

# ゲノム情報解析基礎

## ～ Rで塩基配列解析～

東京大学大学院農学生命科学研究科

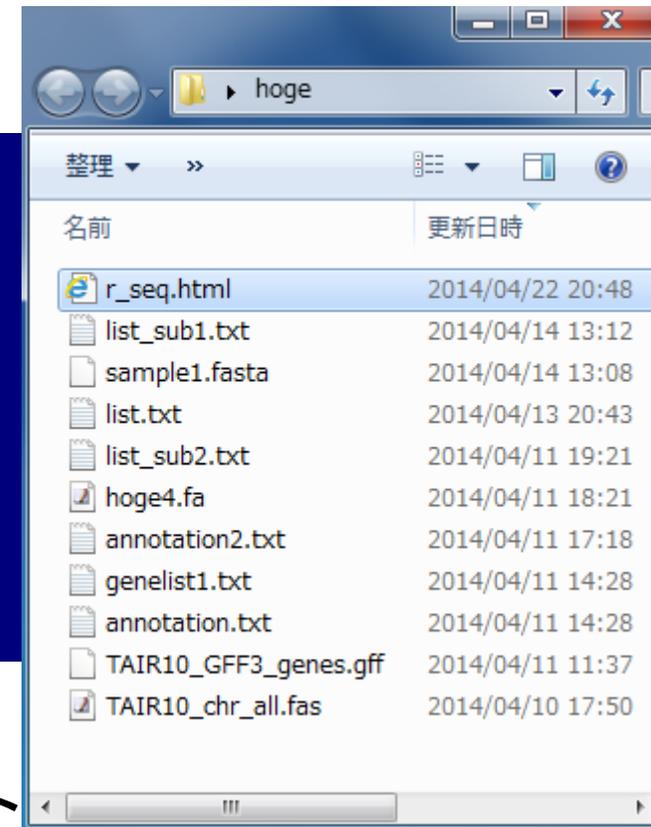
アグリバイオインフォマティクス教育研究ユニット

門田 幸二(かどた こうじ)

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

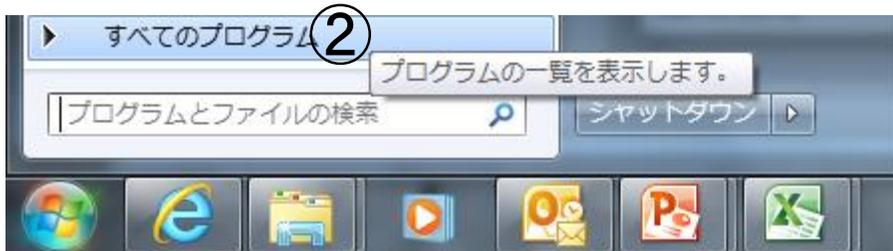
[kadota@iu.a.u-tokyo.ac.jp](mailto:kadota@iu.a.u-tokyo.ac.jp)

講義室後ろにあるUSBメモリ  
中のhogeフォルダをデスクトップ  
にコピーしておいてください。



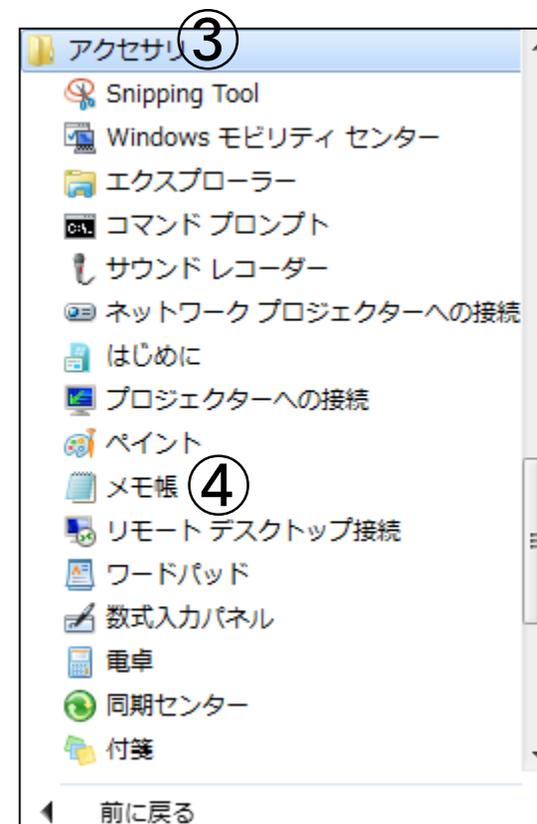
前回のhogeフォルダがデスク  
トップに残っているかもしれ  
ないのでご注意ください。

# 各種ソフトの場所



①

スタート – すべてのプログラム – ...



# ゲノム情報解析

## ■ 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

# (Rで)塩基配列解析

~NGS、RNA-seq、ゲノム、トランスクリプトーム、正規化、発現変動、統計、モデル、バイオインフォマティクス~  
(last modified 2014/04/10, since 2010)

## What's new

- 2014年9月農で開催
- 門田幸二解析を傍向かに、置いた補
- 参考資料
- 私の所属
- フォ関連
- 東大以外
- ノム情報
- 興味ある
- 機能解析
- 解析のみ

- イントロ | 一般 | [逆鎖\(reverse\)を取得](#) (last modified 2013/06/14)
- イントロ | 一般 | [2連続塩基の出現頻度情報を取得](#) (last modified 2014/02/07)
- イントロ | 一般 | [3連続塩基の出現頻度情報を取得](#) (last modified 2013/06/14)
- イントロ | 一般 | [任意の長さの連続塩基の出現頻度情報を取得](#) (last modified 2013/06/14)
- イントロ | 一般 | Tips | [任意の拡張子でファイルを保存](#) (last modified 2013/09/26)
- イントロ | 一般 | Tips | [拡張子は同じで任意の文字を追加して保存](#) (last modified 2013/09/26)
- イントロ | 一般 | 配列取得 | ゲノム配列 | [公共DBから](#) (last modified 2014/04/10) **NEW**
- イントロ | 一般 | 配列取得 | ゲノム配列 | [BSgenome](#) (last modified 2014/04/01) **NEW**
- イントロ | 一般 | 配列取得 | プロモーター配列 | [公共DBから](#) (last modified 2014/04/02) **NEW**

## イントロ | 一般 | 配列取得 | ゲノム配列 | 公共DBから **NEW**

- [UCSCの Sequence and Annotation Downloads](#) (Karolchik et al., Nucleic Acids Res., 2014)
  - [ヒト; Human \(H.sapiens\)](#)
  - [ラット; Rat \(R.norvegicus\)](#)
  - [ネコ; Cat \(F.catus\)](#)
  - [ウサギ; Rabbit \(O.cuniculus\)](#)
  - [ニワトリ; Chicken \(G.gallus\)](#)
  - [イヌ; Dog \(C.familiaris\)](#)
  - [ウマ; Horse \(E.caballus\)](#)
  - ...
- [Helix Systems Scientific Databases](#) (アップデートの日付順になっている。RefSeqやESTなど様々なデータベースを一度にみられる)
- イネ: [RAP-DB](#) (Sakai et al., Plant Cell Physiol., 2013)
  - 「ダウンロード」-「Genome assemblies」の [Download](#)。IRGSP-1.0\_genome.fasta.gz (116MB程度)の圧縮ファイル。
- シロイヌナズナ: [The Arabidopsis Information Resource \(TAIR\)](#) (Lamesch et al., Nucleic Acids Res., 2012)
  - 「ダウンロード」-「Genes」-「TAIR10 genome release」-「TAIR10 chromosome files」の [TAIR10 chr\\_all.fas](#) (120MB程度)

TAIR10のゲノム配列ファイル(TAIR10\_chr\_all.fas)はここからダウンロードしました



02/15/2012 12:00午前	411,136	<a href="#">Ath miRNAs</a>	<a href="#">Chawla 20120215.xls</a>
09/29/2009 12:00午前		ディレクトリ	<a href="#">Gene families</a>
01/19/2013 12:00午前	528,450		<a href="#">Locus Primary Gene Symbol 20130117.txt</a>
02/07/2012 12:00午前		ディレクトリ	<a href="#">OLD</a>
08/23/2011 12:00午前		ディレクトリ	<a href="#">SmallRNAsCarrington</a>
10/24/2013 10:52午後		ディレクトリ	<a href="#">TAIR10 genome release</a>
02/24/2009 12:00午前		ディレクトリ	<a href="#">TAIR6 genome release</a>
08/05/2009 12:00午前		ディレクトリ	<a href="#">TAIR7 genome release</a>
11/30/2009 12:00午前		ディレクトリ	<a href="#">TAIR8 genome release</a>

10/05/2011 12:00	08/22/2012 12:00午前	5,545	<a href="#">README TAIR10.txt</a>
08/23/2011 12:00	08/22/2012 12:00午前	3,964,120	<a href="#">TAIR10-Subcellular Predictions.xlsx</a>
08/23/2011 12:00	08/23/2011 12:00午前	ディレクトリ	<a href="#">TAIR10 NCBI mapping files</a>
08/23/2011 12:00	08/23/2011 12:00午前	792,935	<a href="#">TAIR10 TAIRAccessionID AGI mapping.txt</a>
02/07/2012 12:00	04/13/2012 12:00午前	1,868,951	<a href="#">TAIR10 TAIRlocusaccessionID AGI mapping.txt</a>
01/30/2013 12:00	08/23/2011 12:00午前	47	<a href="#">TAIR10 blastsets</a>
10/24/2013 10:5	08/23/2011 12:00午前	ディレクトリ	<a href="#">TAIR10 chromosome files</a>
	08/27/2010 12:00午前	2,608,703	<a href="#">TAIR10 domain architectures.txt</a>
	01/16/2013 12:00午前	25,396,877	<a href="#">TAIR10 functional descriptions</a>
	11/23/2010 12:00午前	25,396,966	<a href="#">TAIR10 functional descriptions.bk</a>
	10/24/2013 10:51午後	25,874,762	<a href="#">TAIR10 functional descriptions 20130831.txt</a>
	08/23/2011 12:00午前	ディレクトリ	<a href="#">TAIR10 gene confidence ranking</a>
	08/23/2011 12:00午前	ディレクトリ	<a href="#">TAIR10 gene lists</a>
	08/23/2011 12:00午前	ディレクトリ	<a href="#">TAIR10 gene transcript associations</a>
	11/18/2010 12:00午前	30	<a href="#">TAIR10 gff3</a>
	11/23/2010 12:00午前	2,053,133	<a href="#">TAIR10 locushistory.txt</a>
	12/07/2010 12:00午前		
	08/23/2011 12:00午前	26,977,690	<a href="#">NCBI Chr1.tbl</a>
	11/23/2010 12:00午前	15,361,650	<a href="#">NCBI Chr2.tbl</a>
		18,699,824	<a href="#">NCBI Chr3.tbl</a>
		14,999,724	<a href="#">NCBI Chr4.tbl</a>
		22,424,246	<a href="#">NCBI Chr5.tbl</a>
		44	<a href="#">TAIR10 chr all.fas</a>

	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%

11/10/2010 12:00午前	26,977,690	<a href="#">NCBI Chr1.tbl</a>
11/10/2010 12:00午前	15,361,650	<a href="#">NCBI Chr2.tbl</a>
11/10/2010 12:00午前	18,699,824	<a href="#">NCBI Chr3.tbl</a>
11/10/2010 12:00午前	14,999,724	<a href="#">NCBI Chr4.tbl</a>
11/10/2010 12:00午前	22,424,246	<a href="#">NCBI Chr5.tbl</a>
08/23/2011 12:00午前	44	<a href="#">TAIR10 chr all.fas</a>

```

TAIR10_chr_all.fas x
>1 CHROMOSOME dumped from ADB: Feb/3/09 16:9; last updated: 2009-02-02↓
CCCTAAACCCTAAACCCTAAACCCTAAACCTCTGAATCCTTAATCCCTAAATCCCTAAATCTTTAAATCCTACATCCAT
GAATCCCTAAATACCTAATTCCCTAAACCCGAAACCGGTTTCTCTGG
ATCGTTTTTATGTAATTGCTTATTGTTGTGTAGATTTTTTAAAAA
TGTGGTTTTCTTTCCTTCACTTAGCTATGGATGGTTTATCTTCATTT
CATTTGGGAATGTGAGTCTCTTATTGTAACCTTAGGGTTGGTTTATC
TGTTTGGACATTTATTGTCATTCTTACTCCTTTGTGGAAATGTTTGT
TAGTTGTAGGGATGAAGTCTTCTTCTGTTGTTGTTAGCCTTGTCACTCACTCTCAATGATATGGGATGGTCTTTAG

```

TAIR10のゲノム配列ファイル (TAIR10\_chr\_all.fas) のdescription 行の記述はこんな感じです

# (Rで)塩基配列解析

~NGS、RNA-seq、ゲノム、トランスクリプトーム、正規化、発現変動、統計、モデル、バイオインフォマティクス~  
(last modified 2014/04/10, since 2010)

What's

2014  
農で  
門田  
解析  
向ナ  
置い  
参考  
私の  
フォ  
東大  
ノム  
興味  
機能  
解析

- イントロ | NGS | [様々なプラットフォーム](#) (last modified 2013/06/12)
- イントロ | NGS | [qPCRやmicroarrayなどとの比較](#) (last modified 2010/12/16)
- イントロ | NGS | [Viewer](#) (last modified 2014/01/29)
- イントロ | NGS | 配列取得 | FASTQ or SRALite | [公共DBから](#) (last modified 2014/03/27) **NEW**
- イントロ | NGS | 配列取得 | FASTQ or SRALite | [SRADB\(Zhu 2013\)](#) (last modified 2014/04/01) **NEW**
- イントロ | NGS | [アノテーション情報取得](#) | [|について](#) (last modified 2014/03/26) **NEW**
- イントロ | NGS | [アノテーション情報取得](#) | [GFF/GTF形式ファイル](#) (last modified 2014/04/10) **NEW**
- イントロ | NGS | [アノテーション情報取得](#) | [refFlat形式ファイル](#) (last modified 2013/09/25)
- イントロ | NGS | [アノテーション情報取得](#) | [biomaRt\(Durinck 2009\)](#) (last modified 2013/09/26)

## イントロ | NGS | アノテーション情報取得 | GFF/GTF形式ファイル **NEW**

多くの生物種について [Ensembl \(Flicek et al., Nucleic Acids Res., 2013\)](#) の [FTPサイト](#) から GTF形式 (GFF ver. 2) の遺伝子アノテーションファイルを得ることができます。refFlat形式同様、どの領域にどの遺伝子があるのかという座標 (Coordinates) 情報を含みます。ゲノム配列のバージョンと同じであることを確認した上で用いましょう。

### • [Ensembl \(Flicek et al., Nucleic Acids Res., 2013\)](#)

圧縮 (gzip) ファイル形式です。基本は [FTPサイト](#) です。代表的なものを以下にリストアップしています。

- ヒト: [Human \(H.sapiens\)](#)
- ラット: [Rat \(R.norvegicus\)](#)
- ネコ: [Cat \(F.catus\)](#)
- ウサギ: [Rabbit \(O.cuniculus\)](#)
- ニワトリ: [Chicken \(G.gallus\)](#)
- イヌ: [Dog \(C.familiaris\)](#)
- ウマ: [Horse \(E.caballus\)](#)
- ゼブラフィッシュ: [Zebrafish \(D. rerio\)](#)
- ...

