

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

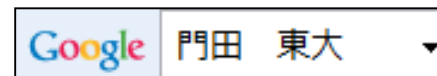
3. データ解析基礎 | 3-3. R各種パッケージ

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)

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

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



Contents

- 3-3. R 各種パッケージ、2014/09/08 15:00-18:15、中級、実習
 - パッケージ
 - 「(Rで)塩基配列解析」のインストール手順おさらい
 - CRANとBioconductor
 - 代表的なパッケージBiostringsの利用法: library, search, objects関数
 - 作業スペース(workspace)の概念
 - Biostringsパッケージで利用可能な関数を概観
 - source関数の利用



パッケージ

(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 2014/08/08, since 2010)

(Rで)塩基配列解析の推奨インストール手順をおさらい。大まかな手順は、①R本体のインストール、および②CRANおよびBioconductorから提供されている各種パッケージのインストール

What's new?

- このウェブページはフリーソフトRと必要な
- 1. [Rのインストールと起動](#)および2. [基本的](#)
- 2014年10月 [HPCIワークショップ「医用フォマティクス」講習会@仙台国際セン](#)
- 門田幸二 著 [シリーズ Useful R 第7巻](#) [トラ](#)

- [日本乳酸菌学会誌のNGS関連連載の第1](#)
- [参考資料\(講義、講習会、本など\)の項目](#)

- [はじめに](#) (last modified 2014/01/30)
- [参考資料\(講義、講習会、本など\)](#) (last mo
- [過去のお知らせ](#) (last modified 2014/08/03)
- [Rのインストールと起動](#) (last modified 2014
- [基本的な利用法](#) (last modified 2014/07/20)
- [サンプルデータ](#) (last modified 2014/07/17)
- [バイオインフォマティクス人材育成カリキュ](#)
- [書籍|トランスクリプトームについて](#) (last m

Rのインストールと起動 NEW

基本的には[こちら](#)または[こちら](#)をご覧ください。

よく分からない人でWindowsユーザーの方は以下を参考にしてください。2014年7月31日にアップデートしたWindows用のインストール手順は[こちら](#)。2014年5月14日にアップデートしたMac版のインストール手順[こちら](#)(by 孫建強氏)もあります。注意点は、「Mac OS Xのバージョンに関わらず R-3.1.0-snowleopard.pkgをインストールしたほうがよい」です。

1. Windows release版のインストールの場合:

1. [Rのインストーラ](#)を「実行」
2. 聞かれるがままに「次へ」などを押してとにかくインストールを完了させる
3. **Windows Vistaの人**は(パッケージのインストール中に書き込み権限に関するエラーが出るのを避けるために)「コントロールパネル」-「ユーザーアカウント」-「ユーザーアカウント 制御の有効化または無効化」で、「ユーザーアカウント 制御(UAC)を使ってコンピュータの保護に役立たせる」の**チェックをあらかじめ外しておく**ことを強くお勧めします。
4. インストールが無事完了したら、デスクトップに出現する「R3.X.Y(32 bitの場合; XやY中の数値はバージョンによって異なります)または「R x64 3.X.Y(64 bitの場合)」アイコンをダブルクリックして起動
5. 以下を、「R コンソール画面」上でコピー&ペーストする。10GB程度のディスク容量を要しますが一番お手軽です。(どこからダウンロードするか?と聞かれるので、その場合は自分のいる場所から近いサイトを指定)

```
install.packages(available.packages()[,1], dependencies=TRUE)#CRAN中にある全てのバック
source("http://www.bioconductor.org/biocLite.R")#おまじない
biocLite(all_group())#Bioconductor中にある全てのパッケージをインストール
biocLite("BSgenome.Athaliana.TAIR.TAIR9", suppressUpdates=TRUE)#Bioconductor中にある
```

6. 「コントロールパネル」-「デスクトップのカスタマイズ」-「フォルダオプション」-「表示(タブ)」-「詳細設定」のところで、「登録されている拡張子は表示しない」のチェックを外してください。

R本体とパッケージの関係

- パソコンを購入しただけの状態では、できることが限られています。
 - 通常は、Officeやウイルス撃退ソフトなどをインストールして利用します。
- Linuxをインストールしただけの状態では、できることが限られています。
 - 通常は、マッピングなど各種プログラムをインストールして利用します。
- R本体をインストールしただけの状態では、できることが限られています。
 - NGS解析を行う各種パッケージ(またはライブラリ)をインストールして利用します。

「R本体」と「パッケージ」の関係は、「パソコン」と「ソフト」、「Microsoft EXCEL」と「アドイン」、「Cytoscape」と「プラグイン」のようなものという理解でよい。

CRANとBioconductor

- R上で利用可能なパッケージの2大リポジトリ(貯蔵庫)
 - CRAN (The Comprehensive R Archive Network): 5,802パッケージ
 - Bioconductor: 824パッケージ

- 解析 | 制限酵素切断部位(RECS)地図 | [REDseq\(Zhu 201X\)](#) (last modified 2011/12)
- 解析 | small RNA | [segmentSeq\(Hardcastle 2012\)](#) (last modified 2014/02/04)
- 作図 | [ggplot2](#)について (last modified 2012/09/10)
- 作図 | [M-A plot](#)(基本編) (last modified 2012/10/01)
- 作図 | [M-A plot](#)(ggplot2編) (last modified 2012/10/01)
- 作図 | [ROC曲線](#) (last modified 2012/10/01)
- 作図 | [SplicingGraphs](#) (last modified 2012/10/01)
- [パイプライン](#) | [ggplot2](#)について (last modified 2012/10/01)
- [パイプライン](#) | ゲノム | 発現変動 | [segmentSeq](#) (last modified 2014/02/04)
- [パイプライン](#) | ゲノム | 機能解析 | [segmentSeq](#) (last modified 2014/02/04)
- [パイプライン](#) | ゲノム | 機能解析 | [segmentSeq](#) (last modified 2014/02/04)
- [パイプライン](#) | ゲノム | small RNA | [segmentSeq](#) (last modified 2014/02/04)
- [リンク集](#) (last modified 2012/03/29)

リンク集

- [R](#)
- [Bioconductor](#): [Gentleman et al., Genome Biol., 2004](#)
- [CRAN](#)
- [RjpWiki](#)
- [R Tips](#)(竹澤様)
- [BioEdit](#)(フリーの配列編集ソフト)
- [BioMart](#): [Smedley et al., BMC Genomics, 2009](#)
- [DDBJ Read Annotation Pipeline](#): [Nagasaki et al., DNA Res., 2013](#)
- [EMBOSS explorer](#) (EMBOSSのウェブ版)
- [Biostar](#): [Parnell et al., PLoS Comput Biol., 2011](#)
- [SEQanswers](#): [Li et al., Bioinformatics, 2012](#)
- [NGS WikiBook](#): [Li et al., Brief Bioinform., 2013](#)
- [HT Sequence Analysis with R and Bioconductor](#)

CRANは生命科学分野に限らず様々な分野で利用されるパッケージを含む。NGS解析は、主にBioconductorから提供されているパッケージを利用します。

パッケージ

■ 「(Rで)塩基配列解析」のインストール手順おさらい

Rのインストールと起動 **NEW**

基本的には[こちら](#)または[こちら](#)をご覧ください。
よく分からない人でWindowsユーザーの方は以下を参考にし
14日にアップデートしたMac版のインストール手順[こちら](#)(by
インストールしたほうがよい)です。

1. Windows release版のインストールの場合:

1. [Rのインストーラ](#)を「実行」
2. 聞かれるがままに「次へ」などを押してとにかくインストールを完了させる
3. **Windows Vista**の人は(パッケージのインストール中に書き込み権限に関するエラーが出るのを避けるために)「コントロールパネル」-「ユーザーアカウント」-「ユーザーアカウント制御の有効化または無効化」で、「ユーザーアカウント制御(UAC)を使ってコンピュータの保護に役立たせる」の**チェックをあらかじめ外しておくことを強くお勧め**します。
4. インストールが無事完了したら、デスクトップに出現する「R3.X.Y(32 bitの場合; XやY中の数値はバージョンによって異なります)」または「R x64 3.X.Y(64 bitの場合)」アイコンをダブルクリックして起動
5. 以下を、「R コンソール画面」上でコピー&ペーストする。10GB程度のディスク容量を要しますが一番お手軽です。(どこからダウンロードするか?と聞かれるので、その場合は自分のいる場所から近いサイトを指定)

(Rで)塩基配列解析の推奨インストール手順は、CRANおよびBioconductorから提供されている各種パッケージを予めインストールしておくこと。たった3行のコードで2つのリポジトリから提供されているパッケージ群を一度にインストール。有線LAN接続環境でも数時間程度かかるのは数千ものパッケージをダウンロードしてインストールする作業に相当する部分だから。

```
install.packages(available.packages()[,1], dependencies=TRUE)#CRAN中にある全てのパッケージをインストール
source("http://www.bioconductor.org/biocLite.R")#おまじない
biocLite(all_group())#Bioconductor中にある全てのパッケージをインストール
biocLite("BSgenome.Athaliana.TAIR.TAIR9", suppressUpdates=TRUE)#Bioconductor中にあるBSgenome.Athaliana.TAIR.TAIR9パ
```

