーケンスを始めよう・セッション1:実験手法, イルミナ株式会社(東京), 2011.9.8
15. 題目:「RNAseqによる定量的解析とqPCR, マイクロアレイなどとの比較」, 新学術領域研究「複合適応形質進化の遺伝子基盤解明」・複合適応形質進化インフォマティクスオープンセミナー, 金沢大学(石川), 2010.12.28
14. 題目:「トランスクリプトーム解析の手段としての次世代シーケンサーデータ解析～実演を中心に～」, 日本バイオインフォマティクス学会・第2回アグリバイオインフォマティクス研究会第2部, 琉球大学(沖縄), 2010.9.17
13. 題目:「トランスクリプトーム解析におけるバイオインフォマティクス要素技術～私の相場観～」, 日本バイオインフォマティクス学会・第2回アグリバイオインフォマティクス研究会第1部, 琉球大学(沖縄), 2010.9.17
12. 題目:「マイクロアレイデータ解析結果の正しい?解釈について」, 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プログラム・マイクロアレイデータ解析講習会, 東京大学(東京), 2009.11.20 (第1回), 2009.11.24 (第2回)
11. 題目:「トランスクリプトームデータの解析戦略とその周辺」, 東京大学大学院農学生命科学研究科第36回アグリバイオインフォマティクスセミナー, 東京大学(東京), 2009.10.21
10. 題目:「マイクロアレイを用いた遺伝子発現解析」, 基礎生物学研究所 バイオインフォマティクス・トレーニングコース

分布、モデル、実験デザイン、なぜ倍率変化がダメかななどの補完的な内容を含むものはこちら

・名前
門田 幸二(かどた こうじ)

・所属
[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](#)

・身分
特任准教授

・研究分野
バイオインフォマティクス、特にトランスクリプトーム解析

・所属学会
[日本バイオインフォマティクス学会](#)
[日本分子生物学会](#)

・研究テーマ (last modified: 2013.05.23)
トランスクリプトーム解析手法の開発。本ユニットでは、マイクロアレイ、一次元
サーなどによって得られる 様々なトランスクリプトームデータの解析や新規解析
生命科学への応用を目指します。「数式を並べ立てた難解な方法を凌駕する」
ットーとしています。これまでの主な研究成果を三つのカテゴリーで分けていま
プトーム解析」でひとまとめにできます。また、実験系の方でも気軽に研究成果
マイクロアレイデータ解析」と「(Rで)塩基配列解析」上にも下記開発手法中の一部
述しています。

・次世代シーケンサー(NGS)解析関連:

1. 二群間比較用データ正規化法。RNA-seqに代表されるダイナミック
ータをもとに二群間比較を行うときに、既存のRパッケージ中のデフ
を組合せたほうがよりよい結果が得られることを示した(TbT: [Kadot](#)
実装したRパッケージ [TOC](#), ver. 1.1.99が最新です！)



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



[受講生の方へ](#)



[研究者の方へ](#)

[ホーム](#) > [教育プログラム](#) > [各講義のページ](#) > 11. 農学生命情報科学特論I



11. 農学生命情報科学特論I

概要

次世代シーケンサーの普及により、以前は主にゲノム解析系で必要とされていた(塩基)配列解析のためのスキルがトランスクリプトーム解析においても要求される時代になりつつあります。本科目では、様々な局面で応用可能な配列解析系のスキルアップを目指し、RNAシーケンス(RNA-Seq)に基づく(非モデル生物の)トランスクリプトーム解析を題材とした実習を含む講義を行います。

担当教員

門田幸二 (東大・農・アグリバイオ)
西 達也 (株式会社ジナリス、東大・農・アグリバイオ / 特任教授)

お知らせ

講義日程が6月26日から6月27日になっています。ご注意ください。(平成25年4月5日更新)。
6/27, 7/3, 7/4の講義では、Rの様々なパッケージを利用します。持ち込み用PC利用希望者はRのインストールと起動を参考にして必要なパッケージをインストールしておいてください。

講義日程 (平成25年度)

1. 平成25年06月19日17:15-20:30 (PC不使用) 講師: 西達也
2. 平成25年06月27日17:15-20:30 (PC使用) 講師: 門田幸二 **日程変更しました!**
3. 平成25年07月03日17:15-20:30 (PC使用) 講師: 門田幸二
4. 平成25年07月04日17:15-20:30 (PC使用) 講師: 門田幸二

より詳細はこの講義科目で話します。外部の方も受講可能です。

謝辞

共同研究者

清水 謙多郎 先生(東京大学・大学院農学生命科学研究科)

西山 智明 先生(金沢大学・学際科学実験センター)

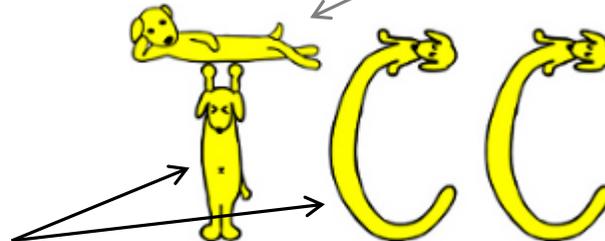
孫 建強 氏(東京大学・大学院農学生命科学研究科・大学院生)

グラント

- 基盤研究(C)(H24-26年度):「シーケンスに基づく比較トランスクリプトーム解析のためのガイドライン構築」(代表)
- 新学術領域研究(研究領域提案型)(H22年度-):「非モデル生物におけるゲノム解析法の確立」(分担;研究代表者:西山智明)

挿絵やTCCのロゴなど

(有能な秘書の)三浦 文さま作



(妻の)門田 雅世さま作