### • イネ: [RAP-DB \(Sakai et al., Plant Cell Physiol., 2013\)](#)

- 「[ダウンロード](#)」 - 「Gene set」 - 「Gene structure and function in IRGSP-1.0\_representative\_2014-03-05.tar.gz (12.4MB程度)」の圧縮ファイルが待っています

### • シロイヌナズナ: [The Arabidopsis Information Resource \(TAIR\) \(Lamesch et al., Nucleic Acids Res., 2012\)](#)

- 「[ダウンロード](#)」 - 「Genes」 - 「TAIR10 genome release」 - 「[TAIR10 gff3](#)」の [TAIR10 GFF3 genes.gff](#) (42MB程度)

TAIR10のアノテーションファイル (TAIR10\_GFF3\_genes.gff) は  
ここからダウンロードしました

Home Help Contact About Us Login/Register

Search Browse Tools Portals Download Submit News ABRC Stocks

**The Arabidopsis Information Resource**

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 expression, DNA microarray data, and information on protein domains. The database is updated every week with new gene annotations and submissions. TAIR resources include:

- Genes
- GO and PO annotations
- Maps

Download Overview  
 ABRC Documents  
 Genes  
 GO and PO Annotations  
 Maps

Subscribe to news feed  
 Follow our Twitter feed  
 Join our Facebook group

**Breaking News**

02/15/2012 12:00午前 411,136 [Ath miRNAs Konika Chawla 20120215.xls](#)

09/29/2009 12:00午前 ディレクトリ [Gene families](#)

01/19/2013 12:00午前 528,450 [Locus Primary Gene Symbol 20120117.txt](#)

02/07/2012 12:00午前 ディレクトリ [OLD](#)

08/23/2011 12:00午前 ディレクトリ [SmallRNAsCarrington](#)

10/24/2013 10:52午後 ディレクトリ [TAIR10 genome release](#)

02/24/2009 12:00午前 ディレクトリ [TAIR6 genome release](#)

08/05/2009 12:00午前 ディレクトリ [TAIR7 genome release](#)

11/30/2009 12:00午前 ディレクトリ [TAIR8 genome release](#)

10/05/2011 12:00午前 ディレクトリ [TAIR9 genome release](#)

08/22/2012 12:00午前 5,545 [README TAIR10.txt](#)

08/22/2012 12:00午前 3,964,120 [TAIR10-Subcellular Predictions.xls](#)

08/23/2011 12:00午前 ディレクトリ [TAIR10 NCBI mapping files](#)

08/23/2011 12:00午前 792,935 [TAIR10 TAIRAccessionID AGI mapping.txt](#)

02/07/2011 04/13/2012 12:00午前 1,868,951 [TAIR10 TAIRlocusaccessionID AGI mapping.txt](#)

10/24/2011 08/23/2011 12:00午前 47 [TAIR10 blastsets](#)

08/23/2011 12:00午前 ディレクトリ [TAIR10 chromosome files](#)

08/27/2010 12:00午前 2,608,703 [TAIR10 domain architectures.tab.t10](#)

01/16/2013 12:00午前 25,396,877 [TAIR10 functional descriptions](#)

11/23/2010 12:00午前 25,396,966 [TAIR10 functional descriptions.bk](#)

10/24/2013 10:51午後 25,874,762 [TAIR10 functional descriptions 20130831.txt](#)

08/23/2011 12:00午前 ディレクトリ [TAIR10 gene confidence ranking](#)

08/23/2011 12:00午前 ディレクトリ [TAIR10 gene lists](#)

08/23/2011 12:00午前 ディレクトリ [TAIR10 gene transcript associations](#)

11/18/2010 12:00午前 30 [TAIR10 gff3](#)

11/23/2010 12:00午前 2,053,133 [TAIR10 locusshifter.txt](#)

12/07/2010 12:00午前 904 [TAIR10 sequence edits.txt](#)

08/23/2011 12:00午前 ディレクトリ [TAIR10 transposable elements](#)

11/23/2010 12:00午前 ディレクトリ [TAIR10 YAC](#)

一般的には、Community annotationのほうがより情報量も多くていいらしい

01/04/2011 12:00午前	2,662,626	<a href="#">Blist TAIR10.gff</a>
02/15/2012 12:00午前	ディレクトリ	<a href="#">Community annotation GFF</a>
03/31/2011 12:00午前	ディレクトリ	<a href="#">DNA replication origin</a>
04/06/2011 12:00午前	885	<a href="#">README TAIR10 GFF3.txt</a>
01/04/2011 12:00午前	266	<a href="#">README gbrowse data.txt</a>
01/04/2011 12:00午前	43,226,098	<a href="#">Spliced Junctions clustered.gff</a>
12/14/2010 12:00午前	44,139,005	<a href="#">TAIR10 GFF3 genes.gff</a>
12/14/2010 12:00午前	49,811,410	<a href="#">TAIR10 GFF3 genes transcribed.gff</a>
01/04/2011 12:00午前	121,054	<a href="#">TAIR10 Models obsoleted.gff</a>
01/04/2011 12:00午前	3,115,098	<a href="#">TAIR10 unconfirmed exons.gff</a>
10/18/2011 12:00午前	39,825,304	<a href="#">TAIR GFF3 ssrs.gff</a>
07/13/2011 12:00午前	134,836	<a href="#">arabidopsis.conf</a>

TAIRウェブインターフェースからアノテーションファイル (TAIR10\_GFF3\_genes.gff) を取得する際のイメージ

# GFF/GTF形式ファイルの例

## GFF3形式ファイルの例(シロイヌナズナ; TAIR10\_GFF3\_genes.gff)

	A	B	C	D	E	F	G	H	I
1	Chr1	TAIR10	chromosome	1	30427671	.	.	.	ID=Chr1;Name=Chr1
2	Chr1	TAIR10	gene	3631	5899	.	+	.	ID=AT1G01010;Note=protein_coding_gene;Name=AT1G01010
3	Chr1	TAIR10	mRNA	3631	5899	.	+	.	ID=AT1G01010.1;Parent=AT1G01010;Name=AT1G01010.1;Index=1
4	Chr1	TAIR10	protein	3760	5630	.	+	.	ID=AT1G01010.1-Protein;Name=AT1G01010.1;Derives_from=AT1G01010.1
5	Chr1	TAIR10	exon	3631	3913	.	+	.	Parent=AT1G01010.1
6	Chr1	TAIR10	five_prime_UTR	3631	3759	.	+	.	Parent=AT1G01010.1
7	Chr1	TAIR10	CDS	3760	3913	.	+	0	Parent=AT1G01010.1,AT1G01010.1-Protein;
8	Chr1	TAIR10	exon	3996	4276	.	+	.	Parent=AT1G01010.1
9	Chr1	TAIR10	CDS	3996	4276	.	+	2	Parent=AT1G01010.1,AT1G01010.1-Protein;
10	Chr1	TAIR10	exon	4486	4605	.	+	.	Parent=AT1G01010.1
11	Chr1	TAIR10	CDS	4486	4605	.	+	0	Parent=AT1G01010.1,AT1G01010.1-Protein;
12	Chr1	TAIR10	exon	4706	5095	.	+	.	Parent=AT1G01010.1

遺伝子ごとに、どの染色体のどの座標上に存在するのかなどの情報を含むタブ区切りテキストファイル

## GTF形式ファイルの例(ゼブラフィッシュ; Danio\_rerio.Zv9.75.gtf)

	A	B	C	D	E	F	G	H	I
1									#!genome-build Zv9
2									#!genome-version Zv9
3									#!genome-date 2010-04
4									#!genome-build-accession NCBI:GCA_000002035.2
5									#!genebuild-last-updated 2014-02
6	7	protein_coding	gene	100958	101715	.	+	.	gene_id "ENSDARG00000076051"; gene_name "CABZ01062994.1"; gene
7	7	protein_coding	transcript	100958	101715	.	+	.	gene_id "ENSDARG00000076051"; transcript_id "ENS DART00000113409
8	7	protein_coding	exon	100958	100975	.	+	.	gene_id "ENSDARG00000076051"; transcript_id "ENS DART00000113409
9	7	protein_coding	CDS	100958	100975	.	+	0	gene_id "ENSDARG00000076051"; transcript_id "ENS DART00000113409
10	7	protein_coding	exon	101077	101715	.	+	.	gene_id "ENSDARG00000076051"; transcript_id "ENS DART00000113409
11	7	protein_coding	CDS	101077	101715	.	+	0	gene_id "ENSDARG00000076051"; transcript_id "ENS DART00000113409
12	7	protein_coding	gene	116160	117573	.	+	.	gene_id "ENSDARG00000088691"; gene_name "BX511027.1"; gene_sour
13	7	protein_coding	transcript	116160	117573	.	+	.	gene_id "ENSDARG00000088691"; transcript_id "ENS DART00000129330

# タブ区切りテキストファイルからの情報抽出

入力: アノテーションファイル  
(`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

出力: `hoge1.txt`



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

入力: リストファイル (`genelist1.txt`)

	A
1	gene1
2	gene7
3	gene9

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NGSやRNA-seq解析)

(last modified 2014/02/06, since 2010)

・ [イントロ](#) | [一般](#) | [任意のキーワードを含む行を抽出\(基礎\)](#)

## What's new?

- 項目名の整理を行っています。3CやBS-seq周辺についても少し言及してあります。(2014/02/06) **NEW**
- 一連の解析パイプライン(RNA-seqデータ取得->マッピング->カウントデータやRPKMデータ取得->サンプル間クラスターリングや発現変動解析およびM-A plot描画まで)をアップデートしました。項目名の一番下のほうです。(2013/10/19)
- 発現変動解析用RパッケージTCC (ver. 1.2.0; [Sun et al., BMC Bioinformatics, 2013](#))がBioconductorよりリリースされました。最新版を利用したい方は、R (ver. 3.0.2)をインストールしたのち、Bioconductor (ver. 2.13)をインストールしてください。(2013/10/17)
- どのブラウザからでもエラーなく見られる(W3C validation)ように((Rで)マイクロアレイデータ解析も含めて)リニューアルしました。(2013/07/30)
- 2013年7月29日まで公開していた以前の「(Rで)塩基配列解析」のウェブページや関連ファイルは [Rdeenni.zip](#) からダウンロード可能です(110MB程度)。(2013/07/30)

- [はじめに](#) (last modified 2014/01/30) **NEW**
- [Rのインストールと起動](#) (last modified 2013/09/27)
- [サンプルデータ](#) (last modified 2014/02/06) **NEW**
- [イントロ](#) | [一般](#) | [ランダムに行を抽出](#) (last modified 2013/10/16)
- [イントロ](#) | [一般](#) | [任意の文字列を最初の挿入](#) (last modified 2013/10/16)
- [イントロ](#) | [一般](#) | [任意のキーワードを含む行を抽出\(基礎\)](#) (last modified 2013/10/16)
- [イントロ](#) | [一般](#) | [ランダムな塩基配列を生成](#) (last modified 2013/10/16)
- [イントロ](#) | [一般](#) | [任意の長さの可能な全ての塩基配列を作成](#) (last modified 2013/10/16)
- [イントロ](#) | [一般](#) | [任意の位置の塩基を置換](#) (last modified 2013/10/16)
- [イントロ](#) | [一般](#) | [指定した範囲の配列を取得](#) (last modified 2013/10/16)
- [イントロ](#) | [一般](#) | [翻訳配列\(translate\)を取得](#) (last modified 2013/10/16)
- [イントロ](#) | [一般](#) | [相補鎖\(complement\)を取得](#) (last modified 2013/10/16)
- [イントロ](#) | [一般](#) | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/10/16)
- [イントロ](#) | [一般](#) | [逆鎖\(reverse\)を取得](#) (last modified 2013/06/14)
- [イントロ](#) | [一般](#) | [2連続塩基の出現頻度情報を取得](#) (last modified 2013/06/14)
- [イントロ](#) | [一般](#) | [3連続塩基の出現頻度情報を取得](#) (last modified 2013/06/14)
- [イントロ](#) | [一般](#) | [任意の長さの連続塩基の出現頻度情報を取得](#) (last modified 2013/06/14)
- [イントロ](#) | [一般](#) | [Tips](#) | [拡張子は同じで任意の文字を追加して](#) (last modified 2013/06/14)
- [イントロ](#) | [一般](#) | [配列取得](#) | [ゲノム配列](#) | [公共DBから](#) (last modified 2013/06/14)
- [イントロ](#) | [一般](#) | [配列取得](#) | [ゲノム配列](#) | [BSgenome](#) (last modified 2013/06/14)

## イントロ | 一般 | 任意のキーワードを含む行を抽出(基礎) **NEW**

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 ([annotation.txt](#)) 中の第1列目をキーとして、リストファイル ([genelist1.txt](#)) 中のものが含まれる行全体を出力したい場合:

```

in_f1 <- "annotation.txt" #入力ファイル名を指定してin_f1に格納(アノテーションファイル)
in_f2 <- "genelist1.txt" #入力ファイル名を指定してin_f2に格納(リストファイル)
out_f <- "hogel1.txt" #出力ファイル名を指定してout_fに格納
param <- 1 #アノテーションファイル中の検索したい列番号を指定

#入力ファイルの読み込み
data <- read.table(in_f1, header=TRUE, sep="\t", quote="") #in_f1で指定したファイルの読み込み
keywords <- readLines(in_f2) #in_f2で指定したファイルの読み込み
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) #条件を満たすかどうかを判定した結果をobjに格納
out <- data[obj,] #objがTRUEとなる行のみ抽出した結果をoutに格納
dim(out) #オブジェクトoutの行数と列数を表示

#ファイルに保存
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

```

# 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(genelist1.txt)中のものが含まれる行全体を出力したい場合:

```
in_f1 <- "annotation.txt" #入力ファイル名を指定してin_f1に格納(アノテーションファイル)
in_f2 <- "genelist1.txt" #入力ファイル名を指定してin_f2に格納(リストファイル)
out_f <- "hoge1.txt" #出力ファイル名を指定してout_fに格納
param <- 1 #アノテーションファイル中の検索したい列番号を指定

#入力ファイルの読み込み
data <- read.table(in_f1, header=TRUE, sep="\t", quote="") #in_f1で指定したファイルの読み込み
keywords <- readLines(in_f2) #in_f2で指定したファイルの読み込み
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) #条件を満たすかどうかを判定した結果をobjに格納
out <- data[obj,] #objがTRUEとなる行のみ抽出した結果をoutに格納
dim(out) #オブジェクトoutの行数と列数を表示

#ファイルに保存
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
```