6. 「コントロールパネル」-「デスクトップのカスタマイズ」-「フォルダオプション」-「表示(タブ)」-「詳細設定」のところで、「登録されている拡張子は表示しない」のチェックを外してください。

- 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 2013/06/14)
- イントロ | 一般 | [任意の位置の塩基を置換](#) (last modified 2013/09/12)
- イントロ | 一般 | [指定した範囲の配列を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [指定したID\(染色体やdescription\)の配列を取得](#) (last modified 2013/06/14)
- 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)
- イントロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- 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 2014/07/18) NEW
- 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)

塩基配列を入力としてコピーで翻訳配列を得ることができたのは、Bioconductorから提供されているBiostringsというパッケージを予めインストールしておいたからです。

イントロ | 一般 | 翻訳配列(translate)を取得 NEW

塩基配列を読み込んでアミノ酸配列に翻訳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(sample1.fasta)の場合:

```

in_f <- "sample1.fasta"           #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.fasta"           #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 <- translate(fasta)         #fastaをアミノ酸配列に翻訳した結果をhogeに格納
names(hoge) <- names(fasta)     #現状では翻訳した結果のオブジェクトhogeのdescription
fasta <- hoge                   #hogeの中身を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イルに保存

```


塩基配列を読み込んでアミノ酸配列に翻訳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

塩基配列を入力としてコピーで翻訳配列を得ることができたのは、Bioconductorから提供されているBiostringsというパッケージを予めインストールしておいたからです。

1. FASTA形式ファイル(sample1.fasta)の場合:

```

in_f <- "sample1.fasta"
out_f <- "hoge1.fasta"

#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で指定したファイルの読み込み

#本番
hoge <- translate(fasta)
names(hoge) <- names(fasta)
fasta <- hoge
fasta

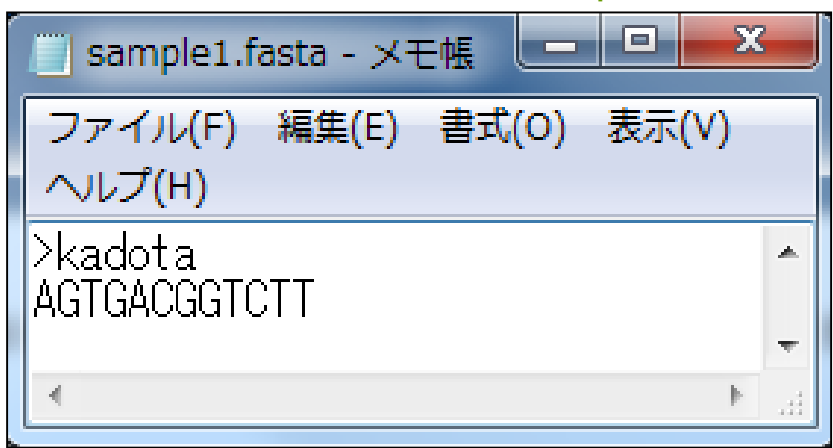
#ファイルに書き出す
writeXStringSet(fasta, out_f)
    
```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納

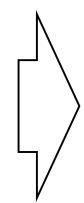
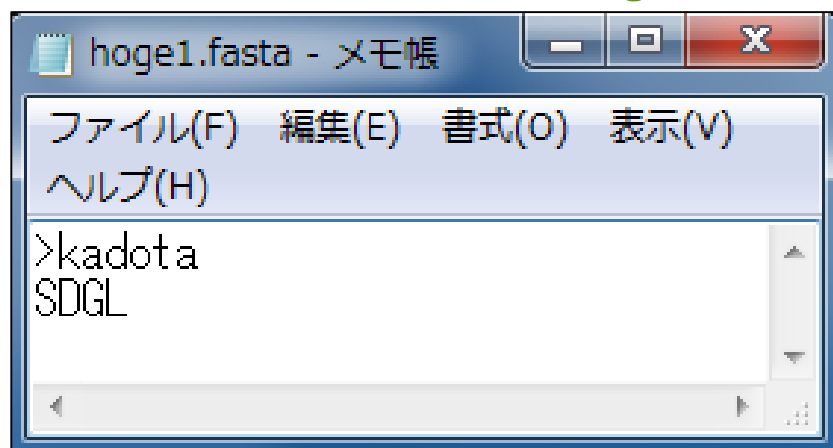
#パッケージの読み込み

#fastaをアミノ酸配列に翻訳した結果をhogeに格納
#現状では翻訳した結果のオブジェクトhogeのdescript
#hogeの中身をfastaに格納
#確認してるだけです

入力: 塩基配列ファイル(sample1.fasta)



出力: アミノ酸配列ファイル(hoge1.fasta)



塩基配列を読み込んでアミノ酸配列に翻訳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(sample1.fasta)の場合:

```

in_f <- "sample1.fasta"      #入力ファイル名を指
out_f <- "hoge1.fasta"      #出力ファイル名を指

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

#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 <- translate(fasta)    #fastaをアミノ酸配列に翻訳した結果をhogeに格納
names(hoge) <- names(fasta) #現状では翻訳した結果のオブジェクトhogeのdescript
fasta <- hoge               #hogeの中身を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
    
```

Biostringsパッケージが提供している様々な関数を利用するという宣言。Biostringsパッケージのインストールができていない状態でlibrary関数を実行しても、「Biostringsという名前のパッケージはありません」という類のエラーが出て、目的の結果は得られません。

塩基配列を読み込んでアミノ酸配列に翻訳するやり方を示し「ファイル」-「ディレクトリの変更」で解析したいファイルを選び

1. FASTA形式ファイル(sample1.fasta)の場合:

```
in_f <- "sample1.fasta"
out_f <- "hoge1.fasta"

#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(

#本番
hoge <- translate(fasta)
names(hoge) <- names(fasta)
fasta <- hoge
fasta

#ファイルに保存
writeXStringSet(fasta, fil
```

#入力
#出力

Biostringsパッケージがインストールされていても、library関数でBiostringsパッケージを読み込まないと、Biostringsパッケージで提供されている関数やサンプルデータを利用することができない。以下のコードは、library(Biostrings)の左側に#を入れて、実行されないように変更したもの。

```
#####↓
### 翻訳配列取得(Biostringsパッケージの読み込みなし)↓
#####↓
in_f <- "sample1.fasta" #入力ファイル名を指定してin_fに格納↓
out_f <- "hoge1.fasta" #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
↓
#本番↓
hoge <- translate(fasta) #fastaをアミノ酸配列に翻訳した結果をh
names(hoge) <- names(fasta) #現状では翻訳した結果のオブジェクトh
fasta <- hoge #hogeの中身を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
↓
↓
```

原因既知状態でエラーを出す

(すでに開いているRがあれば一旦終了し)
R起動直後の状態で下記コードをコピー

```
R Console
R version 3.1.0 (2014-04-10) -- "Spring Dance"
Copyright (C) 2014 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R は、自由なソフトウェアであり、「完全」
一定の条件に従えば、自由にこれを再
配布条件の詳細に関しては、'lic

R は多くの貢献者による共同プロジェクト
詳しくは 'contributors()'
また、R や R のパッケージを出版物
'citation()' と入力してください

'demo()' と入力すればデモをみ
'help()' とすればオンラインヘル
'help.start()' で HTML
'q()' と入力すれば R を終了し

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

```
#####↓
### 翻訳配列取得(Biostringsパッケージの読み込みなし)↓
#####↓
in_f <- "sample1.fasta" #入力ファイル名を指定してin_fに格納↓
out_f <- "hoge1.fasta" #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
↓
#本番↓
hoge <- translate(fasta) #fastaをアミノ酸配列に翻訳した結果をh
names(hoge) <- names(fasta) #現状では翻訳した結果のオブジェクトhoge
fasta <- hoge #hogeの中身を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
↓
↓
```

原因既知状態でエラーを出す

ファイル名: rcode_20140908.txt

```
#####↓
### 翻訳配列取得(Biostrings/
#####↓
in_f <- "sample1.fasta"
out_f <- "hogel.fasta"
↓
#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_
↓
#本番↓
hoge <- translate(fasta)
names(hoge) <- names(fasta)
fasta <- hoge
fasta
↓
#ファイルに保存↓
writeXStringSet(fasta, file=o
↓
↓
```

```
R Console
> in_f <- "sample1.fasta"
> out_f <- "hogel.fasta"
>
> #必要なパッケージをロード
> #library(Biostrings) #パッケージの読み$
>
> #入力ファイルの読み込み
> fasta <- readDNAStringSet(in_f, format="fasta") #in_fで指定$
エラー: 関数 "readDNAStringSet" を見つけることができません$
>
> #本番
> hoge <- translate(fasta) #fastaをアミノ酸配$
エラー: 関数 "translate" を見つけることができませんでした
> names(hoge) <- names(fasta) #現状では翻訳した結$
エラー: オブジェクト 'fasta' がありません
> fasta <- hoge #hogeの中身をfasta$
エラー: オブジェクト 'hoge' がありません
> fasta #確認してるだけです
エラー: オブジェクト 'fasta' がありません
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width=5$
エラー: 関数 "writeXStringSet" を見つけることができません$
> |
```

3つの関数(readDNAStringSet, translate, and writeXStringSet)がBiostringsパッケージから提供されているものであることが分かる。また、names関数はBiostringsパッケージとは無関係であることもわかる。

```

> fasta <- readDNASTringSet(in_f, format="fasta") #in_fで指定し$
エラー: 関数 "readDNASTringSet" を見つけることができませんで$
>
> #本番
> hoge <- translate(fasta) #fastaをアミノ
エラー: 関数 "translate" を見つけることができませんでした #現状では翻訳した結果$
> names(hoge) <- names(fasta) #hogeの中身を
エラー: オブジェクト 'fasta' がありません #確認してるだけ
> fasta <- hoge
エラー: オブジェクト 'hoge' がありません
> fasta
エラー: オブジェクト 'fasta' がありません
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width=50) # $
エラー: 関数 "writeXStringSet" を見つけることができませんで$
>
> fasta
エラー: オブジェクト 'fasta' がありません
> search()
[1] ".GlobalEnv"      "package:stats"    "package:graphics"
[4] "package:grDevices" "package:utils"    "package:datasets"
[7] "package:methods" "Autoloads"        "package:base"
> objects()
[1] "in_f" "out_f"
> in_f
[1] "sample1.fasta"
> out_f
[1] "hogel.fasta"
> |

```

names関数がないと文句を言われてはいないが、fastaオブジェクトがないと文句を言われている。

in_fで指定したFASTA形式ファイルを読み込むためのreadDNASTringSet関数実行段階でこけており、読み込み結果のfastaオブジェクトが作成されていないため。

```

> fasta <- readDNASTringSet(in_f, format="fasta")#in_fで指定し$
エラー: 関数 "readDNASTringSet" を見つけることができませんで$
>
> #本番
> hoge <- translate(fasta)                                #fastaをアミノ
エラー: 関数 "translate" を見つけることができませんでした
> names(hoge) <- names(fasta)                             #現状では翻訳し
エラー: オブジェクト 'fasta' がありません
> fasta <- hoge                                           #hogeの中身を
エラー: オブジェクト 'hoge' がありません
> fasta                                                    #確認してるだけ
エラー: オブジェクト 'fasta' がありません
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width=50)#$
エラー: 関数 "writeXStringSet" を見つけることができませんで$
>
> fasta
エラー: オブジェクト 'fasta' がありません
> search()
[1] ".GlobalEnv"      "package:stats"      "package:graphics"
[4] "package:grDevices" "package:utils"      "package:datasets"
[7] "package:methods" "Autoloads"          "package:base"
> objects()
[1] "in_f"  "out_f"
> in_f
[1] "sample1.fasta"
> out_f
[1] "hogel.fasta"
> |

```

search関数は現在読み込んでいるパッケージをリストアップ。このR環境は、明示的に何も読み込んでいないのにいくつか表示されている。これは、R起動時に基本的なパッケージ(stats, graphics, utils, baseなど)を自動的に読み込んでいるから。ちなみにnames関数はbaseパッケージで提供されているのでいつでも利用可能。


```

> fasta <- readDNASTringSet(in_f, format="fasta") #in_fで指定し$
エラー: 関数 "readDNASTringSet" を見つけることができませんで$
>
> #本番
> hoge <- translate(fasta) ← #fastaをアミノ酸配列に変換
エラー: 関数 "translate" を見つけることができませんでした
> names(hoge) <- names(fasta) ← #現状では翻訳結果の名前
エラー: オブジェクト 'fasta' がありません
> fasta <- hoge ← #hogeの中身をfastaに代入
エラー: オブジェクト 'hoge' がありません
> fasta ← #確認してるだけです
エラー: オブジェクト 'fasta' がありません
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width=50) # $
エラー: 関数 "writeXStringSet" を見つけることができませんでし$
>
> fasta ← #確認してるだけです
エラー: オブジェクト 'fasta' がありません
> search()
[1] ".GlobalEnv"      "package:stats"      "package:graphics"
[4] "package:grDevices" "package:utils"      "package:datasets"
[7] "package:methods" "Autoloads"          "package:base"
> objects()
[1] "in_f" "out_f"
> in_f
[1] "sample1.fasta"
> out_f
[1] "hoge1.fasta"
> |

```

translate関数もないしfastaオブジェクトもないので、translate(fasta)実行結果のhogeオブジェクトもない。ではどんなオブジェクトが利用可能かを調べるのがobjects関数。ls関数でもよい。


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

?namesを実行することで、names関数がbaseパッケージから提供されていることがわかる。

names {base} R Documentation

The Names of an Object

Description

Functions to get or set the names of an object.

Usage

```
names(x)
names(x) <- value
```

Arguments

x an R object.
value a character vector of up to the same length as **x**, or **NULL**.

Details

`names` is a generic accessor function, and `names<-` is a generic replacement function. The default methods get and set the "names" attribute of a vector (including a list) or pairlist.

If `value` is shorter than `x`, it is extended by character `NA`s to the length of `x`.

R Console

```
> search()
[1] ".GlobalEnv"          "package:stats"
[3] "package:graphics"    "package:grDevices"
[5] "package:utils"        "package:datasets"
[7] "package:methods"     "Autoloads"
[9] "package:base"
```

```
> ?readDNAStringSet
```

```
No documentation for 'readDNAStringSet' in specified packages and libraries:
you could try '??readDNAStringSet'
```

```
> library(Biostrings)
```

```
要求されたパッケージ BiocGenerics をロード中です
```

```
要求されたパッケージ parallel をロード中です
```

```
次のパッケージを付け加えます: 'BiocGenerics'
```

```
以下のオブジェクトはマスクされています (from 'package:parallel') :
```

```
clusterApply, clusterApplyLB, clusterCall, clusterEvalQ,
clusterExport, clusterMap, parApply, parCapply, parLapply,
parLapplyLB, parRapply, parSapply, parSapplyLB
```

```
以下のオブジェクトはマスクされています (from 'package:stats') :
```

```
xtabs
```

?関数名でhtmlマニュアルを見ることが
できるは、search()で表示されているパッ
ケージで提供されている関数のみ。
Biostringsパッケージを読み込んでない段
階では、Biostringsが提供している関数の
マニュアルを見ることはできない。

R Console

```
anyDuplicated, append, as.data.frame, as.  
cbind, colnames, do.call, duplicated, eval  
Filter, Find, get, intersect, is.unsorted,  
Map, mapply, match, mget, order, paste, p  
pmax.int, pmin, pmin.int, Position, rank,  
Reduce, rep.int, rownames, sapply, setdiff  
table, tapply, union, unique, unlist
```

要求されたパッケージ IRanges をロード中です

要求されたパッケージ XVector をロード中です

```
> search()
```

```
[1] ".GlobalEnv"           "package:Biostrings"  
[3] "package:XVector"      "package:IRanges"  
[5] "package:BiocGenerics" "package:parallel"  
[7] "package:stats"        "package:graphics"  
[9] "package:grDevices"    "package:utils"  
[11] "package:datasets"     "package:methods"  
[13] "Autoloads"            "package:base"  
> |
```

library(Biostrings)実行によって、search()で表示されているパッケージリストにBiostringsを含むいくつかのパッケージが追加されていることが分かる。この段階で?readDNAStringSetが有効となり、htmlマニュアルが開くようになる。

```
> search()
[1] ".GlobalEnv"          "package:Biostrings"
[3] "package:XVector"     "package:IRanges"
[5] "package:BiocGenerics" "package:parallel"
[7] "package:stats"       "package:graphics"
[9] "package:grDevices"   "package:utils"
[11] "package:datasets"    "package:methods"
[13] "Autoloads"           "package:base"
> ?readDNAStringSet
starting httpd help server ... done
> |
```

library(Biostrings)実行によって、search()で表示されているパッケージリストにBiostringsを含むいくつかのパッケージが追加されていることが分かる。この段階で?readDNAStringSetが有効となり、htmlマニュアルが開くようになる。

XStringSet-io {Biostrings}

R Documentation

Read/write an XStringSet object from/to a file

Description

Functions to read/write an [XStringSet](#) object from/to a file.

Usage

```
## Read FASTA (or FASTQ) files in an XStringSet object:
```

```
readBStringSet(filepath, format="fasta",
               nrec=-1L, skip=0L, seek.first.rec=FALSE, use.names=TRUE)
readDNAStringSet(filepath, format="fasta",
                 nrec=-1L, skip=0L, seek.first.rec=FALSE, use.names=TRUE)
readRNAStringSet(filepath, format="fasta",
                 nrec=-1L, skip=0L, seek.first.rec=FALSE, use.names=TRUE)
readAAStringSet(filepath, format="fasta",
                 nrec=-1L, skip=0L, seek.first.rec=FALSE, use.names=TRUE)
```

```
## Extract basic information about FASTA (or FASTQ) files
```

```
## without actually loading the sequence data:
```

```
fasta_info(filepath)
```

もう一度コピー

ファイル名: rcode_20140908.txt

#でコメントアウトしているにも関わらずエラーが出ないのは、以前にBiostringsパッケージを読み込んでいるから。

```
#####↓
### 翻訳配列取得(Biostrings/パッケージの読み込みなし)↓
#####↓
in_f <- "sample1.fasta"
out_f <- "hogel.fasta"
↓
#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(
↓
#本番↓
hoge <- translate(fasta)
names(hoge) <- names(fasta)
fasta <- hoge
fasta
↓
#ファイルに保存↓
writeXStringSet(fasta, fil
↓
↓
```