## 入力1: 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

## 入力2: genelist1.txt

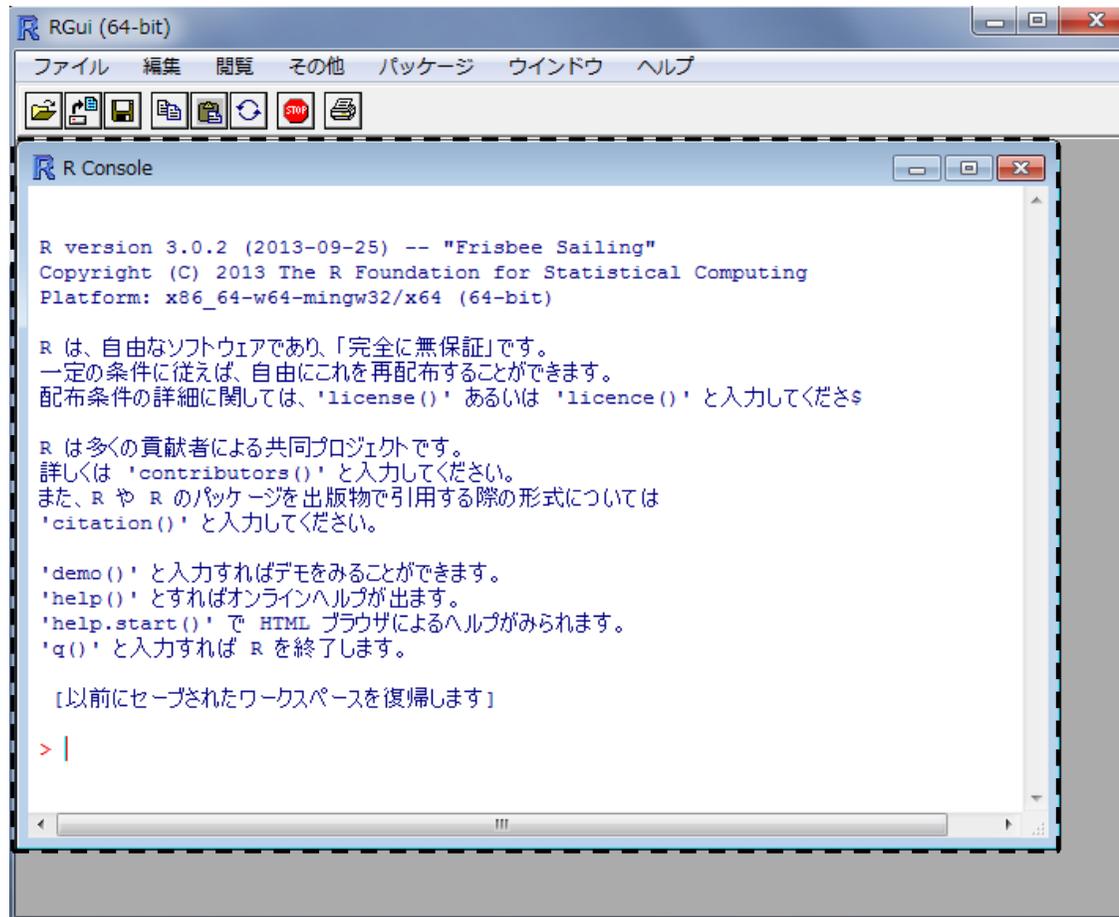
	A
1	gene1
2	gene7
3	gene9

## 出力: 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genelist1.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

# Rの起動



デスクトップにあるhogeフォルダ中のファイルを解析

# 作業ディレクトリの変更

「Windows(C:)」となっている場合もあるが、気にしない

貸与PCは「iu」

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 コードのソースを読み込み...  
新しいスクリプト  
スクリプトを開く...  
ファイルの表示...  
作業スペースの読み込み...  
作業スペースの保存...  
履歴の読み込み...  
履歴の保存...  
ディレクトリの変更... ①  
印刷...  
ファイルを保存...  
終了

作業ディレクトリの変更  
C:\

コンピューター  
ローカル ディスク (C:) ②  
SD Card (E:)

空き領域: 280 GB  
合計サイズ: 453 GB

フォルダー(F): ローカル ディスク (C:)

新しいフォルダーの作成(N) OK キャンセル

```
'help.start()'でHTMLブラウザによるヘルプ  
'q()'と入力すればRを終了します。  
> |
```

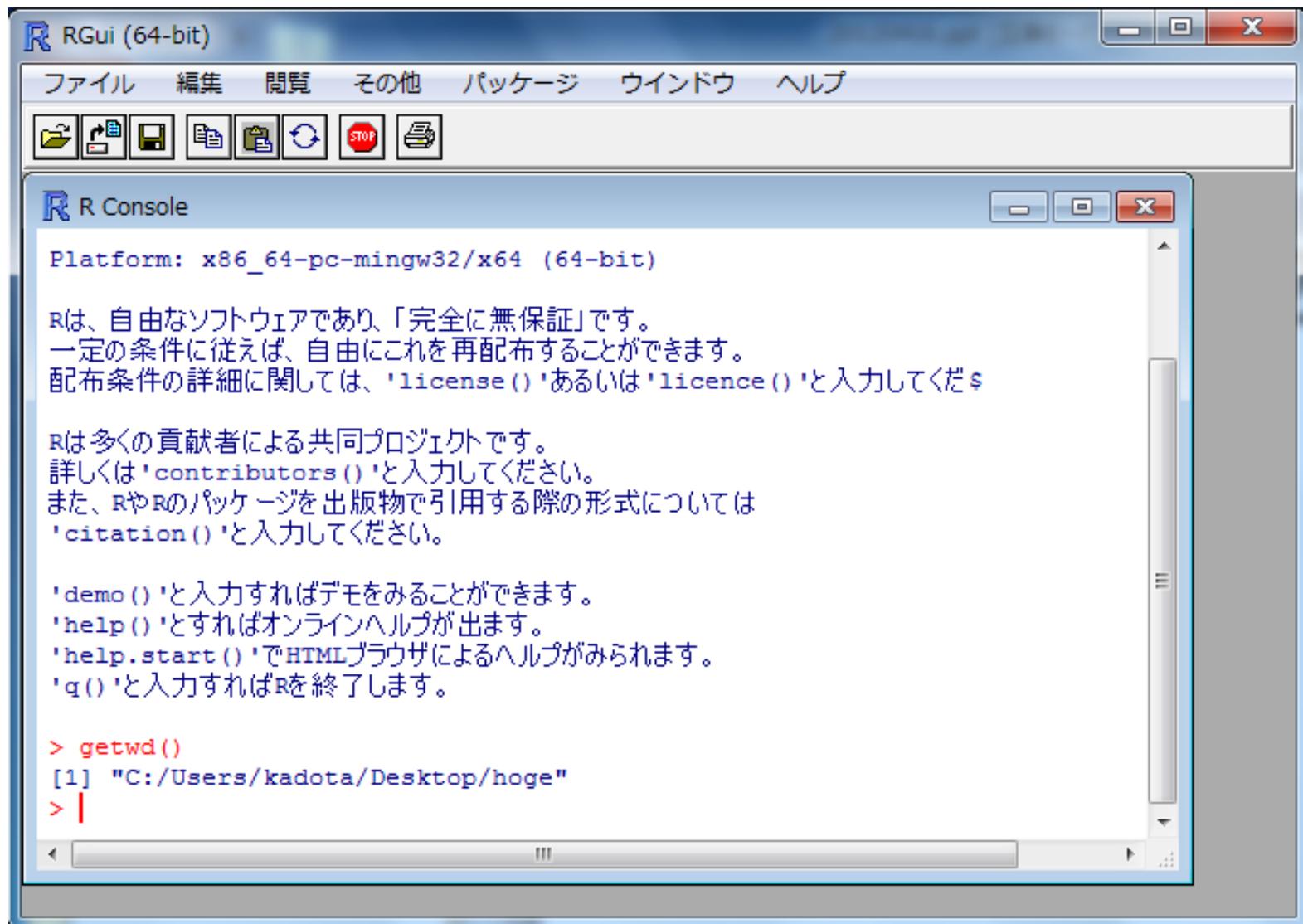
作業ディレクトリの変更  
C:\Users\kadota\Desktop\hoge

Users ③  
Default  
kadota ④  
AppData  
Dropbox  
Roaming  
アドレス帳  
お気に入り  
ダウンロード  
デスクトップ ⑤  
hoge ⑥

フォルダー(F): hoge

新しいフォルダーの作成(N) OK ⑦ キャンセル

# getwd() と打ち込んで確認



```
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 | 一般 | 任意のキーワードを含む行を抽出(基礎) **NEW**

例えばタブ区切りテキストファイルが手元があり、この中からリストファイル中の文字列を含む行を抽出するやり方 (UNIX) の `grep` コマンドのようなものであり、`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 (`annotation.txt`) 中の行全体を出力したい場合:

```
f1 <- "annotation.txt"
f2 <- "genelist1.txt"
out_f <- "hogel.txt"
param <- 1

#入力ファイルの読み込み
data <- read.table(in = f1, header = TRUE, as.is = TRUE)
keywords <- readLines(in = f2)
dim(data)

#本番
obj <- is.element(as.character(keywords), data[,1])
out <- data[obj,]
dim(out)

#ファイルに保存
write.table(out, out_f, sep = "\t", as.is = TRUE)
```



2013年7月以降のリニューアルで、コードのコピーがやりやすくなっています。**CTRLとALTキー**を押しながらコードの枠内で**左クリック**すると、**全選択**できます。

The screenshot shows the RGui (64-bit) window. The R 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には、
'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ることができます。
'q()'と入力すればRを終了します。

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

A context menu is open over the console, with the 'Copy' option selected. The menu items are:

- 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)
- ウィンドウを常にトップに置く

- 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) 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

hoge

名前	日付時刻
genelist1.txt	2012/03/28 16:41
annotation.txt	2012/03/28 16:41

## 実行後のhogeフォルダ

hoge

名前	日付時刻
hoge1.txt	2012/03/28 16:49
genelist1.txt	2012/03/28 16:41
annotation.txt	2012/03/28 16:41

# (Rで)塩基配列解析

～NGS、RNA-seq、ゲノム、トランスクリプトーム、正規化、発現変動、統計、モデル、バイオインフォマティクス～  
(last modified 2014/04/10, since 2010)

## What's new?

- 2014年9月1日～12日に「バイオインフォマティクス人材育成カリキュラム(次世代シーケンサ)速習コース」を東大農で開催します。近いうちに詳細を公開しますので興味ある方は予定を開けといてください。(2014/04/05) **NEW**
- 門田幸二 著 [シリーズ Useful R 第7巻トランスクリプトーム解析](#)が共立出版から出ました。マイクロレイとRNA-seq解析を例としてRを用いてトランスクリプトーム解析を行うための体系的な本としてまとめました。数式が苦手なヒト向けに、重みつき平均の具体的な計算例などを挙げてオプションの意味などがわかるような中身の理解に重点を置いた構成にしています。(2014/04/10) **NEW**
- [参考資料\(講義、講習会、本など\)](#)の項目を追加しました。(2014/04/10) **NEW**
- 私の所属する[アグリバイオインフォマティクス教育研究プログラム](#)では、平成26年度も(東大生に限らず)バイオインフォ関連講義を行います。4/9に私の第一回目の講義がありましたが、過去最高の122名の出席がありました。例年東大以外の企業の方、研究員、学生が二割程度は受講しております。このウェブページと直接関連する講義は「[ゲノム情報解析基礎](#)」と「[農学生命情報科学特論I](#)」ですが、背景理論の説明などは「[機能ゲノム学](#)」でも行います。興味ある科目のみの受講も可能ですので、お気軽にどうぞ。(2014/04/10) **NEW**
- 機能解析の遺伝子オントロジー(GO)解析とパスウェイ(Pathway)解析周辺を更新し、SeqGSEAパッケージを用いた解析のみですが一通りできるようにしました。(2014/03/30) **NEW**

-  [はじめに](#) (last modified 2014/01/30)
- [参考資料\(講義、講習会、本など\)](#) (last modified 2014/04/10) **NEW**
- [過去のお知らせ](#) (last modified 2014/03/23) **NEW**
- [Rのインストールと起動](#) (last modified 2013/09/27)
- [サンプルデータ](#) (last modified 2014/04/01) **NEW**
- イントロ | 一般 | [ランダムに行を抽出](#) (last modified 2013/10/10)
- イントロ | 一般 | [任意の文字列を行の最初に挿入](#) (last modified 2013/10/10)
- イントロ | 一般 | [任意のキーワードを含む行を抽出\(基礎\)](#) (last modified 2013/10/10)
- イントロ | 一般 | [ランダムな塩基配列を生成](#) (last modified 2013/10/10)
- イントロ | 一般 | [任意の長さの可能な全ての塩基配列を作成](#) (last modified 2013/10/10)
- イントロ | 一般 | [任意の位置の塩基を置換](#) (last modified 2013/10/10)
- イントロ | 一般 | [指定した範囲の配列を取得](#) (last modified 2013/10/10)
- イントロ | 一般 | [指定したID\(染色体やdescription\)の配列を取得](#) (last modified 2014/03/10)
- イントロ | 一般 | [翻訳配列\(translate\)を取得](#) (last modified 2013/06/14)
- イントロ | 一般 | [相補鎖\(complement\)を取得](#) (last modified 2013/06/14)

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

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

```

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

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

#本番
obj <- is.element(as.character(data[,param]), keywords)#条件を満たすかどうかを判定した結果をobjに格納
out <- data[obj,]
dim(out)

#ファイルに保存
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

```

#入力ファイル名を指定してin\_f1に格納(アノテーションファイル)  
 #入力ファイル名を指定してin\_f2に格納(リストファイル)  
 #出力ファイル名を指定してout\_fに格納  
 #アノテーションファイル中の検索したい列番号を指定

#in\_f2で指定したファイルの読み込み  
 #オブジェクトdataの行数と列数を表示

#objがTRUEとなる行のみ抽出した結果を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

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



# 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 Console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

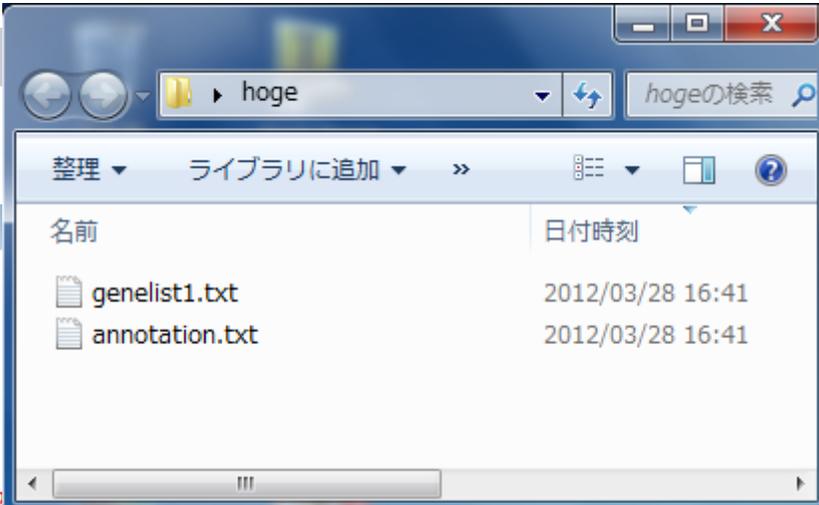
#ファイルの読み込み
data <- read.table(in_f1, header=TRUE, sep="#t", quote="")
keywords <- readLines(in_f2)
dim(data)

#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1つの  
ファイルとして保存することをお勧め

# 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

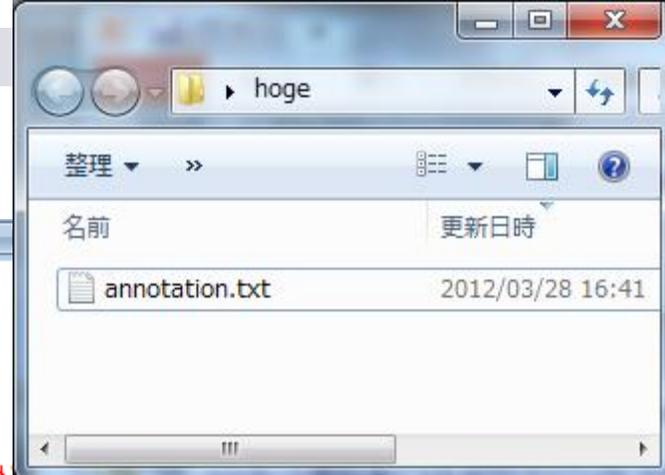
#入力ファイル\$  
#入力ファイル\$

#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)

# 改行を入れておいたほうが警告が出ない

```
list.txt x
nuclear↓
membrane↓
↓
```

```
list.txt * x
nuclear↓
membrane↓
```

```
R Console
> in_f1 <- "annotation.txt" #入力$
> in_f2 <- "list.txt" #入力$
> out_f <- "hoge2.txt" #出力$
> param <- 4 #アノ$
>
> #入力ファイルの読み込み
> data <- read.table(in_f1, header=TRUE, sep="$"
> keywords <- readLines(in_f2) #in_f$
> dim(data) #オブ$
[1] 11 4
>
> #本番
> obj <- is.element(as.character(data[,param])$
> out <- data[obj,] #obj$
> dim(out) #オブ$
[1] 7 4
>
> #ファイルに保存
> write.table(out, out_f, sep="\t", append=F, $
> |
```

```
R Console
> in_f1 <- "annotation.txt" #入力ファイル名を指定して in_f$
> in_f2 <- "list.txt" #入力ファイル名を指定して in_f$
> out_f <- "hoge2.txt" #出力ファイル名を指定して out_$
> param <- 4 #アノテーションファイル中の検$
>
> #入力ファイルの読み込み
> data <- read.table(in_f1, header=TRUE, sep="\t", quote="") #in_f1で指$
> keywords <- readLines(in_f2) #in_f2で指定したファイルの読$
警告メッセージ:
In readLines(in_f2) : 'list.txt' で不完全な最終行が見つかりました
> dim(data) #オブジェクト dataの行数と列数$
[1] 11 4
>
> #本番
> obj <- is.element(as.character(data[,param]), keywords) #条件を満たす$
> out <- data[obj,] #objがTRUEとなる行のみ抽出し$
> dim(out) #オブジェクト outの行数と列数$
[1] 7 4
>
> #ファイルに保存
> write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F) #ou$
> |
```

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(左)と入れない場合(右)の挙動の違いを把握し、後学のために警告メッセージの意味を理解しておくとい。この場合は結果には影響していないことがわかる。

# 読み込み

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

#入力ファイルの読み込み ①

```
data <- read.table(in_f1, header=TRUE, sep="\t", quote="") #in_f1で指定した
```

②

③

	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	age ne1	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一部分です
- ③ ファイルの区切り文字はタブです

# 行列data

	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", qu
>
> 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

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

#入力ファイルの読み込み

```
data <- read.table(in_f1, header=TRUE)
keywords <- readLines(in_f2)
```

```
dim(data)
```

#本番

```
obj <- is.element
out <- data[obj,]
dim(out)
```

#ファイルに保存

```
write.table(out,
```

annotation.txt

	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

R Console

```
> data
```

```

  genename accession de
1   gene1   hoge01 plasma_mem      nuclear
2   gene2   hoge02   hohinu      membrane
3   gene3   hoge03   agribio      endoplasmic
4   gene4   hoge04   genesi
5   gene5   hoge05     kam
6   gene6   hoge06   netteb
7   gene7   hoge07   tebasak
8   gene8   hoge08     biiru
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

```
R Console
10 gene10 hoge10 agene1 membrane
11 gene11 hoge11 iyaaaa endoplasmic
> dim(data)
[1] 11 4
> data[6,4]
[1] humei
Levels: endoplasmic humei membrane nuclear
> data[2,]
 gene10 accession description subcellular_location
2 gene2 hoge02 hohinu membrane
> 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	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

```
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**

# やりたかったことをおさらい

genelist1.txt

	A
1	gene1
2	gene7
3	gene9

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

#入力ファイル名を指定し  
#入力ファイル名を指定し

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

```
#本番
obj <- is.elementary(data)
out <- data[obj,]
dim(out)
```

```
#ファイルに保存
write.table(out, out_f, as.is=TRUE)
```

```
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     biiru             nuclear
9   gene9   hoge09   nihonshu            nuclear
10  gene10   hoge10     agene1
11  gene11   hoge11     iyaaaa

> 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

# 論理値ベクトルを理解

genelist1.txt

	A
1	gene1
2	gene7
3	gene9

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

#入力ファイルの読み込み

```
data <- read.table(in_f1, head=1)
keywords <- readLines(in_f2)
dim(data)
```

#本番

```
obj <- is.element(as.character(data[,param]), keywords) #条件を満たすかどうか
out <- data[obj,] #objがTRUEとなる要素のみを表示
dim(out) #オブジェクトの行数と列数を表示
```

ここも分かりづらいところではある...

#ファイルに保存

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

```
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試す

#本番

```
obj <- is.element(as.character(data[,param]), keywords) #条件を満たすかどうか
out <- data[obj,] #objがTRUEとなる行のみ抽出した結果を表示
dim(out) #オブジェクトoutの行数と列数を表示
```

#ファイルに保存

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

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([genelist1.txt](#))中のものが含まれる行全体を出力したい場合:

genelist1.txt

	A
1	gene1
2	gene7
3	gene9

```
in_f1 <- "annotation.txt" #入力ファイル名を指定してin_f1に格納(アノテーションファイル)
in_f2 <- "genelist1.txt" #入力ファイル名を指定してin_f2に格納(リストファイル)
out_f <- "hoge1.txt" #出力ファイル名を指定してout_fに格納
param <- 1 #アノテーションファイル中の検索したい列番号を指定
```

```
#入力ファイルの読み込み
data <- read.table(in_f1, header=TRUE, sep="\t", quote="") #in_f1で指定したファイルの読み込み
keywords <- readLines(in_f2) #in_f2で指定したファイルの読み込み
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) #条件を満たすかどうかを判定した結果をobjに格納
out <- data[obj,] #objがTRUEとなる行のみ抽出した結果をoutに格納
dim(out) #オブジェクトoutの行数と列数を表示
```

1と12は手順が異なるだけで実質的に同じです

12. 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、param2で指定した文字列が含まれる行全体を出力したい場合:

```
#ファイルに保存
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F) #outの中身を指定したファイル名で保存
```

```
in_f <- "annotation.txt" #入力ファイル名を指定してin_fに格納(アノテーションファイル)
out_f <- "hoge12.txt" #出力ファイル名を指定してout_fに格納
param1 <- 1 #アノテーションファイル中の検索したい列番号を指定
param2 <- c("gene1", "gene7", "gene9") #検索したい文字列を指定
```

```
#入力ファイルの読み込み
data <- read.table(in_f, header=TRUE, sep="\t", quote="") #in_f1で指定したファイルの読み込み
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[,param1]), param2) #条件を満たすかどうかを判定した結果をobjに格納
out <- data[obj,] #objがTRUEとなる行のみ抽出した結果をoutに格納
dim(out) #オブジェクトoutの行数と列数を表示
```

```
#ファイルに保存
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F) #outの中身を指定したファイル名で保存
```

12. 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、param2で指定した文字列が含まれる行全体を出力したい場合:  
・ イントロ | 一般 | [任意のキーワードを含む行を抽出\(基礎\)](#)

```
in_f <- "annotation.txt" #入力ファイル名を指定してin_fに格納(アンテーションファイル)  
out_f <- "hoge12.txt" #出力ファイル名を指定してout_fに格納  
param1 <- 1 #アンテーションファイル中の検索したい列番号を指定  
param2 <- c("gene1", "gene7", "gene9") #検索したい文字列を指定  
  
#入力ファイルの読み込み  
data <- read.table(in_f, header=TRUE, sep="\t", quote="")#in_fで指定したファイルの読み込み  
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[,param1]), param2)#条件を満たすかどうかを判定した結果をobjに格納  
out <- data[obj,] #objがTRUEとなる行のみ抽出した結果をoutに格納  
dim(out) #オブジェクトoutの行数と列数を表示  
  
#ファイルに保存  
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F)#outの中身を指定したファイル名で保存
```

### 入力: 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

このコードはヘッダー行  
がある場合のものです

### 出力: hoge12.txt

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

### 13. 目的のタブ区切りテキストファイル(annotation2.txt)中の第1列目をキーとして、param2で指定した文字列が含まれる行全体を出力したい場合:

• イントロ | 一般 | [任意のキーワードを含む行を抽出\(基礎\)](#)

入力ファイル中にヘッダー行がない場合の読み込み例です。

```
in_f <- "annotation2.txt" #入力ファイル名を指定してin_fに格納(アノテーションファイル)
out_f <- "hoge13.txt" #出力ファイル名を指定してout_fに格納
param1 <- 1 #アノテーションファイル中の検索したい列番号を指定
param2 <- c("gene1", "gene7", "gene9") #検索したい文字列を指定

#入力ファイルの読み込み
data <- read.table(in_f, header=F, sep="\t", quote="") #in_fで指定したファイルの読み込み
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[,param1]), param2) #条件を満たすかどうかを判定した結果をobjに格納
out <- data[obj,] #objがTRUEとなる行のみ抽出した結果をoutに格納
dim(out) #オブジェクトoutの行数と列数を表示

#ファイルに保存
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F, col.names=F) #outの中身を指定したファイル名で保存
```

入力: annotation2.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	iyaaaa	endoplasmic



出力: hoge13.txt

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

このコードはヘッダー行がない場合のものです

## 12. 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、param2で指定した文字列が含まれる行全体を出力したい場合:

• イントロ | 一般 | [任意のキーワードを含む行を抽出\(基礎\)](#)

```
in_f <- "annotation.txt" #入力ファイル名を指定してin_fに格納(アノテーションファイル)
out_f <- "hoge12.txt" #出力ファイル名を指定してout_fに格納
param1 <- 1 #アノテーションファイル中の検索したい列番号を指定
param2 <- c("gene1", "gene7", "gene9") #検索したい文字列を指定

#入力ファイルの読み込み
data <- read.table(in_f, header=TRUE, sep="\t", quote="") #in_fで指定したファイルの読み込み
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[,param1]), param2) #条件を満たすかどうかを判定した結果をobjに格納
out <- data[obj,] #objがTRUEとなる行のみ抽出した結果をoutに格納
dim(out) #オブジェクトoutの行数と列数を表示

#ファイルに保存
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F) #outの中身を指定したファイル名で保存
```

ヘッダー行がある場合

## 13. 目的のタブ区切りテキストファイル(annotation2.txt)中の第1列目をキーとして、param2で指定した文字列が含まれる行全体を出力したい場合:

入力ファイル中にヘッダー行がない場合の読み込み例です。

```
in_f <- "annotation2.txt" #入力ファイル名を指定してin_fに格納(アノテーションファイル)
out_f <- "hoge13.txt" #出力ファイル名を指定してout_fに格納
param1 <- 1 #アノテーションファイル中の検索したい列番号を指定
param2 <- c("gene1", "gene7", "gene9") #検索したい文字列を指定

#入力ファイルの読み込み
data <- read.table(in_f, header=F, sep="\t", quote="") #in_fで指定したファイルの読み込み
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[,param1]), param2) #条件を満たすかどうかを判定した結果をobjに格納
out <- data[obj,] #objがTRUEとなる行のみ抽出した結果をoutに格納
dim(out) #オブジェクトoutの行数と列数を表示

#ファイルに保存
write.table(out, out_f, sep="\t", append=F, quote=F, row.names=F, col.names=F) #outの中身を指定したファイル名で保存
```

ヘッダー行がない場合

# 課題1

- シロイヌナズナのGFF3形式ファイル(TAIR10\_GFF3\_genes.gff)を入力として下記の解析結果を示せ。また、3と4の結果の違いについて論ぜよ
  1. 染色体数(3列目が"chromosome"となっている行数)
  2. mRNAの個数(3列目が"mRNA"となっている行数)
  3. 1列目が"Chr1"および"ChrM"となっているトータルの行数
  4. 1列目が"chromosome 1"および"ChrM"となっているトータルの行数

	A	B	C	D	E	F	G	H	I
1	Chr1	TAIR10	chromosome	1	30427671	.	.	.	ID=Chr1;Name=Chr1
2	Chr1	TAIR10	gene	3631	5899	.	+	.	ID=AT1G01010;Note=protein_coding_gene;Name=AT1G01010
3	Chr1	TAIR10	mRNA	3631	5899	.	+	.	ID=AT1G01010.1;Parent=AT1G01010;Name=AT1G01010.1;Index=1
4	Chr1	TAIR10	protein	3760	5630	.	+	.	ID=AT1G01010.1-Protein;Name=AT1G01010.1;Derives_from=AT1G01010.1
5	Chr1	TAIR10	exon	3631	3913	.	+	.	Parent=AT1G01010.1
6	Chr1	TAIR10	five_prime_UTR	3631	3759	.	+	.	Parent=AT1G01010.1
7	Chr1	TAIR10	CDS	3760	3913	.	+	0	Parent=AT1G01010.1,AT1G01010.1-Protein;
8	Chr1	TAIR10	exon	3996	4276	.	+	.	Parent=AT1G01010.1
9	Chr1	TAIR10	CDS	3996	4276	.	+	2	Parent=AT1G01010.1,AT1G01010.1-Protein;
10	Chr1	TAIR10	exon	4486	4605	.	+	.	Parent=AT1G01010.1
11	Chr1	TAIR10	CDS	4486	4605	.	+	0	Parent=AT1G01010.1,AT1G01010.1-Protein;
12	Chr1	TAIR10	exon	4706	5095	.	+	.	Parent=AT1G01010.1