```
R Console
> in_f <- "sample1.fasta" #入力ファイル名を指定して$
> out_f <- "hogel.fasta" #出力ファイル名を指定して$
> #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$
> #本番
> hoge <- translate(fasta) #fastaをアミノ酸配列に翻$
> names(hoge) <- names(fasta) #現状では翻訳した結果のオ$
> fasta <- hoge #hogeの中身をfastaに格納
> fasta #確認してるだけです
A AASTringSet instance of length 1
width seq names
[1] 4 SDGL kadota
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width=50)#fas$
> |
```

利用可能なオブジェクトを表示

```

R Console
> hoge <- translate(fasta) #fastaをアミノ酸配列に翻訳
> names(hoge) <- names(fasta) #現状では翻訳した結果のオブジェクト
> fasta <- hoge #hogeの中身をfastaに格納
> fasta #確認してるだけです
A AAStringSet instance of length 1
  width seq          names
[1]      4 SDGL      kadota
>
> #ファイルに保存
> writeXStringSet(fasta, file=out f, format="fasta", width=50)#fasta
> objects()
[1] "fasta" "hoge" "in_f" "out_f"
> fasta
A AAStringSet instance of length 1
  width seq          names
[1]      4 SDGL      kadota
> hoge
A AAStringSet instance of length 1
  width seq          names
[1]      4 SDGL      kadota
> |

```

翻訳配列取得がエラーなく終了すると4つのオブジェクトが利用可能なようです。

Contents

- 3-3. R 各種パッケージ、2014/09/08 15:00-18:15、中級、実習
 - パッケージ
 - 「(Rで)塩基配列解析」のインストール手順おさらい
 - CRANとBioconductor
 - 代表的なパッケージBiostringsの利用法: library, search, objects関数
 - 作業スペース(workspace)の概念
 - Biostringsパッケージで利用可能な関数を概観
 - source関数の利用



Rの終了と作業スペース

The screenshot shows the R Console window with the following code and output:

```
> names(hoge) <- names(fasta)
> fasta <- hoge
> fasta
  A AAStringSet instance of length 1
    width seq          names
[1]      4 SDGL      kadota
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_1.fasta)
> objects()
[1] "fasta" "hoge"  "in_f"  "out_1"
> fasta
  A AAStringSet instance of length 1
    width seq          names
[1]      4 SDGL      kadota
> hoge
  A AAStringSet instance of length 1
    width seq          names
[1]      4 SDGL      kadota
> q()
```

Red annotations on the right side of the console:

- #現状では翻訳した結果の才\$
- #hogeの中身をfastaに格納
- #確認してるだけです

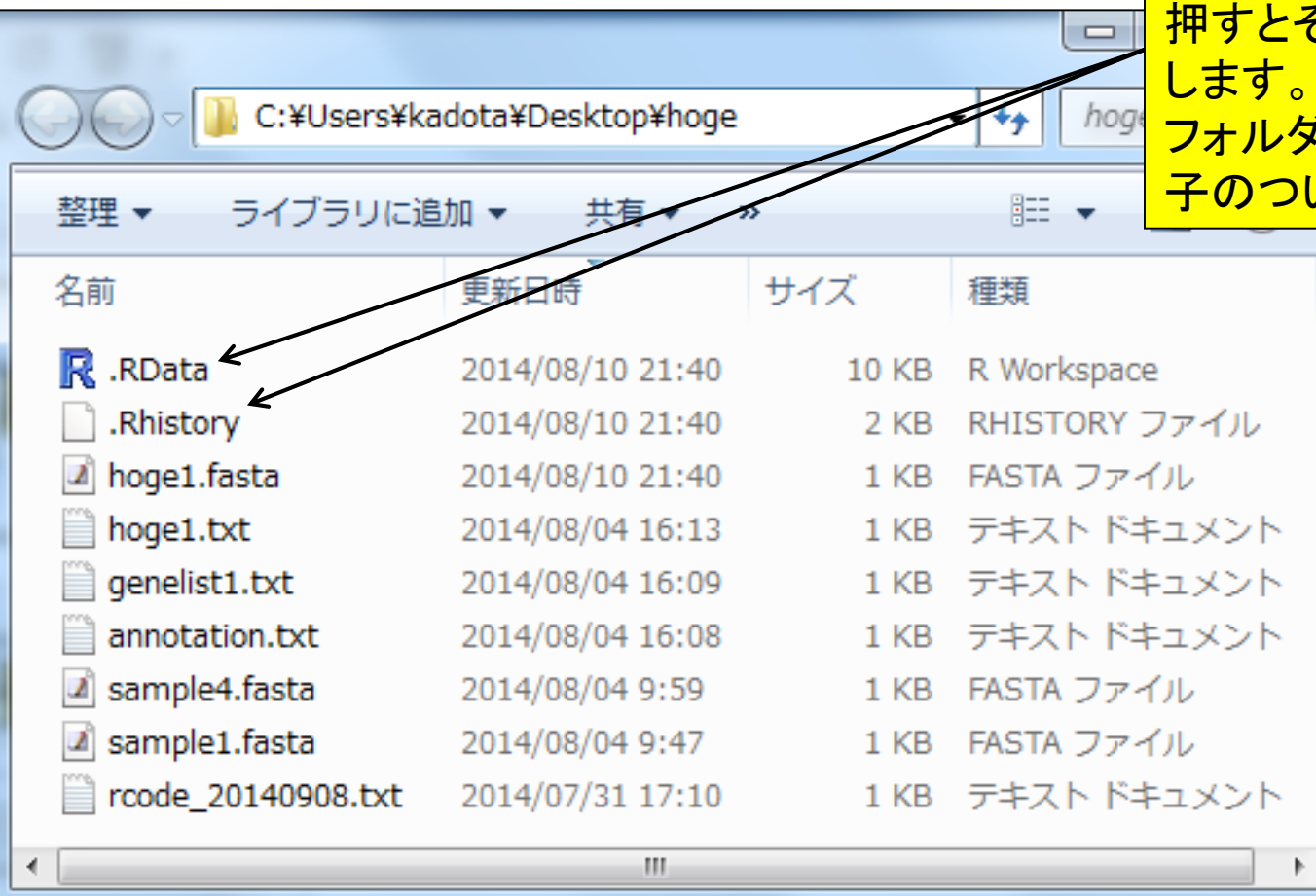
A dialog box titled "質問" (Question) is overlaid on the console, asking "作業スペースを保存しますか?" (Save workspace?). The "はい(Y)" (Yes) button is highlighted with a red arrow.

Annotations on the screenshot:

- A yellow box on the right says: "作業スペース(workspace) 保存の意味を説明します" (Explains the meaning of saving the workspace).
- A white box at the bottom left says: "R Gui画面右上の×ボタン以外にq()でもRを終了させることができます。" (Besides the × button in the top right of the R Gui window, you can also exit R with q()).

Rの終了と作業スペース

R環境(バージョンやOS)によって若干違うかもしれませんが、私の環境では「はい」を押すとそれ以上何も聞かれることなく終了します。そして、作業ディレクトリだったhogeフォルダ中に.RDataと.Rhistoryという拡張子のついたファイルが自動生成されます。



Rの起動と作業スペースの読み込み

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... (3)
 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...
 印刷...
 ファイルを保存...
 終了

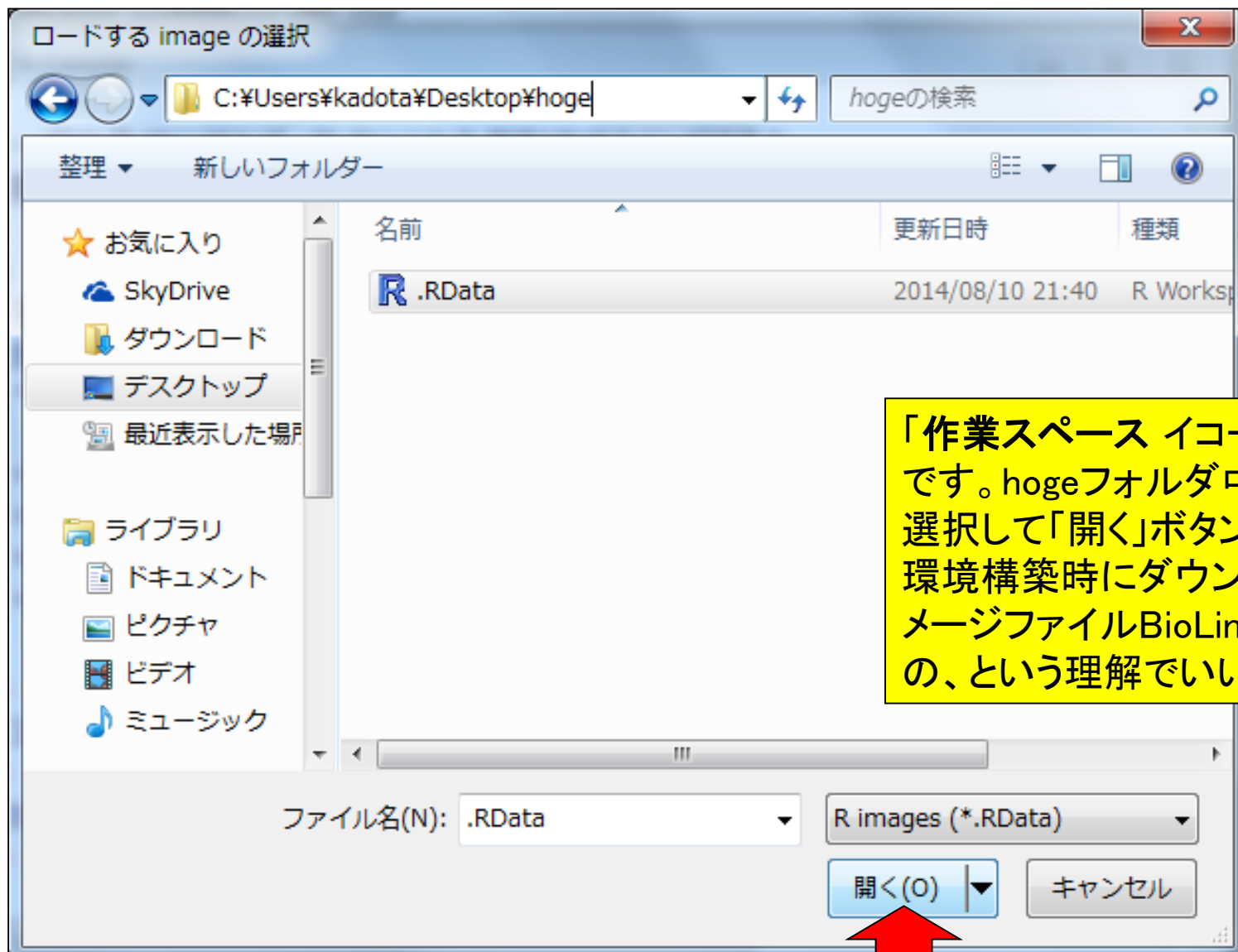
```

> getwd() (1)
[1] "C:/Users/kadota/Desktop"
> objects() (2)
character(0)
> |
    
```

再配布することができ\$
 ense()' あるいは 'lice\$
 クトです。
 と入力してください。
 物で引用する際の形式\$
 ない。
 ることができます。
 プが出ます。
 ブラウザによるヘルプがみられ\$
 'q()' と入力すれば R を終了します。

Rを起動して、(.RDataファイルが存在する
 hogeフォルダでなくてもいいことを示すため
 に)①作業ディレクトリをデスクトップにした状
 態で②利用可能なオブジェクトを調べます。
 R起動直後は利用可能なオブジェクトがない
 ことを確認した上で、③「ファイル - 作業ス
 ペースの読み込み」を選択します。

Rの起動と作業スペースの読み込み



「作業スペース イコール .RDataファイル」です。hogeフォルダ中の.RDataファイルを選択して「開く」ボタンを押す。これはLinux環境構築時にダウンロードしてもらったイメージファイルBioLinux.ovaと同じようなもの、という理解でいいです。

作業スペースをロードした直後の状態

```

R Console

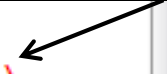
配布条件の詳細に関しては、'license()' あるいは 'lice$

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$
'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
'q()' と入力すれば R を終了します。

> getwd()
[1] "C:/Users/kadota/Desktop"
> objects()
character(0)
> load("C:\\Users\\kadota\\Desktop\\hoge\\.RData")
> |
    
```

自動的にloadという関数と.RDataが存在する場所(パス)をC:ドライブから指示したものが表示されます。つまり、「ファイル - 作業スペースの読み込み」がload関数の記述に相当し、その後の.Rdataの場所を選んで開くボタンを押す作業がload関数の括弧の中身の記述に相当する、ということです。



作業スペースロード後の状態を確認

```

R Console
> getwd()
[1] "C:/Users/kadota/Desktop"
> objects()
character(0)
> load("C:\\Users\\kadota\\Desktop\\hoge\\.RData")
> getwd()
[1] "C:/Users/kadota/Desktop"
> objects()
[1] "fasta" "hoge"  "in_f"  "out_f"
> in_f
[1] "sample1.fasta"
> out_f
[1] "hoge1.fasta"
> search()
[1] ".GlobalEnv"          "package:stats"
[3] "package:graphics"   "package:grDevices"
[5] "package:utils"      "package:datasets"
[7] "package:methods"    "Autoloads"
[9] "package:base"
> |
    
```

作業スペースロード後でも、作業ディレクトリ自体は以前行っていたhogeフォルダに変更されるわけではない。

作業スペースをロードする主なご利益は、以前の作業スペースで利用可能だったオブジェクトが今の作業スペース中でも利用可能になる点です。得るのに非常に計算時間がかかるオブジェクトを再利用したい場合などに便利。

search()実行結果からわかるように、以前行っていたlibrary(Biostrings)は反映されないようです。

作業スペースロード後の状態を確認

```
R Console
> pre <- search()
> pre
[1] ".GlobalEnv"          "package:stats"
[3] "package:graphics"   "package:grDevices"
[5] "package:utils"      "package:datasets"
[7] "package:methods"    "Autoloads"
[9] "package:base"
> fasta
Loading required package: Biostrings
Loading required package: BiocGenerics
Loading required package: parallel

Attaching package: 'BiocGenerics'

以下のオブジェクトはマスクされています (from 'pac$
  clusterApply, clusterApplyLB,
  clusterCall, clusterEvalQ, clusterExport,
  clusterMap, parApply, parCapply,
  parLapply, parLapplyLB, parRapply,
  parSapply, parSapplyLB

以下のオブジェクトはマスクされています (from 'pac$
```

```
R Console
xta
以下のオブ
any
as.v
dupl
get,
map
pmax
rbind, reduce, rep.int, rownames, sapply,
setdiff, sort, table, tapply, union,
unique, unlist

Loading required package: IRanges
Loading required package: XVector
  A AAStringSet instance of length 1
    width seq          names
[1] 4 SDGL             kadota
> fasta
  A AAStringSet instance of length 1
    width seq          names
[1] 4 SDGL             kadota
> |
```

fastaオブジェクトも以前と同じものが利用可能であることが分かる。fastaオブジェクトはBiostringsパッケージ中の関数を用いて得られたものであり、Biostringsパッケージ中で定義された専用の形式(AAStringSet)で記述されている。それゆえ、fastaの中身を表示する際に、自動的に必要なパッケージをロードしてくれるようだ。

作業スペースロード後の状態を確認

```
R Console
> post <- search()
> post
 [1] ".GlobalEnv"          "package:Biostrings"  "package:XVector"
 [4] "package:IRanges"    "package:BiocGenerics" "package:parallel"
 [7] "package:stats"      "package:graphics"    "package:grDevices"
[10] "package:utils"      "package:datasets"    "package:methods"
[13] "Autoloads"          "package:base"
> pre
 [1] ".GlobalEnv"          "package:stats"        "package:graphics"
 [4] "package:grDevices"  "package:utils"        "package:datasets"
 [7] "package:methods"    "Autoloads"            "package:base"
> setdiff(post, pre)
 [1] "package:Biostrings"  "package:XVector"      "package:IRanges"
 [4] "package:BiocGenerics" "package:parallel"
> setdiff(pre, post)
character(0)
> intersect(pre, post)
 [1] ".GlobalEnv"          "package:stats"        "package:graphics"
 [4] "package:grDevices"  "package:utils"        "package:datasets"
 [7] "package:methods"    "Autoloads"            "package:base"
> |
```

fastaオブジェクト表示の際に必要なパッケージをロードする前後のパッケージリストをpreおよびpostオブジェクトで保持している。集合演算は、ロード後にどんなパッケージが利用可能になっているかを調べる場合にも便利。setdiff関数は、左側のベクトル中で右側のベクトル中の要素として含まれないものを表示。intersect関数は共通要素を表示。

Contents

- 3-3. R 各種パッケージ、2014/09/08 15:00-18:15、中級、実習
 - パッケージ
 - 「(Rで)塩基配列解析」のインストール手順おさらい
 - CRANとBioconductor
 - 代表的なパッケージBiostringsの利用法: library, search, objects関数
 - 作業スペース(workspace)の概念
 - Biostringsパッケージで利用可能な関数を概観
 - source関数の利用



Biostringsで利用可能な関数を概観

- 書籍 | 日本乳酸菌学会誌 | [第1回イントロダクション](#) (last modified 2014/07/07) **NEW**
- イントロ | 一般 | [ランダムに行を抽出](#) (last modified 2014/07/17) **NEW**
- イントロ | 一般 | [任意の文字列を行の最初に挿入](#) (last modified 2014/07/17) **NEW**
- 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 2014/04/11)
- イントロ | 一般 | [ランダムな塩基配列を生成](#) (last modified 2014/06/16)
- 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 2013/06/14)
- イントロ | 一般 | [任意の位置の塩基を置換](#) (last modified 2013/09/12)
- イントロ | 一般 | [指定した範囲の配列を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [指定したID\(染色体やdescription\)の配列を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [翻訳配列\(translate\)を取得](#) (last modified 2014/03/08) **NEW**
- イントロ | 一般 | [相補鎖\(complement\)を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [逆鎖\(reverse\)を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [2連続塩基の出現頻度情報を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [3連続塩基の出現頻度情報を取得](#) (last modified 2014/03/08)
- 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 2014/03/08)
- イントロ | 一般 | [Tips | 任意の拡張子でファイルを保存](#) (last modified 2014/03/08)
- イントロ | 一般 | [Tips | 拡張子は同じで任意の文字列を保存](#) (last modified 2014/03/08)
- イントロ | 一般 | [配列取得 | ゲノム配列 | 公共DBから取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [配列取得 | ゲノム配列 | BSgenomeパッケージ](#) (last modified 2014/03/08)
- イントロ | 一般 | [配列取得 | プロモーター配列 | 公共DBから取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [配列取得 | プロモーター配列 | BSgenomeパッケージ](#) (last modified 2014/03/08)

Biostringsパッケージは、翻訳配列取得以外にも様々な塩基配列解析用の関数を提供しています。

イントロ | 一般 | [翻訳配列\(translate\)を取得](#) **NEW**

塩基配列を読み込んでアミノ酸配列に翻訳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(sample1.fasta)の場合:

```

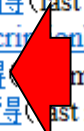
in_f <- "sample1.fasta" #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.fasta" #出力ファイル名を指定してout_fに格納

#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 <- readDNAStrngSet(in_f, format="fasta")#in_fで指定したファイルの読み込み

#本番
hoge <- translate(fasta) #fastaをアミノ酸配列に翻訳した結果をhogeに格納
names(hoge) <- names(fasta) #現状では翻訳した結果のオブジェクトhogeのdescription
fasta <- hoge #hogeの中身を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イルに保存
    
```



Biostringsで利用可能な関数を概観

- 書籍 | 日本乳酸菌学会誌 | [第1回イントロダクション](#) (last modified 2014/07/07) **NEW**
- イントロ | 一般 | [ランダムに行を抽出](#) (last modified 2014/07/17) **NEW**
- イントロ | 一般 | [任意の文字列を行の最初に挿入](#) (last modified 2014/07/17) **NEW**
- 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 2014/04/11)
- イントロ | 一般 | [ランダムな塩基配列を生成](#) (last modified 2014/06/16)
- 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 2013/06/14)
- イントロ | 一般 | [任意の位置の塩基を置換](#) (last modified 2013/09/12)
- イントロ | 一般 | [指定した範囲の配列を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [指定したID\(染色体やdescription\)の配列を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [翻訳配列\(translate\)を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [相補鎖\(complement\)を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [逆鎖\(reverse\)を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [2連続塩基の出現頻度情報を取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [3連続塩基の出現頻度情報を取得](#) (last modified 2014/03/08)
- 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 2014/03/08)
- イントロ | 一般 | [Tips | 任意の拡張子でファイルを保存](#) (last modified 2014/03/08)
- イントロ | 一般 | [Tips | 拡張子は同じで任意の文字列を保存](#) (last modified 2014/03/08)
- イントロ | 一般 | [配列取得 | ゲノム配列 | 公共DBから取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [配列取得 | ゲノム配列 | BSgenomeパッケージ](#) (last modified 2014/03/08)
- イントロ | 一般 | [配列取得 | プロモーター配列 | 公共DBから取得](#) (last modified 2014/03/08)
- イントロ | 一般 | [配列取得 | プロモーター配列 | BSgenomeパッケージ](#) (last modified 2014/03/08)

相補鎖を取得する場合も、基本的な枠組みは翻訳配列取得と同じ。主要な関数がtranslateからcomplementに変わったただけです。

イントロ | 一般 | 相補鎖(complement)を取得

FASTA形式ファイルを読み込んで相補鎖を得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(sample1.fasta)の場合:

```

in_f <- "sample1.fasta" #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.fasta" #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
#確認してるだけです

#本番
fasta <- complement(fasta) #fastaオブジェクトの相補鎖をfastaに格納
#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
    
```

Biostringsで利用可能な関数を概観

- 書籍 | 日本乳酸菌学会誌 | [第1回イントロダクション](#) (last modified 2014/07/07) **NEW**
- インポート | 一般 | [ランダムに行を抽出](#) (last modified 2014/07/17) **NEW**
- インポート | 一般 | [任意の文字列を行の最初に挿入](#) (last modified 2014/07/17) **NEW**
- 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 2014/04/11)
- インポート | 一般 | [ランダムな塩基配列を生成](#) (last modified 2014/06/16)
- 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 2014/03/08)
- インポート | 一般 | [指定したID\(染色体やdescription\)の配列を取得](#) (last modified 2014/03/10)
- インポート | 一般 | [翻訳配列\(translate\)を取得](#) (last modified 2013/06/14)
- インポート | 一般 | [相補鎖\(complement\)を取得](#) (last modified 2013/06/14)
- インポート | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- インポート | 一般 | [逆鎖\(reverse\)を取得](#) (last modified 2013/06/14)
- インポート | 一般 | [2連続塩基の出現頻度情報を取得](#) (last modified 2014/07/18) **NEW**
- インポート | 一般 | [3連続塩基の出現頻度情報を取得](#) (last modified 2013/06/14)
- インポート | 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)
- インポート | 一般 | **Tips** | [任意の拡張子でファイルを保存](#) (last modified 2013/09/26)
- インポート | 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)
- インポート | 一般 | [配列取得 | ゲノム配列 | 公共DBから](#) (last modified 2014/05/28)
- インポート | 一般 | [配列取得 | ゲノム配列 | BSgenome](#) (last modified 2014/06/28) **NEW**
- インポート | 一般 | [配列取得 | プロモーター配列 | 公共DBから](#) (last modified 2014/04/02)
- インポート | 一般 | [配列取得 | プロモーター配列 | BSgenome](#) (last modified 2014/04/25)
- インポート | 一般 | [配列取得 | プロモーター配列 | GenomicFeatures\(Lawrence 2013\)](#) (last modified 2014/04/23)
- インポート | 一般 | [配列取得 | トランスクリプトーム配列 | 公共DBから](#) (last modified 2014/04/02)
- インポート | 一般 | [配列取得 | トランスクリプトーム配列 | biomaRt\(Durinck 2009\)](#) (last modified 2013/09/25)
- インポート | NGS | [様々なプラットフォーム](#) (last modified 2014/06/10)
- インポート | NGS | [qPCRやmicroarrayなどとの比較](#) (last modified 2014/07/11) **NEW**
- インポート | NGS | [可視化\(ゲノムブラウザやViewer\)](#) (last modified 2014/06/25) **NEW**
- インポート | NGS | [配列取得 | FASTQ or SRALite | 公共DBから](#) (last modified 2014/06/28) **NEW**
- インポート | NGS | [配列取得 | FASTQ or SRALite | SRADB\(Zhu 2013\)](#) (last modified 2014/06/26) **NEW**
- **イントロ** | NGS | [配列取得 | シミュレーションデータ | について](#) (last modified 2014/06/25) **NEW**
- インポート | NGS | [配列取得 | シミュレーションデータ | ランダムな塩基配列の生成から](#) (last modified 2014/06/25)

(Rで)塩基配列解析上で示している項目はごく一部です。Biostringsパッケージが提供している全ての関数を網羅しているわけでもありません。

Biostringsで利用可能な関数を概観

イントロ | 一般 | 相補鎖(complement)を取得

FASTA形式ファイルを読み込んで相補鎖を得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([sample1.fasta](#))の場合:

```
in_f <- "sample1.fasta"      #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.fasta"      #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
fasta                                #確認してるだけです

#本番
fasta <- complement(fasta)   #fastaオブジェクトの相補鎖を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
```

2. multi-FASTA形式ファイル([h_rna.fasta](#))の場合:

```
in_f <- "h_rna.fasta"      #入力ファイル名を指定してin_fに格納
out_f <- "hoge2.fasta"    #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
fasta                                #確認してるだけです

#本番
fasta <- complement(fasta)   #fastaオブジェクトの相補鎖を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
```

「必要なパッケージをロード」のところで、Biostringsを明示的に読み込んでいる場合は、項目の最後のところにBioconductorのBiostringsパッケージのウェブサイトリンクを張っている。そこを眺めるとBiostringsが提供している関数のリストを概観できます。

[Home](#) » [Bioconductor 2.14](#) » [Software Packages](#) » [Biostrings](#)

Biostrings

String objects representing biological sequences, and matching algorithms

Bioconductor version: Release (2.14)

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. Abouyou, R. Gentleman, and S. DebRoy

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

Citation (from within R, enter `citation("Biostrings")`):

Pages H, Abouyou P, Gentleman R and DebRoy S. *Biostrings: String objects representing biological sequences, and matching algorithms*. R package version 2.32.1.

Installation

To install this package, start R and enter:

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

Documentation

To view documentation for the version of this package installed in your system, start R and enter:

```
browseVignettes("Biostrings")
```

PDF	Script	A short presentation of the basic classes defined in Biostrings 2
PDF	Script	Biostrings Quick Overview
PDF	Script	Handling probe sequence information
PDF	R Script	Multiple Alignments
PDF	R Script	Pairwise Sequence Alignments
PDF		Reference Manual
Text		NEWS

Workflows »

Common Bioconductor workflows include:

- [Oligonucleotide Arrays](#)
- [High-throughput Sequencing](#)
- [Counting Reads for Differential Expression](#) (parathyroidSE vignette)
- [Annotation](#)
- [Annotating Variants](#)
- [Annotating Ranges](#)
- [Flow Cytometry](#) and other assays
- [Candidate Binding Sites for Known Transcription Factors](#)
- [Cloud-enabled cis-eQTL search and annotation](#)

Mailing Lists »

Post questions about Bioconductor packages to our mailing lists. Read the [posting guide](#) before posting!

- [bioconductor](#)
- [bioc-devel](#)

(実習用PC環境はインストール済みだが)Biostringsパッケージを個別にインストールしたい場合

下から2番目のReference Manualは、おそらくどのパッケージにも存在する。一番下のNEWSがないパッケージもある。

Biostringsはここが5つ存在するが、パッケージによっては1つだけの場合もある。ここは開発者次第。

Biostringsで利用可能な関数を概観

Biostrings Quick Overview

Hervé Pagès
Fred Hutchinson Cancer Research Center
Seattle, WA

July 4, 2014

DNA配列を入力としてアミノ酸配列を出力することが分かっている **translate**関数を眺めるとFunction列とDescription列の意味がよくわかる。

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

Function	Description
length	Return the number of sequences in an object.
names	Return the names of the sequences in an object.
[Extract sequences from an object.
head, tail	Extract the first or last sequences from an object.

rev	Reverse the sequences in an object.
c	Concatenate sequences in an object.
width, nchar	Return the width of the sequences in an object.
==, !=	Element-wise comparison of sequences in an object.
match, %in%	Match sequences in an object against a set of sequences.
duplicated, unique	Identify duplicated or unique sequences in an object.
sort, order	Sort or order sequences in an object.
relist, split, extractList	Reorganize sequences in an object.

Table 1: Low-level manipulation

Function	Description
reverse	Compute the reverse, complement, or reverse-complement, of a set of DNA sequences.
complement	
reverseComplement	
translate	Translate a set of DNA sequences into a set of Amino Acid sequences.
chartr	Translate the letters in a set of sequences.
subseq, subseq<-	Extract/replace arbitrary substrings from/in a string or set of strings.
extractAt, replaceAt	
replaceLetterAt	Replace the letters specified by a set of positions by new letters.
padAndClip, stackStrings	Pad and clip strings.
unstrsplit	A fast implementation of <code>sapply(x, paste0, collapse=sep)</code> for collapsing the list elements of a <i>DNASetList</i> or <i>AASetList</i> object.

Table 3: Sequence transformation and editing.

Contents

- 3-3. R 各種パッケージ、2014/09/08 15:00-18:15、中級、実習
 - パッケージ
 - 「(Rで)塩基配列解析」のインストール手順おさらい
 - CRANとBioconductor
 - 代表的なパッケージBiostringsの利用法: library, search, objects関数
 - 作業スペース(workspace)の概念
 - Biostringsパッケージで利用可能な関数を概観
 - source関数の利用



source関数を利用

イントロ | 一般 | [翻訳配列\(translate\)を取得](#) **NEW**

塩基配列を読み込んでアミノ酸配列に翻訳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

1. FASTA形式ファイル(sample1.fasta)の場合:

```
in_f <- "sample1.fasta"
out_f <- "hoge1.fasta"

#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")#

#本番
hoge <- translate(fasta)
names(hoge) <- names(fasta)
fasta <- hoge
fasta

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

#入力ファイル名を指定してin_fに
#出力ファイル名を指定してout_fに格納

#パッケージ

#fastaをア
#現状では番
#hogeの中
#確認して

RでもLinuxのシェルスクリプトっぽいことができます。基本手順は①コピーする一連のコードからなるファイル(rcode_translate.txt)を入力ファイル(sample1.fasta)と同じフォルダ上に保存し、②作業ディレクトリの変更を行って、③source("rcode_translate.txt")を実行。

```
rcode_translate.txt x
#####↓
### 翻訳配列取得↓
#####↓
in_f <- "sample1.fasta"
out_f <- "hoge1.fasta"
↓
#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で指定したファイルの読み込
↓
#本番↓
hoge <- translate(fasta)
names(hoge) <- names(fasta)
fasta <- hoge
fasta
↓
#ファイルに保存↓
writeXStringSet(fasta, file=out_f, format="fasta", width=50)#fastaの中身を指定
↓
```

#入力ファイル名を指定してin_fに格納↓
#出力ファイル名を指定してout_fに格納↓

#パッケージの読み込み↓

#fastaをアミノ酸配列に翻訳した結果をhoge
#現状では翻訳した結果のオブジェクトhoge
#hogeの中身をfastaに格納↓
#確認してるだけです↓

source関数を利用

① rcode_translate.txtと入力ファイル sample1.fastaをデスクトップ上のhoge2フォルダ中に保存し、②Rを起動して作業ディレクトリの変更を行って、③source("rcode_translate.txt")を実行した結果。Source関数実行後に出カファイルhoge1.fastaが生成されていることがわかる

以下の

以下のオブジェクトはマスクされています (from \$

```
anyDuplicated, append, as.data.frame,
as.vector, cbind, colnames, do.call,
duplicated, eval, evalq, Filter, Find,
get, intersect, is.unsorted, lapply,
Map, mapply, match, mget, order,
paste, pmax, pmax.int, pmin, pmin.int,
Position, rank, rbind, Reduce,
rep.int, rownames, sapply, setdiff,
sort, table, tapply, union, unique,
unlist
```

要求されたパッケージ IRanges をロード中です

要求されたパッケージ XVector をロード中です

```
> list.files()
[1] "hoge1.fasta"           "rcode_translate.txt"
[3] "sample1.fasta"
> |
```

```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge2"
> list.files()
[1] "rcode_translate.txt" "sample1.fasta"
> source("rcode_translate.txt")
要求されたパッケージ BiocGenerics をロード中です
要求されたパッケージ parallel をロード中です

次のパッケージを付け加えます: 'BiocGenerics'

以下のオブジェクトはマスクされています (from $

  clusterApply, clusterApplyLB,
  clusterCall, clusterEvalQ,
  clusterExport, clusterMap, parApply,
  parCapply, parLapply, parLapplyLB,
  parRapply, parSapply, parSapplyLB

以下のオブジェクトはマスクされています (from $

  xtabs

以下のオブジェクトはマスクされています (from $
```

source関数を利用

うまくいく例

```
R Console
> list.files()
[1] "rcode_translate.txt"
[2] "sample1.fasta"
> source("rcode_translate.txt")
> list.files()
[1] "hogel.fasta"
[2] "rcode_translate.txt"
[3] "sample1.fasta"
>
```

うまくいかない例

```
R Console
> list.files()
[1] "rcode_translate.txt"
[2] "sample1.fasta"
> source("rcode_translate.txt")
> list.files()
[1] "rcode_translate.txt"
[2] "sample1.fasta"
> |
```

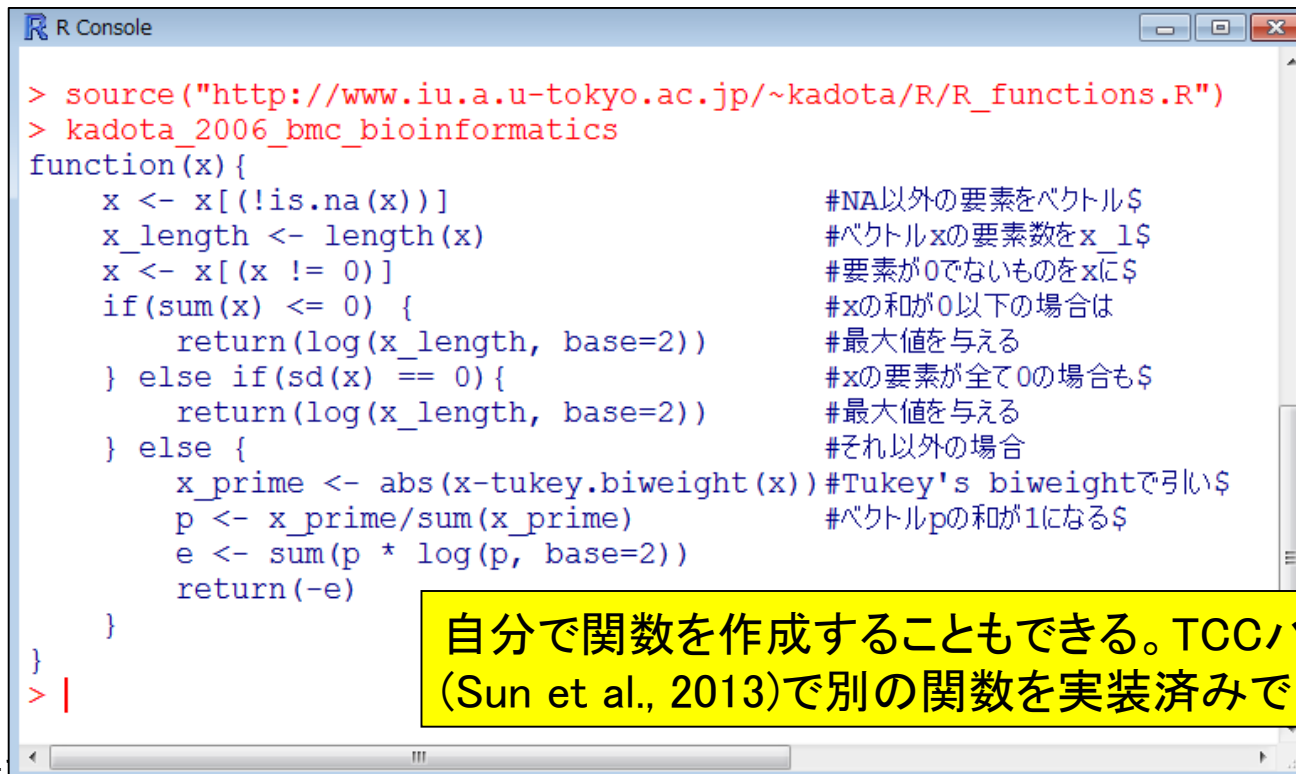
二重クォーテーションに注意!!
 「全角」ではなく「半角」です。丸括弧もそうです。エラーメッセージが出ないので、こういうこともあるという経験をしておかないと対処法が分からず苦悩します。