入力ファイルのヘッダー行の有無に気をつけよう

# 課題1の3.と4.のヒント

```
R Console
> in_f <- "TAIR10_GFF3_genes.gff"           #入力ファイル名を指定してin_fに格納(アノテーションファイル)
> data <- read.table(in_f, header=F, sep="\t", quote="") #in_fで指定したファイルの読み込み
> dim(data)                                 #オブジェクトdataの行数と列数を表示
[1] 590264      9
> head(data)                                #最初の6行分を表示
  V1      V2      V3      V4      V5 V6 V7 V8      V9
1 Chr1 TAIR10 chromosome 1 30427671 . . . ID=Chr1;Name=Chr1
2 Chr1 TAIR10 gene 3631 5899 . + . ID=AT1G01010;Note=protein_coding_gene;Name=AT1G01010
3 Chr1 TAIR10 mRNA 3631 5899 . + . ID=AT1G01010.1;Parent=AT1G01010;Name=AT1G01010.1;Index=1
4 Chr1 TAIR10 protein 3760 5630 . + . ID=AT1G01010.1-Protein;Name=AT1G01010.1;Derives_from=AT1G01010.1
5 Chr1 TAIR10 exon 3631 3913 . + . Parent=AT1G01010.1
6 Chr1 TAIR10 five_prime_UTR 3631 3759 . + . Parent=AT1G01010.1
> head(data[,1])
[1] Chr1 Chr1 Chr1 Chr1 Chr1 Chr1
Levels: Chr1 Chr2 Chr3 Chr4 Chr5 ChrC ChrM
> table(data[,1])

  Chr1  Chr2  Chr3  Chr4  Chr5  ChrC  ChrM
157712 91857 113968 90371 135017 616 723
> 157712+91857+113968+90371+135017+616+723
[1] 590264
> |
```

Levelsのところはこのベクトル中に存在する文字列が表示されている。この場合全部で7種類。

table関数実行結果は文字列ごとの出現回数



# multi-FASTAファイルからの各種情報抽出

## FASTAフォーマット [\[編集\]](#)

FASTAでは、シーケンスデータの記述形式としてFASTAフォーマットという形式を使う。FASTAフォーマットはブレンテキストである。1つのシーケンスのデータは、">"で始まる1行のヘッダ行と、2行目以降の実際のシーケンス文字列で構成される。ヘッダ行では、">"の次にシーケンスデータを識別するための文字列を記述し、続けてそのシーケンスデータを説明する文字列を記述する(両方とも省略してよい)。ヘッダ行の">"と識別文字列の間にスペースを入れてはいけない。FASTAフォーマットの全ての行は、80文字未満とすることが推奨される。">"で始まる別の行が出現すると、そこでシーケンスデータが区切られ、別のシーケンスデータが始まる。

FASTA ファイルフォーマットの例を示す。

```
>gi|5524211|gb|AAD44166.1| cytochrom
LCLYTHIGRNIYYGSYLYSETWNTGIMLLLITMATA
EWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHL
LLILILLLLLLLALLSPDMLGDPDNHMPADPLNTPH
GLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTW
IENY
```

Rでmulti-FASTAファイルを読み込んで自在に解析できます

# コピー (CTRL+ALT+左クリック) & ペースト

イントロ | NGS | 読み込み | FASTA形式 | 基本情報を取得 **NEW**

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。  
「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. イントロ | 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)の場合:

```

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

```

```

f <- "hoge4.fa"
out_f <- "hoge1.txt"

#必要なパッケージをロー
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNASTrin

#本番(基本情報取得)
Total_len <- sum(widt
Number_of_contigs <-
Average_len <- mean(w
Median_len <- median(
Max_len <- max(width(
Min_len <- min(width(

#本番(N50情報取得)
sorted <- rev(sort(wi
N50 <- sorted[cumsum(sorted) >= Total_len/2][1]#「

#本番(GC含量情報取得)

```

1

```

RGui (64-bit)
ファイル 編集 閲覧 その他
R Console
> #GC含量(GC content)計算のと
> count <- alphabetFrequency
> 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
> 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()$
>
> |

```

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

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

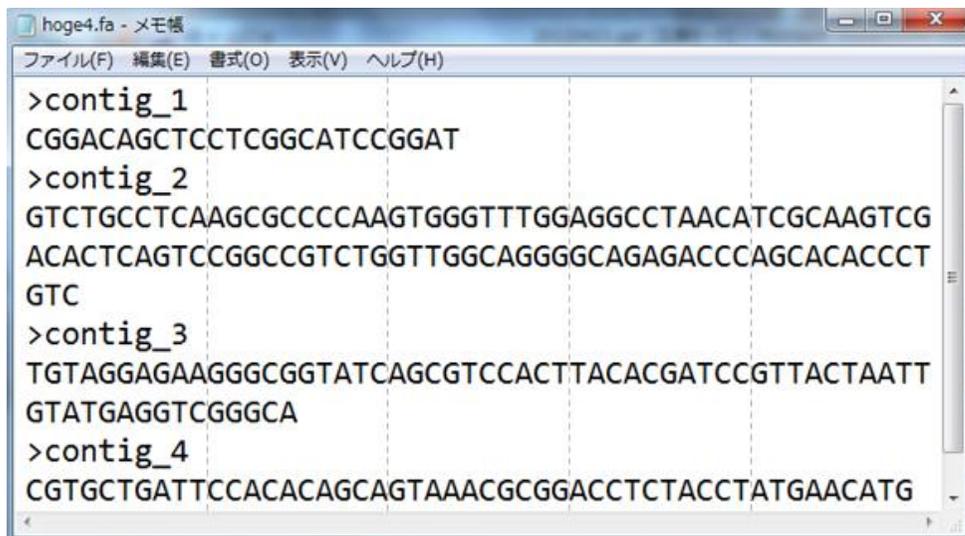
イントロ | NGS | 読み込み | FASTA形式 | 基本情報を取得 **NEW**

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。「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. [イントロ](#) | [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](#))の場合:

```
in_f <- "hoge4.fa" #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.txt" #出力ファイル名を指定してout_fに格納
```

入力: hoge4.fa



出力: hoge1.txt

	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

## ■ アセンブル結果の評価基準の一つ

- 長いコンティグから足していってTotal\_lengthの50%に達したときのコンティグの長さ
- 一般に数値が大きいほどよい

	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

ID	Length
contig_1	24
contig_2	103
contig_3	65
contig_4	49

contig\_2 (103 bp)

contig\_3 (65 bp)

contig\_4 (49 bp)

contig\_1 (24 bp)

Total\_length / 2 (120.5 bp)

Total\_length (241 bp)

averageだと外れ値の影響を受けやすく、medianだと短いコンティグが多くを占める場合に不都合らしい。

# 情報抽出手順の一部

[イントロ](#) | [NGS](#) | [読み込み](#) | [FASTA形式](#) | [基本情報を取得](#) **NEW**

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。  
「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. [イントロ](#) | [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](#))の場合:

```
in_f <- "hoge4.fa"           #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.txt"        #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
fasta <- readDNAStringSet(in_f, format="fasta")#in_fで指定したファイルの読み込み

#本番(基本情報取得)
Total_len <- sum(width(fasta)) #コンティグの「トータル長さ」を取得
Number_of_contigs <- length(fasta) #「コンティグ数」を取得
Average_len <- mean(width(fasta)) #コンティグの「平均長さ」を取得
Median_len <- median(width(fasta)) #コンティグの「中央値」を取得
Max_len <- max(width(fasta)) #コンティグの長さの最大値を取得
Min_len <- min(width(fasta)) #コンティグの長さの最小値を取得
```

```
R Console
> fasta
A DNAStringSet instance of length 4
  width seq
[1] 24 CGGACAGCTCCTCGGCATCCGGAT contig_1
[2] 103 GTCTGCCTCAAGCGCC...CCCAGCACACCCTGTC contig_2
[3] 65 TGTAGGAGAAGGGCGG...TGTATGAGGTCGGGCA contig_3
[4] 49 CGTGCTGATTCCACAC...TCTACCTATGAACATG contig_4
> width(fasta)
[1] 24 103 65 49
> sum(width(fasta))
[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

ID	Length
contig_1	24
contig_2	103
contig_3	65
contig_4	49

R Console

```
> fasta[1]
  A DNASTringSet instance of length 1
  width seq                               names
[1]    24 CGGACAGCTCCTCGGCATCCGGAT        contig_1
> fasta[2]
  A DNASTringSet instance of length 1
  width seq                               names
[1]   103 GTCTGCCTCAAGCGCC...CCCAGCACACCCTGTC contig_2
> width(fasta) >= 50
[1] FALSE TRUE TRUE FALSE
> fasta[width(fasta) >= 50]
  A DNASTringSet instance of length 2
  width seq                               names
[1]   103 GTCTGCCTCAAGCGCC...CCCAGCACACCCTGTC contig_2
[2]    65 TGTAGGAGAAGGGCGG...TGTATGAGGTCGGGCA contig_3
> |
```

上矢印キーを押すと直前に打ったコマンドが出る。有効に利用し最小限の労力で打つ。

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!



# コードの中身が分かると応用範囲が拡大

- 前処理 | クオリティチェック | [配列長分布を調べる](#) (last modified 2013/06/18)
- 前処理 | フィルタリング | [PHREDスコアが低い塩基をNに置換](#) (last modified 2013/06/15)
- 前処理 | フィルタリング | [PHREDスコアが低い配列\(リード\)を除去](#) (last modified 2013/06/15)
- 前処理 | フィルタリング | [ACGTのみからなる配列を抽出](#) (last modified 2013/06/18)
- 前処理 | フィルタリング | [ACGT以外のcharacter "-"をNに変換](#) (last modified 2013/06/18)
- 前処理 | フィルタリング | [ACGT以外の文字数が閾値以下の配列を抽出](#) (last modified 2013/09/27)
- 前処理 | フィルタリング | [重複のない配列セットを作成](#) (last modified 2013/06/18)
- 前処理 | フィルタリング | [指定した長さ以上の配列を抽出](#) (last modified 2013/06/18)
- 前処理 | フィルタリング | [指定した長さの範囲の配列を抽出](#) (last modified 2013/06/18)

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

## 前処理 | フィルタリング | 指定した長さ以上の配列を抽出 NEW

FASTA形式やFASTQ形式ファイルを入力として、指定した配列長以上の配列を抽出するやり方を示します。  
「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

### 1. FASTA形式ファイル(hoge4.fa)の場合:

```

in_f <- "hoge4.fa" #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.txt" #出力ファイル名を指定してout_fに格納
param <- 50 #配列長の閾値を指定

#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
fasta <- readDNASTringSet(in_f, format="fasta") #in_fで指定したファイルの読み込み
fasta #確認してるだけです

#本番
obj <- as.logical(width(fasta) >= param) #条件を満たすかどうかを判定した結果をobjに格納
fasta <- fasta[obj] #objがTRUEとなる要素のみ抽出した結果をfastaに格納
fasta #確認してるだけです

#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta", width=50) #fastaの中身を指定したファ
    
```

- [はじめに](#) (last modified 2014/01/30) **NEW**
- [Rのインストールと起動](#) (last modified 2013/09/27)
- [サンプルデータ](#) (last modified 2014/02/06) **NEW**
- インポート | 一般 | [ランダムに行を抽出](#) (last modified 2013/10/10)
- インポート | 一般 | [任意の文字列を行の最初に挿入](#) (last modified 2013/10/10)
- インポート | 一般 | [任意のキーワードを含む行を抽出\(基礎\)](#) (last modified 2014/02/06) **NEW**
- インポート | 一般 | [ランダムな塩基配列を生成](#) (last modified 2013/09/29)
- インポート | 一般 | [任意の長さの可能な全ての塩基配列を作成](#) (last modified 2013/06/14)
- インポート | 一般 | [任意の位置の塩基を置換](#) (last modified 2013/09/12)
- インポート | 一般 | [指定した範囲の配列を取得](#) (last modified 2014/02/07) **NEW**

## イントロ | 一般 | 指定した範囲の配列を取得 **NEW**

single-FASTA形式やmulti-FASTA形式ファイルから様々な部分配列を取得するやり方を示します。  
 「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

### 1. (single-)FASTA形式ファイル(sample1.fasta)の場合:

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(hoge1.txt)に出力するやり方です。

```

in_f <- "sample1.fasta"      #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.txt"        #出力ファイル名を指定してout_fに格納
param <- c(3, 9)           #抽出したい範囲の始点と終点を指定

#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
fasta <- readDNASTringSet(in_f, format="fasta") #in_fで指定したファイルの読み込み
fasta                                     #確認してるだけです

#本番
fasta <- subseq(fasta, param[1], param[2]) #paramで指定した始点と終点の範囲の配列を抽出し
fasta                                     #確認してるだけです

#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta", width=50) #fastaの中身を指定したファイル名で保存
  
```

入力: sample1.fasta

```
>kadota
AGTGACGGTCTT
```

出力: hoge1.txt

```
>kadota
TGACGGT
```

# 1. (single-)FASTA形式ファイル(sample1.fasta)の場合:

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(hoge1.txt)に保存する

```
in_f <- "sample1.fasta" #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.txt"    #出力ファイル名を指定してout_fに格納
param <- c(3, 9)        #抽出したい範囲の始点と終点を指定

#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
fasta <- readDNAStringSet(in_f, format="fasta") #in_fで指定したファイルの読み込み
fasta                                         #確認してるだけです

#本番
fasta <- subseq(fasta, param[1], param[2]) #paramで指定した始点と終点の範囲の配列を抽出
fasta                                         #確認してるだけです

#ファイルに保存
writeXStringSet(fasta, file=out_f)
```

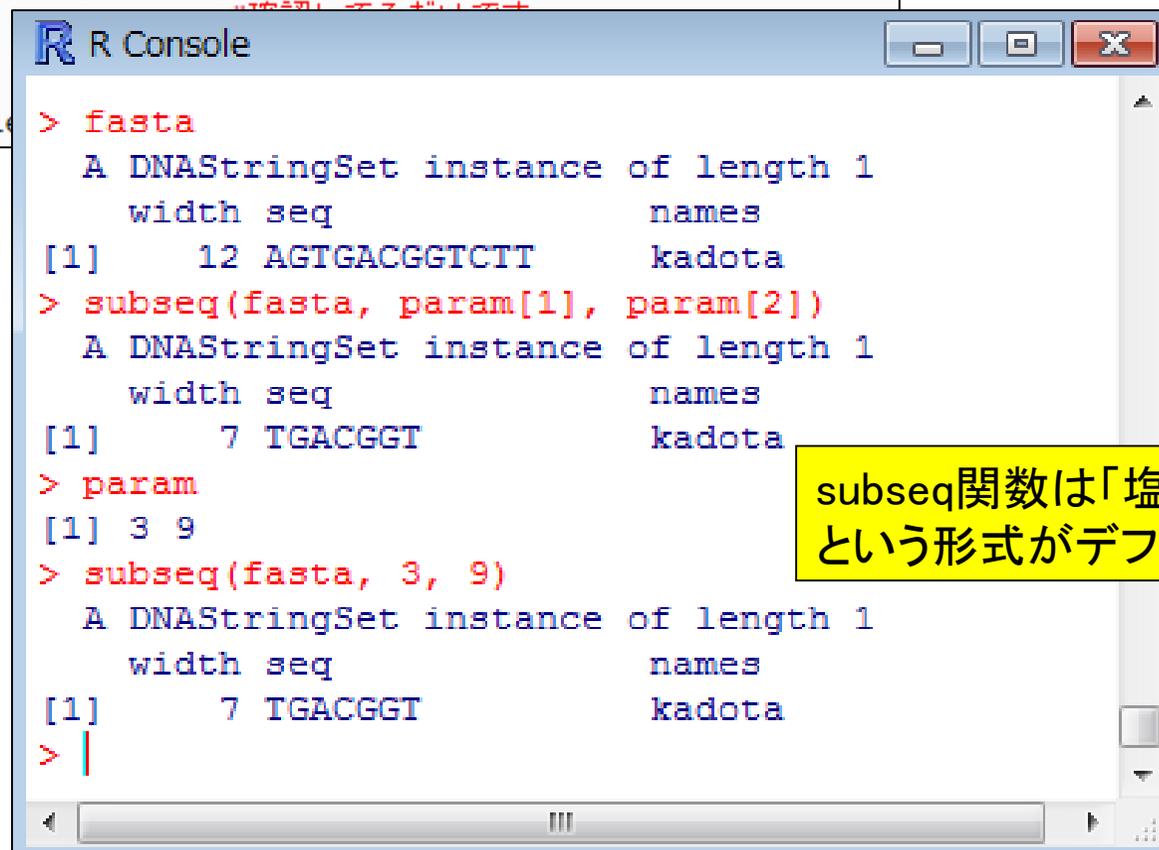
• イントロ | 一般 | [指定した範囲の配列を取得](#)

入力: sample1.fasta

```
>kadota
AGTGACGGTCTT
```

出力: hoge1.txt

```
>kadota
TGACGGT
```



```
R Console
> fasta
A DNAStringSet instance of length 1
  width seq          names
[1]    12 AGTGACGGTCTT kadota
> subseq(fasta, param[1], param[2])
A DNAStringSet instance of length 1
  width seq          names
[1]     7 TGACGGT      kadota
> param
[1] 3 9
> subseq(fasta, 3, 9)
A DNAStringSet instance of length 1
  width seq          names
[1]     7 TGACGGT      kadota
> |
```

subseq関数は「塩基配列, start, end」という形式がデフォルトのようだ

# 関数の使用法について

The screenshot shows the R Console on the left and the R Documentation window on the right. The console shows the command `?subseq` being entered. The documentation window displays the title "XVector objects" and a description of the XVector class, including a list of methods like `subseq(x4, start=10)` and `x3[length(x3):1]`.

R Console

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

R Documentation

## XVector objects

### Description

The XVector virtual class is a general container for storing an "external vector". It inherits from the [Vector](#), which has a very rich interface.

The following classes 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]
```

[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
- ページの下のほうに、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

# 原因既知状態で意図的にエラーを出す

```
R Console
> fasta
A DNASTringSet instance of length 1
  width seq                      names
[1]    12 AGTGACGGTCTT           kadota
> subseq(fasta, start=3, end=9)
A DNASTringSet instance of length 1
  width seq                      names
[1]     7 TGACGGT                kadota
> subseq(fasta, start=3, width=11)
以下にエラー .Call2("solve_user_SEW", refwidths, start, end, width, transla$
solving row 1: 'allow.nonnarrowing' is FALSE and the solved end (13) is > $
> subseq(fasta, start=3, width=10)
A DNASTringSet instance of length 1
  width seq                      names
[1]    10 TGACGGTCTT           kadota
> |
```

入力: `sample1.fasta`

```
>kadota
AGTGACGGTCTT
```

出力: `hoge1.txt`

```
>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. インポート一般 | 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)の場合:

目的のaccession番号が複数ある場合に対応したものです。予め用意しておいた「1列目:accession, 2列目:start位置, 3列目:end位置」からなるリストファイル(list\_sub2.txt)を読み込ませて、目的の配列のmulti-fastaファイルを取得するやり方です。

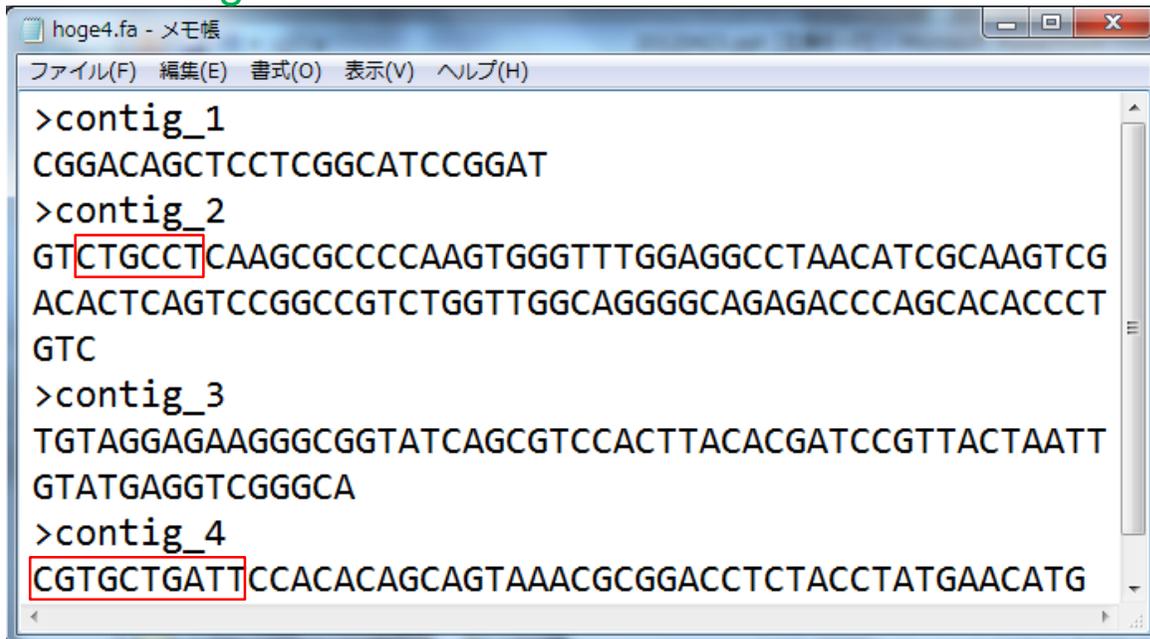
```
in_f1 <- "hoge4.fa" #入力ファイル名を指定してin_f1に格納(multi-fastaファイル)
in_f2 <- "list_sub2.txt" #入力ファイル名を指定してin_f2に格納(リストファイル)
out_f <- "hoge4.txt" #出力ファイル名を指定してout_fに格納

#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
fasta <- readDNASTringSet(in_f1, format="fasta") #in_f1で指定したファイルの読み込み
posi <- read.table(in_f2) #in_f2で指定したファイルの読み込み
fasta #確認してるだけです

#本番
hoge <- NULL #最終的に得る結果を格納するためのプレースホルダhogeを作成して
for(i in 1:nrow(posi)){ #length(posi)回だけループを回す
  obj <- names(fasta) == posi[i,1] #条件を満たすかどうかを判定した結果をobjに格納
  hoge <- append(hoge, subseq(fasta[obj], start=posi[i,2], end=posi[i,3])) #subseq関数を用いて
```

#### 入力1: hoge4.fa



```
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACACCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG
```

#### 入力2: list\_sub2.txt

contig_4	1	10
contig_2	3	8

#### 出力: hoge4.txt

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

- 正規化 | サンプル間 | 3群間 | 複製あり | [IDEGES/edgeR\(Sun 2013\)](#)(last modified 2013/09/15)推奨
- 正規化 | サンプル間 | 3群間 | 複製あり | [TMM\(Robinson 2010\)](#)(last modified 2013/09/16)
- 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 2013/06/19)
- 解析 | 一般 | [GC含量 \(GC contents\)](#)(last modified 2013/06/24)
- 解析 | 一般 | [Sequence logos\(Schneider 1990\)](#)(last modified 2012/06/27)
- 解析 | 一般 | 上流配列解析 | [LDSS\(Yamamoto 2007\)](#)(last modified 2012/07/17)
- 解析 | 一般 | 上流配列解析 | [Relative](#)
- 解析 | 基礎 | [平均-分散プロット\(Techn](#)
- 解析 | 基礎 | [平均-分散プロット\(Biolog](#)
- 解析 | [クラスターリング | について](#)(last modified 2013/06/19)
- 解析 | [クラスターリング | サンプル間 | he](#)
- 解析 | [クラスターリング | 遺伝子間 | MB](#)
- 解析 | 発現変動 | [ポアソン分布 | シミュ](#)
- 解析 | 発現変動 | [負の二項分布 | シミュ](#)
- 解析 | [発現変動 | について](#)(last modified 2013/06/19)
- 解析 | [発現変動 | 2群間 | 対応なし | 複](#)
- 解析 | [発現変動 | 2群間 | 対応なし | 複](#)

配列ごとのGC含量を計算したいとき

## 解析 | 一般 | GC含量 (GC contents)

multi-fasta形式ファイルを読み込んで配列ごとのGC含量 (GC contents)を出力するやり方を示します。出力ファイルは、「description」「CGの総数」「ACGTの総数」「配列長」「%GC含量」としています。尚、%GC含量は「CGの総数/ACGTの総数」で計算しています。

「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. [イントロ](#) | 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](#))の場合:

```

in_f <- "hoge4.fa"           #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.txt"        #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
fasta <- readDNAStringSet(in_f, format="fasta")#in_fで指定したファイルの読み込み

#本番
hoge <- alphabetFrequency(fasta) #A,C,G,T,..の数を各配列ごとにカウントした結果を
CG <- rowSums(hoge[,2:3])       #C,Gの総数を計算してCGIに格納
ACGT <- rowSums(hoge[,1:4])     #A,C,G,Tの総数を計算してACGTIに格納
GC_content <- CG/ACGT*100      #%GC含量を計算してGC_contentIに格納

#ファイルに保存
tmp <- cbind(names(fasta), CG, ACGT, width(fasta), GC_content)#ファイルに保存したい列
colnames(tmp) <- c("description", "CG", "ACGT", "Length", "%GC_contents")#列名情報
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F, col.names=T)#tmp

```

multi-fasta形式ファイルを読み込んで配列ごとのGC含量 (GC contents)を出力するやり方を示します。出力ファイルは、「description」「CGの総数」「ACGTの総数」「配列長」「%GC含量」としています。尚、%GC含量は「CGの総数/ACGTの総数」で計算しています。

「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ペ。

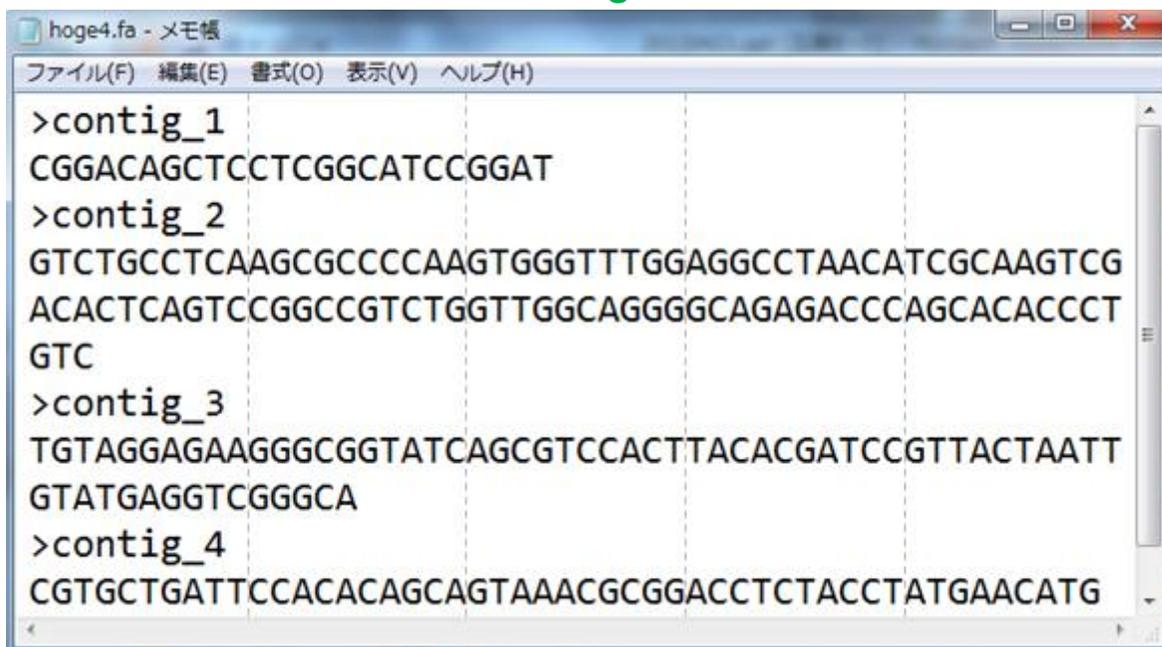
1. [イントロ | 一般 | ランダムな塩基配列を作成](#)の4.を実行して得られたmulti-fastaファイル([hoge4.fa](#))の場合:

```
in_f <- "hoge4.fa"           #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.txt"        #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
fasta <- readDNAStringSet(in_f, format="fasta") #in_fで指定したファイルの読み込み
```

入力: hoge4.fa



出力: hoge1.txt

	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

このコードをテンプレートとすればゲノム配列を読み込んで、染色体ごとのGC含量を算出可能

# 課題2

1. シロイヌナズナのゲノム配列ファイル([TAIR10\\_chr\\_all.fas](#))を入力として染色体ごとの%GC含量計算結果を示せ(小数点2ケタ程度でよい)。
2. 右表はシロイヌナズナゲノム配列決定に関する原著論文中のGC含量の数値である。発表当時(2000年ごろ)と2010年ごろのバージョン(TAIR10)の違いを含めて簡単に考察せよ。

原著論文の%GC含量

	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%

入力ファイルはhoge  
フォルダ内にあります

# GC含量

```

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
fasta <- readDNASTringSet(in_f, format="fasta")#in_fで指定し

#本番
hoge <- alphabetFrequency(fasta)
CG <- rowSums(hoge[,2:3])
ACGT <- rowSums(hoge[,1:4])
GC_content <- CG/ACGT*100
    
```

#入力ファイル名を指定し  
#出力ファイル名を指定し

詳細を解説します

#パッケージの読み込み

#A,C,G,T...の数を各配  
#C,Gの総数  
#A,C,G,T  
#%GC含量

```

#ファイルに保存
tmp <- cbind(names(fasta), CG, ACGT, width(fasta))
colnames(tmp) <- c("description", "CG", "ACGT", "width")
write.table(tmp, out_f, sep="\t", append=F, quote=F)
    