source関数を利用

- Biostringsのようにパッケージを作成して提供する以外に、自作のR関数をウェブ上で公開しているヒトもいる。

- http://www.iu.a.u-tokyo.ac.jp/~kadota/R/R_functions.R

- 組織特異的発現パターン検出法ROKU (Kadota et al., *BMC Bioinformatics*, 2006)を実行するためのkadota_2006_bmc_bioinformatics関数
- RNA-seq正規化法やシミュレーションデータ作成関数(Kadota et al., *AMB*, 2012)
- etc...



```
> source("http://www.iu.a.u-tokyo.ac.jp/~kadota/R/R_functions.R")
> kadota_2006_bmc_bioinformatics
function(x) {
  x <- x[(!is.na(x))]           #NA以外の要素をベクトル$
  x_length <- length(x)       #ベクトルxの要素数をx_l$
  x <- x[x != 0]              #要素が0でないものをxに$
  if(sum(x) <= 0) {          #xの和が0以下の場合
    return(log(x_length, base=2)) #最大値を与える
  } else if(sd(x) == 0) {    #xの要素が全て0の場合も$
    return(log(x_length, base=2)) #最大値を与える
  } else {                  #それ以外の場合
    x_prime <- abs(x-tukey.biweight(x)) #Tukey's biweightで引い$
    p <- x_prime/sum(x_prime)         #ベクトルpの和が1になる$
    e <- sum(p * log(p, base=2))
    return(-e)
  }
}
```

自分で関数を作成することもできる。TCCパッケージ ver. 1.4 (Sun et al., 2013)で別の関数を実装済みで、現在は不使用。

解析 | 発現変動 | 多群間 | ROKU (Kadota_2006)

TCCパッケージで提供しているROKU法(Kadota et al., 2006)を用いて、遺伝子発現行列中の遺伝子を全体的な組織特異性の度合いでランキングします。出力ファイル中の"modH"列の数値は、「ROKU論文中のAdditional file 1(Suppl.xls)の"H(x)"列の数値」と対応しています。つまり、データ変換後のエントロピー値です。"ranking"列は、modHの値でランキングした結果です。"ranking"列で昇順にソートすることで、全体的な組織特異性の度合いでランキングしていることになります。つまり、上位が「どの組織で特異的かはこのスコアだけでは分からないが」組織特異性が高い遺伝子」ということとなります。残りの結果は「1:特異的高発現、-1:特異的低発現、0:その他」からなる「外れ値行列」です。例えば、組織AとBで1、それ以外の組織で0を示す遺伝子(群)は「AとB特異的高発現遺伝子」と判断します。「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. サンプルデータ21のsample21.txtの場合:

log2変換後のデータであるという前提です。

```

in_f <- "sample21.txt" #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(TCC) #パッケージの読み込み

#入力ファイルの読み込み
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="") #in_fで指定したファイルの読み込み

#本番
hoge <- ROKU(data) #ROKUを実行した結果をhoge1に格納
outlier <- hoge$outlier #外れ値行列をoutlierに格納
modH <- hoge$modH #データ変換後のエントロピー値をmodHに格納(原著論文のH(x')の値に相当)
ranking <- hoge$rank #modHでランキングした結果をrankingに格納

#ファイルに保存
tmp <- cbind(rownames(data), outlier, modH, ranking) #左端の列が遺伝子ID、次にサンプル数分だけの列からなる「外れ値行列」
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F) #tmpの中身を指定したファイル名で保存

```

2. サンプルデータ5のsample5.txtの場合:

おまけのところで、「心臓特異的発現パターン」を示す遺伝子群を抽出するための"理想的なパターン(テンプレート)"を含むファイル(GDS1096_cl_heart.txt)を読み込んでいます。

```

in_f <- "sample5.txt" #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

#必要な関数などをロード
source("http://www.iu.a.u-tokyo.ac.jp/~kadota/R/R_functions.R") #ROKUを実行するkadota_2006_bmc_bioinform
library(affy) #Tukey's Biweightを計算するためのtukey.biweight関数が含まれている
library(som) #各遺伝子発現ベクトルを正規化(平均=0, 標準偏差=1)するためのnormal

#入力ファイルの読み込み
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="") #in_fで指定したファイルの読み込み

#本番
data.z <- normalize(data, byrow=TRUE) #正規化を実行し、結果をdata.zに格納
out <- t(apply(data.z, 1, kadota_2003_physiol_genomics_0.25)) #各遺伝子発現ベクトルについて、Ueda's AIC-t
colnames(out) <- colnames(data) #列名を付与
entropy_score <- apply(data, 1, kadota_2006_bmc_bioinformatics) #一行(一遺伝子)ずつ、遺伝子発現ベクトル

#ファイルに保存
tmp <- cbind(rownames(data), out, entropy_score) #左端の列が遺伝子ID、次にサンプル数分だけの列からなる「外れ値行列」
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F) #tmpの中身を指定したファイル名で保存

#以下は(こんなこともできますという)おまけ

```

- What's
- 門田
- や、
- レイ
- お知
- ども
- はじ
- 過去
- RO
- RO

1. TCCパッケージ中のROKU関数を用いるやり方
2. source関数で読み込んで kadota_2006_bmc_bioinformatics関数を用いるやり方

source関数を利用

(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 2014/08/08, since 2010)

What's new?

- このウェブページはフリーソフトRと必要な
- 1. [Rのインストールと起動](#)および2. [基本的](#)
- 2014年10月 [HPCワークショップ「医](#)
[フォマティクス」講習会@仙台国際セン](#)
- 門田幸二 著 [シリーズ Useful R 第7巻](#) [トラ](#)
- [日本乳酸菌学会誌のNGS関連連載の第1](#)
- [参考資料\(講義、講習会、本など\)](#)の項目

- [はじめに](#) (last modified 2014/01/30)
- [参考資料\(講義、講習会、本など\)](#) (last mo
- [過去のお知らせ](#) (last modified 2014/08/03)
- [Rのインストールと起動](#) (last modified 2014
- [基本的な利用法](#) (last modified 2014/07/20)
- [サンプルデータ](#) (last modified 2014/07/17)
- バイオインフォマティクス人材育成カリキュ
- [書籍|トランスクリプトームについて](#) (last m

<http://www.bioconductor.org/biocLite.R>
をsource関数で読み込むことで、ネット
経由で個別のパッケージのインストール
を行うためのbiocLite関数を利用できる。

Rのインストールと起動 NEW

基本的には[こちら](#)または[こちら](#)をご覧ください。

よく分からない人でWindowsユーザーの方は以下を参考にしてください。2014年7月31日にアップデートしたWindows用のインストール手順は[こちら](#)。2014年5月14日にアップデートしたMac版のインストール手順[こちら](#)(by 孫建強氏)もあります。注意点は、「Mac OS Xのバージョンに関わらず R-3.1.0-snowleopard.pkgをインストールしたほうがよい」です。

1. Windows release版のインストールの場合:

- [Rのインストーラ](#)を「実行」
- 聞かれるがままに「次へ」などを押してとにかくインストールを完了させる
- Windows Vistaの人**は(パッケージのインストール中に書き込み権限に関するエラーが出るのを避けるために)「コントロールパネル」-「ユーザーアカウント」-「ユーザーアカウント 制御の有効化または無効化」で、「ユーザーアカウント 制御(UAC)を使ってコンピュータの保護に役立たせる」のチェックをあらかじめ外しておくことを強くお勧めします。
- インストールが無事完了したら、デスクトップに出現する「R3.X.Y(32 bitの場合; XやY中の数値はバージョンによって異なります)」または「R x64 3.X.Y(64 bitの場合)」アイコンをダブルクリックして起動
- 以下を、「R コンソール画面」でコピー&ペーストする。10GB程度のディスク容量を要しますが一番お手軽です。(どこからダウンロードするか?と聞かれるので、その場合は自分のいる場所から近いサイトを指定)

```
install.packages(available.packages()[,1], dependencies=TRUE)#CRAN中にある全てのバック
source("http://www.bioconductor.org/biocLite.R")#おまじない
biocLite(all_group())#Bioconductor中にある全てのパッケージをインストール
biocLite("BSgenome.Athaliana.TAIR.TAIR9", suppressUpdates=TRUE)#Bioconductor中にある
```

- 「コントロールパネル」-「デスクトップのカスタマイズ」-「フォルダオプション」-「表示(タブ)」-「詳細設定」のところで、「登録されている拡張子は表示しない」のチェックを外してください。