```

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

> fasta
A DNASTringSet instance of length 4
width seq names
[1] 24 CGGACAGCTCCTCGGCATCCGGAT contig_1
[2] 103 GTCTGCCTCAAGC...GCACACCCTGTC contig_2
[3] 65 TGTAGGAGAAGGG...TGAGGTCGGGCA contig_3
[4] 49 CGTGCTGATTCCA...CCTATGAACATG contig_4

> hoge
      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
starting httpd help server ... done

> |
    
```

# 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  
 どのような出力結果を返すの  
 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)という事実は常識だから書かれていないのか?!

```
R Console
> alphabetFrequency(fasta)
      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(fasta, baseOnly=T)
      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(引数)とは、入力や出力の形式などを指定するオプションであり、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](#), [XStringViews](#) or [MaskedXString](#) object for `alphabetFrequency`, `letterFrequency`, or `uniqueLetters`.

DNA or RNA input for `hasOnlyBaseLetters`.

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

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

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

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

```
R Console
> alphabetFrequency(fasta, baseOnly=T)
      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(fasta, baseOnly=T, as.prob=T)
      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

```

in_f <- "hoge4.fa"           #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.txt"        #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
fasta <- readDNASTringSet(in_f, format="fasta")#in_fに格納

#本番
hoge <- alphabetFrequency(fasta) #A,C,G,T,...の出現頻度を算出
CG <- rowSums(hoge[,2:3])        #C,Gの総数を算出
ACGT <- rowSums(hoge[,1:4])      #A,C,G,Tの総数を算出
GC_content <- CG/ACGT*100        #%GC含量を計算

#ファイルに保存
tmp <- cbind(names(fasta), CG, ACGT, width(fasta),
colnames(tmp) <- c("description", "CG", "ACGT", "Length")
write.table(tmp, out_f, sep="\t", append=F, quote=F)

```

「コンティグごとのCとGの出現頻度の和」は、「行の和」として得ることができる

```

R Console
> hoge
      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
[2,] 20 34 31 18  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
[4,] 14 15 10 10  0  0  0  0  0  0  0  0  0  0  0  0
> dim(hoge)
[1]  4 17
> hoge[,2:3]
      C  G
[1,]  9  7
[2,] 34 31
[3,] 13 20
[4,] 15 10
> rowSums(hoge[,2:3])
[1] 16 65 33 25
> |

```

# GC含量

```

in_f <- "hoge4.fa"           #入力ファイル
out_f <- "hoge1.txt"        #出力ファイル

#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
fasta <- readDNAStringSet(in_f, format="fasta")#in_f

#本番
hoge <- alphabetFrequency(fasta) #A,C,G,T,..の総数を計算
CG <- rowSums(hoge[,2:3])        #C,Gの総数を計算
ACGT <- rowSums(hoge[,1:4])      #A,C,G,Tの総数を計算
GC_content <- CG/ACGT*100        #%GC含量を計算

#ファイルに保存
tmp <- cbind(names(fasta), CG, ACGT, width(fasta), GC_content)#ファイルに保存したい
colnames(tmp) <- c("description", "CG", "ACGT", "Length", "%GC_contents")#列名情報
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F, col.names=T)#出力
    
```

```

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ー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

```

in_f <- "hoge4.fa"           #入力ファイル
out_f <- "hoge1.txt"        #出力ファイル

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

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

#本番
hoge <- alphabetFrequency(fasta) #A,C,G,T,..の頻度
CG <- rowSums(hoge[,2:3])        #C,Gの総数を計算
ACGT <- rowSums(hoge[,1:4])     #A,C,G,Tの総数を計算
GC_content <- CG/ACGT*100      #%GC含量を計算

#ファイルに保存
tmp <- cbind(names(fasta), CG, ACGT, width(fasta),
             colnames(tmp) <- c("description", "CG", "ACGT", "Length", "%GC_contents")
write.table(tmp, out_f, sep="\t", append=F, quote=F)
    
```

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

```

R Console
> fasta
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(fasta)
[1] "contig_1" "contig_2" "contig_3" "contig_4"
> CG
[1] 16 65 33 25
> cbind(names(fasta), 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

# GC含量

```

in_f <- "hoge4.fa"           #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.txt"        #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
fasta <- readDNAStringSet(in_f, format="fasta")#in_fで指定したファイルの読み込み

#本番
hoge <- alphabetFrequency(fasta) #A,C,G,T,..の数を各配列ごとにカウントした結果をhogeに格納
CG <- rowSums(hoge[,2:3])        #C,Gの総数を計算してCGに格納
ACGT <- rowSums(hoge[,1:4])     #A,C,G,Tの総数を計算してACGTに格納
GC_content <- CG/ACGT*100      #%GC含量を計算してGC_contentに格納

#ファイルに保存
tmp <- cbind(names(fasta), CG, ACGT, width(fasta), GC_content)#保存したい情報をtmpに格納
colnames(tmp) <- c("description", "CG", "ACGT", "Length", "%GC_contents")#列名を付与
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F, col.names=T)#tmpの中身を指定したファイル名で保存
    
```

cbind関数適用時に、出力したい順番に並べているだけです

	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

#ファイルに保存

```
tmp <- cbind(names(fasta), CG, ACGT, width(fasta), GC_content)#保存したい情報をtmpに格納
colnames(tmp) <- c("description", "CG", "ACGT", "Length", "%GC_contents")#列名を付与
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F, col.names=T)#tmpの中身を指定したファイル名で保存
```

```
R Console
> tmp
      CG  ACGT      GC_content
[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"
> |
```

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

	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

# ヒトゲノム中のCpG出現確率は低い

- 全部で16通りの2連続塩基の出現頻度分布を調べると、CGとなる確率の実測値(0.986%)は期待値(4.2%)よりもかなり低い
- 期待値
  - ゲノム中のGC含量を考慮した場合: 約41%(A:0.295, C:0.205, G: 0.205, T:0.295)なので、 $0.205 \times 0.205 = 4.2\%$
  - ゲノム中のGC含量を考慮しない場合: 50%(A:0.25, C:0.25, G: 0.25, T:0.25)なので、 $0.25 \times 0.25 = 6.25\%$

•	イントロ		一般		<a href="#">任意の位置の塩基を置換</a> (last modified 2013/09/12)
•	イントロ		一般		<a href="#">指定した範囲の配列を取得</a> (last modified 2014/02/07) <b>NEW</b>
•	イントロ		一般		<a href="#">翻訳配列(translate)を取得</a> (last modified 2013/06/14)
•	イントロ		一般		<a href="#">相補鎖(complement)を取得</a> (last modified 2013/06/14)
•	イントロ		一般		<a href="#">逆相補鎖(reverse complement)を取得</a> (last modified 2013/06/14)
•	イントロ		一般		<a href="#">逆鎖(reverse)を取得</a> (last modified 2013/06/14)
•	イントロ		一般		<a href="#">2連続塩基の出現頻度情報を取得</a> (last modified 2014/02/05) <b>NEW</b>
•	イントロ		一般		<a href="#">3連続塩基の出現頻度情報を取得</a> (last modified 2013/06/14)
•	イントロ		一般		<a href="#">任意の長さの連続塩基の出現頻度情報を取得</a> (last modified 2013/06/14)
•	イントロ		一般		Tips   <a href="#">任意の拡張子でファイルを保存</a> (last modified 2013/09/26)
•	イントロ		一般		Tips   <a href="#">拡張子は同じで任意の文字を追加して保存</a> (last modified 2013/09/26)
•	イントロ		一般		配列取得   ゲノム配列   <a href="#">公共DBから</a> (last modified 2013/08/15)
•	イントロ		一般		配列取得   ゲノム配列   <a href="#">BSgenome</a> (last modified 2012/02/05)
•	イントロ		一般		配列取得   プロモーター配列   <a href="#">BSgenome</a> (last modified 2013/10/10)
•	イントロ		一般		配列取得   プロモーター配列   <a href="#">GenomicFeatures(Lawrence 2013)</a> (last modified 2013/10/10)
•	イントロ		一般		配列取得   トランスクリプトーム配列   <a href="#">biomaRt(Durinck 2009)</a> (last modified 2013/09/25)

Rで調べることができます



# 2連続塩基の出現頻度：基本形

イントロ | 一般 | [2連続塩基の出現頻度情報を取得](#) **NEW**

multi-fasta形式ファイルを読み込んで、"AA", "AC", "AG", "AT", "CA", "CC", "CG", "CT", "GA", "GC", "GG", "GT", "TA", "TC", "TG", "TT"の計 $4^2 = 16$ 通りの2連続塩基の出現頻度を調べるやり方を示します。  
ヒトゲノムで"CG"の割合が期待値よりも低い(Lander et al., 2001; Saxonov et al., 2006)ですが、それを簡単に検証できます。  
「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピペ。

## 1. イントロ | 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)の場合:

タイトル通りの出現頻度です。

```
in_f <- "hoge4.fa" #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.txt" #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
fasta <- readDNAStringSet(in_f, format="fasta")#in_fで指定したファイルの読み込み

#本番
out <- dinucleotideFrequency(fasta) #2連続塩基の出現頻度情報をoutに格納

#ファイルに保存
tmp <- cbind(names(fasta), out) #最初の列にID情報、そのあとに出現頻度情報
write.table(tmp, out_f, sep=" ", as.is=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
```

出力:hoge1.txt

	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
contig_1	0	1	1	2	2	2	3	2	2	2	3	0	0	3	0	0
contig_2	4	6	9	1	11	11	5	6	4	9	10	8	1	8	6	3
contig_3	2	4	5	4	4	2	5	2	4	3	7	6	6	4	3	3
contig_4	3	6	2	3	5	3	3	4	3	3	1	2	3	2	4	1

# 2連続塩基の出現確率：基本形

2. [イントロ](#) | [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](#))の場合:

出現頻度ではなく、出現確率を得るやり方です。

```

in_f <- "hoge4.fa"           #入力ファイル名を指定してin_fに格納
out_f <- "hoge2.txt"        #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
fasta <- readDNASTringSet(in_f, format="fasta")#in_fで指定したファイル

#本番
out <- dinucleotideFrequency(fasta, as.prob=T)#2連続塩基の出現確率情報を取得

#ファイルに保存
tmp <- cbind(names(fasta), out) #最初の列にID情報、そのあとに出現確率情報
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F)#tmp
    
```

```

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:hoge2.txt

	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
contig_1	0.0%	4.3%	4.3%	8.7%	8.7%	8.7%	13.0%	8.7%	8.7%	8.7%	13.0%	0.0%	0.0%	13.0%	0.0%	0.0%
contig_2	3.9%	5.9%	8.8%	1.0%	10.8%	10.8%	4.9%	5.9%	3.9%	8.8%	9.8%	7.8%	1.0%	7.8%	5.9%	2.9%
contig_3	3.1%	6.3%	7.8%	6.3%	6.3%	3.1%	7.8%	3.1%	6.3%	4.7%	10.9%	9.4%	9.4%	6.3%	4.7%	4.7%
contig_4	6.3%	12.5%	4.2%	6.3%	10.4%	6.3%	6.3%	8.3%	6.3%	6.3%	2.1%	4.2%	6.3%	4.2%	8.3%	2.1%

# 2連続塩基の出現確率: ヒトゲノム

5. [BSgenome](#)パッケージ中のヒトゲノム配列 ("[BSgenome.Hsapiens.UCSC.hg19](#)")の場合:

出現頻度ではなく、出現確率

出力: hoge5.txt

```
out_f <- "hoge5.txt"
param <- "BSgenome.Hsapiens.UCSC.hg19"

#必要なパッケージをロード
library(Biostrings)
library(param, character.only = TRUE)

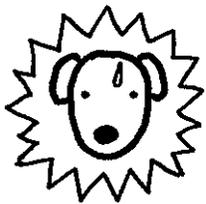
#前処理(paramで指定したファイル名)
tmp <- ls(paste("package", "BSgenome.Hsapiens.UCSC.hg19"))
genome <- eval(parse(text = paste("load(", tmp, ".RData")
fasta <- getSeq(genome, "chr1")
names(fasta) <- seqnames(fasta)

#本番
out <- dinucleotideFrequencies(fasta)

#ファイルに保存
tmp <- cbind(names(fasta), out)
write.table(tmp, out_f, as.is = TRUE)
```

	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
chr1	9.5%	5.0%	7.1%	7.5%	7.3%	5.4%	1.0%	7.1%	6.0%	4.4%	5.4%	5.0%	6.4%	6.0%	7.3%	9.5%
chr2	10.0%	5.0%	7.0%	7.9%	7.2%	5.0%	0.9%	7.0%	5.9%	4.1%	5.0%	5.0%	6.7%	5.9%	7.2%	10.0%
chr3	10.1%	5.0%	6.9%	8.0%	7.2%	4.9%	0.8%	6.9%	5.9%	4.0%	4.9%	5.0%	6.9%	5.9%	7.2%	10.2%
chr4	10.6%	5.0%	6.7%	8.5%	7.1%	4.5%	0.8%	6.7%	5.8%	3.8%	4.5%	5.0%	7.4%	5.8%	7.1%	10.6%
chr5	10.2%	5.0%	6.9%	8.1%	7.2%	4.8%	0.8%	6.9%	5.9%	4.0%	4.8%	5.0%	7.0%	5.9%	7.2%	10.2%
chr6	10.2%	5.0%	6.9%	8.1%	7.2%	4.9%	0.9%	6.9%	5.9%	4.0%	4.9%	5.0%	6.9%	5.9%	7.2%	10.2%
chr7	9.8%	5.0%	7.0%	7.8%	7.2%	5.2%	1.0%	7.0%	5.9%	4.3%	5.2%	5.0%	6.6%	5.9%	7.2%	9.9%
chr8	10.0%	5.1%	6.9%	7.9%	7.2%	5.0%	0.9%	6.9%	5.9%	4.1%	5.0%	5.0%	6.7%	5.9%	7.2%	10.0%
chr9	9.7%	5.0%	7.0%	7.6%	7.3%	5.3%	1.0%	7.0%	5.9%	4.3%	5.3%	5.0%	6.4%	5.9%	7.3%	9.7%
chr10	9.6%	5.0%	7.0%	7.5%	7.3%	5.4%	1.0%	7.1%	6.0%	4.4%	5.4%	5.1%	6.3%	6.0%	7.3%	9.6%
chr11	9.5%	5.0%	7.1%	7.6%	7.3%	5.4%	1.0%	7.1%	6.0%	4.4%	5.4%	5.1%	6.3%	6.0%	7.3%	9.5%
chr12	9.8%	5.0%	7.0%	7.7%	7.2%	5.2%	1.0%	7.0%	5.9%	4.2%	5.2%	5.0%	6.6%	5.9%	7.2%	9.8%
chr13	10.5%	5.0%	6.7%	8.5%	7.1%	4.6%	0.8%	6.7%	5.8%	3.8%	4.6%	5.0%	7.3%	5.8%	7.1%	10.5%
chr14	9.7%	5.0%	7.0%	7.7%	7.2%	5.2%	1.0%	7.0%	5.9%	4.3%	5.2%	5.1%	6.6%	5.9%	7.3%	9.9%
chr15	9.4%	5.1%	7.1%	7.3%	7.3%	5.5%	1.1%	7.2%	6.0%	4.5%	5.5%	5.1%	6.2%	6.0%	7.3%	9.4%
chr16	8.6%	5.1%	7.3%	6.7%	7.5%	6.2%	1.4%	7.3%	6.1%	5.0%	6.2%	5.1%	5.4%	6.1%	7.6%	8.6%
chr17	8.4%	5.0%	7.4%	6.4%	7.4%	6.5%	1.5%	7.4%	6.0%	5.2%	6.5%	5.0%	5.3%	6.1%	7.4%	8.5%
chr18	10.1%	5.0%	6.9%	8.1%	7.2%	4.9%	0.9%	6.9%	5.9%	4.1%	4.9%	5.1%	6.9%	5.9%	7.2%	10.1%
chr19	7.6%	5.1%	7.4%	5.7%	7.6%	7.2%	1.9%	7.4%	6.1%	5.7%	7.3%	5.1%	4.5%	6.1%	7.6%	7.6%
chr20	8.7%	5.0%	7.3%	6.8%	7.5%	6.0%	1.2%	7.3%	6.0%	4.9%	6.1%	5.1%	5.6%	6.1%	7.6%	8.9%
chr21	9.9%	5.1%	6.9%	7.8%	7.3%	5.2%	1.1%	6.8%	5.9%	4.3%	5.2%	5.1%	6.6%	5.9%	7.3%	9.8%
chr22	7.6%	5.1%	7.5%	5.8%	7.7%	7.1%	1.7%	7.5%	6.1%	5.7%	7.1%	5.1%	4.6%	6.1%	7.7%	7.6%
chrX	10.1%	5.0%	6.8%	8.2%	7.2%	4.9%	0.8%	6.8%	5.9%	3.9%	4.9%	5.1%	7.0%	5.9%	7.2%	10.2%
chrY	9.9%	5.1%	6.8%	8.1%	7.3%	4.9%	0.8%	6.8%	6.0%	3.9%	4.9%	5.2%	6.7%	6.0%	7.5%	10.0%

2分強かかります



- インタロ | 一般 | [2連続塩基の出現頻度情報を取得](#) (last modified 2014/02/07)
- インタロ | 一般 | [3連続塩基の出現頻度情報を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [任意の長さの連続塩基の出現頻度情報を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [Tips | 任意の拡張子でファイルを保存](#) (last modified 2013/09/26)
- インタロ | 一般 | [Tips | 拡張子は同じで任意の文字を追加して保存](#) (last modified 2013/09/26)
- インタロ | 一般 | 配列取得 | ゲノム配列 | [公共DBから](#) (last modified 2014/04/10) **NEW**
- インタロ | 一般 | 配列取得 | ゲノム配列 | [BSgenome](#) (last modified 2014/04/01) **NEW**
- インタロ | 一般 | 配列取得 | プロモーター配列 | [公共DBから](#) (last modified 2014/04/02) **NEW**
- インタロ | 一般 | 配列取得 | プロモーター配列 | [BSgenome](#) (last modified 2013/10/10)

**インタロ | 一般 | 配列取得 | ゲノム配列 | BSgenome**

[BSgenome](#)パッケージを用いて様々な生物種のゲノム配列を取得するやり方を示します。ミヤマハタザオ (*A. lyrata*)、セイヨウミツバチ (*A. mellifera*)、シロイヌナズナ (*A. thaliana*)、ウシ (*B. taurus*)、線虫 (*C. elegans*)、犬 (*C. familiaris*)、キイロショウジョウバエ (*D. melanogaster*)、ゼブラフィッシュ (*D. rerio*)、大腸菌 (*E. coli*)、イトヨ (*G. aculeatus*)、セキショクヤケイ (*G. gallus*)、ヒト (*H. sapiens*)、アカゲザル (*M. mulatta*)、マウス (*M. musculus*)、チンパンジー (*P. troglodytes*)、ラット (*R. norvegicus*)、出芽酵母 (*S. cerevisiae*)、トキンプラズマ (*T. gondii*)と実に様々な生物種が利用可能であることがわかります。

`getSeq`関数はBSgenomeオブジェクト中の「single sequences」というあたりにリストアップされているchr...というものを全て抽出しています。したがって、例えばマウスゲノムは「chr1」以外に「chr1\_random」や「chrUn\_random」なども等価に取扱っている点に注意してください。

「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. 利用可能な生物種とRにインストール済みの生物種をリストアップしたい場合:**

```
#必要なパッケージをロード
library(BSgenome) #パッケージの読み込み

#本番 (利用可能なリストアップ; インストール済みとは限らない)
available.genomes() #このパッケージ中で利用可能なゲノムをリストアップ

#本番 (インストール済みの生物種をリストアップ)
installed.genomes() #インストール済みの生物種をリストアップ

#後処理 (パッケージ名でだいたいわかるがproviderやversionを分割して表示したい場合)
installed.genomes(splitNameParts=TRUE) #インストール済みの生物種をリストアップ
```

様々な生物種のゲノム配列がRのパッケージとして提供されています。installed.genomes関数のところで2分程度かかります

BSgenomeパッケージを用いて様々な生物種のゲノム配列を取得するやり方を示します。ミヤマハタザオ (*A. lyrata*)、セイヨウミツバチ (*A. mellifera*)、シロイヌナズナ (*A. thaliana*)、バエ (*D. melanogaster*)、ゼブラフィッシュ (*D. rerio*)、ヒト (*H. sapiens*)、アカゲザル (*M. mulatta*)、マウス (*M. musculus*)、酵母 (*S. cerevisiae*)、トキソプラズマ (*T. gondii*) と実に様々な生物種に対応しています。したがって、例えばマウスゲノムは「chr1」で指定できます。注意してください。

1. 利用可能な生物種とRにインストール済みの生物種

```
#必要なパッケージをロード
library(BSgenome)

#本番 (利用可能なリストアップ; インストール済みの生物種をリストアップ)
available.genomes()

#本番 (インストール済みの生物種をリストアップ)
installed.genomes()

#後処理 (パッケージ名でだいたいわかるがpr
installed.genomes(splitNameParts=TRUE)
```

```
R Console

[24] "BSgenome.Hsapiens.UCSC.hg18"
[25] "BSgenome.Hsapiens.UCSC.hg19"
[26] "BSgenome.Mmulatta.UCSC.rheMac2"
[27] "BSgenome.Mmulatta.UCSC.rheMac3"
[28] "BSgenome.Mmusculus.UCSC.mm10"
[29] "BSgenome.Mmusculus.UCSC.mm8"
[30] "BSgenome.Mmusculus.UCSC.mm9"
[31] "BSgenome.Osativa.MSU.MSU7"
[32] "BSgenome.Ptroglydotes.UCSC.panTro2"
[33] "BSgenome.Ptroglydotes.UCSC.panTro3"
[34] "BSgenome.Rnorvegicus.UCSC.rn4"
[35] "BSgenome.Rnorvegicus.UCSC.rn5"
[36] "BSgenome.Scerevisiae.UCSC.sacCer2"
[37] "BSgenome.Scerevisiae.UCSC.sacCer3"
[38] "BSgenome.Scerevisiae.UCSC.sacCer4"
[39] "BSgenome.Tgondii.UCSC.tg1"
>
> #本番(インストール済みの生物種をリストアップ)
> installed.genomes() #イ$
[1] "BSgenome.Celegans.UCSC.ce2"
[2] "BSgenome.Drerio.UCSC.danRer7"
[3] "BSgenome.Ecoli.NCBI.20080805"
[4] "BSgenome.Hsapiens.UCSC.hg19"
[5] "BSgenome.Hsapiens.UCSC.hg19.Rs20101127"
[6] "BSgenome.Mmusculus.UCSC.mm9"
[7] "BSgenome.Scerevisiae.UCSC.sacCer2"
>
```

ヒトゲノム(BSgenome.Hsapiens.UCSC.hg19)の2連続塩基出現頻度計算ができたのは、このパッケージをインストール済みだからです



BSgenomeパッケージを用いて様々な生物種のゲノム配列を取得するやり方を示します。ミヤマハタザオ (A. lyrata)、セイヨフミツバチ (A. mellifera)、シロイヌナズナ (A. thaliana)、ウシ (B. taurus)、線虫 (C. elegans)、犬 (C. familiaris)、キイロショウジョウバエ (D. melanogaster)、ゼブラフィッシュ (D. rerio)、大腸菌 (E. coli)、イトヨ (G. aculeatus)、セキショクヤケイ (G. gallus)、ヒト (H. sapiens)、アカゲザル (M. mulatta)、マウス (M. musculus)、チンパンジー (P. troglodytes)、ラット (R. norvegicus)、出芽酵母 (S. cerevisiae)、トキソプラズマ (T. gondii) と実に様々な生物種が利用可能であることがわかります。

getSeq関数はBSgenomeオブジェクト中の「single sequence」を取得します。したがって、例えばマウスゲノムは「chr1」に注意してください。

「ファイル」-「ディレクトリの変更」でファイルを保存した

1. 利用可能な生物種とRにインストール済みの生物種

```
#必要なパッケージをロード
library(BSgenome)

#本番 (利用可能なリストアップ; インストール済みの生物種をリストアップ)
available.genomes()

#本番 (インストール済みの生物種をリストアップ)
installed.genomes()

#後処理 (パッケージ名でだいたいわかるがprovided.installed.genomes(splitNameParts=TRUE))
```

```
R Console

[24] "BSgenome.Hsapiens.UCSC.hg18"
[25] "BSgenome.Hsapiens.UCSC.hg19"
[26] "BSgenome.Mmulatta.UCSC.rheMac2"
[27] "BSgenome.Mmulatta.UCSC.rheMac3"
[28] "BSgenome.Mmusculus.UCSC.mm10"
[29] "BSgenome.Mmusculus.UCSC.mm8"
[30] "BSgenome.Mmusculus.UCSC.mm9"
[31] "BSgenome.Osativa.MSU.MSU7"
[32] "BSgenome.Ptroglodytes.UCSC.panTro2"
[33] "BSgenome.Ptroglodytes.UCSC.panTro3"
[34] "BSgenome.Rnorvegicus.UCSC.rn4"
[35] "BSgenome.Rnorvegicus.UCSC.rn5"
[36] "BSgenome.Scerevisiae.UCSC.sacCer2"
[37] "BSgenome.Scerevisiae.UCSC.sacCer3"
[38] "BSgenome.Scerevisiae.UCSC.sacCer4"
[39] "BSgenome.Tgondii.ToxoDB.tgd3"
>
> #本番(インストール済みの生物種をリストアップ)
> installed.genomes() #イ$
[1] "BSgenome.Celegans.UCSC.ce2"
[2] "BSgenome.Drerio.UCSC.danRer7"
[3] "BSgenome.Ecoli.NCBI.20080805"
[4] "BSgenome.Hsapiens.UCSC.hg19"
[5] "BSgenome.Hsapiens.UCSC.hg19.Rbowtie"
[6] "BSgenome.Mmusculus.UCSC.mm9"
[7] "BSgenome.Scerevisiae.UCSC.sacCer2"
>
```

もしゼブラフィッシュゲノム (BSgenome.Drerio.UCSC.danRer7) がインストールされていないならば...

インストール済みの生物種をリストアップした結果、BSgenome.Drerio.UCSC.danRer7 がリストアップされていることが確認できます。

## 2. ゼブラフィッシュ("BSgenome.Drerio.UCSC.danRer7")のゲノム情報をRにインストールしたい場合:

400MB程度あります...

```
param <- "BSgenome.Drerio.UCSC.danRer7"#パッケージ名を指定

#本番
source("http://bioconductor.org/biocLite.R")#おまじない
biocLite(param) #おまじない

#後処理 (インストール済みの生物種をリストアップ)
installed.genomes() #インストール済みの生物種をリストアップ
```

ここを参考にして  
**available.genomes()** でリストアップされて  
いない任意のパッケージ名を指定

## 3. インストール済みのゼブラフィッシュのゲノム配列をmulti-fastaファイルで保存したい場合:

1.4GB程度のファイルが生成されます...

```
out_f <- "hoge3.txt" #出力ファイル名を指定してout_fに格納
param <- "BSgenome.Drerio.UCSC.danRer7"#パッケージ名を指定

#必要なパッケージをロード
library(param, character.only=T) #paramで指定したパッケージの読み込み

#前処理(paramで指定したパッケージ中のオブジェクト名をgenomeに統一)
#tmp <- unlist(strsplit(param, ".", fixed=TRUE))[2]#paramで指定した文字列からオブジェクト名を取得
tmp <- ls(paste("package", param, sep=":"))#paramで指定したパッケージで利用可能なオブジェクト名を1
genome <- eval(parse(text=tmp)) #文字列tmpをRオブジェクトとしてgenomeに格納(パッケージ中に
genome #確認してるだけです

#本番
fasta <- getSeq(genome) #ゲノム塩基配列情報を抽出した結果をfastaに格納
names(fasta) <- seqnames(genome) #description情報を追加している

#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta", width=50)#fastaの中身を指定したファイル名で保
```

**available.genomes()** でリストアップされて  
いるパッケージ名を指定可能です

# multi-FASTAファイルとして保存したい場合

5. インストール済みのヒト("BSgenome.Hsapiens.UCSC hg19")のゲノム配列をmulti-fastaファイルで保存したい場合:  
 3.0GB程度のファイルが生成されます...。ヒトゲノムは、まだ完全に22本の常染色体とX, Y染色体の計24本になっているわけではないことがわかります。

```

out_f <- "hoge5.txt" #出力ファイル名を指定してout_fに格納
param <- "BSgenome.Hsapiens.UCSC.hg19" #パッケージ名を指定

#必要なパッケージ
library(param)

#前処理(paramで指定したパッケージの読み込み)
genome <- BSgenome(param) #ゲノムオブジェクト名をgenomeに統一

#本番
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)

#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta", width=50)
    
```

8分程度かかります。  
 実行はやめましょう。

フリーズするので、得られたhoge5.txtを  
 テキストエディタで開くことはやめましょう

```

R Console
| chrUn_g1000246
| chrUn_g1000247
| chrUn_g1000248
| chrUn_g1000249
|
| multiple sequences (see '?mseqnames') :
|   upstream1000 upstream2000 upstream5000
|
| (use the '$' or '[' operator to access a
| given sequence)
>
> #本番
> fasta <- getSeq(genome) #ゲノム塩基配列情報を抽出し$
> names(fasta) <- seqnames(genome) #description情報を追加して$
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width=50)#fasta$
> |
    
```

内部的には、  
 ①writeXStringSet関数を用いて、  
 ②fastaオブジェクトを、  
 ③out\_fで指定したファイル名で、  
 ④FASTA形式で、  
 ⑤1行あたりの文字数を50文字  
 にして保存しています。