

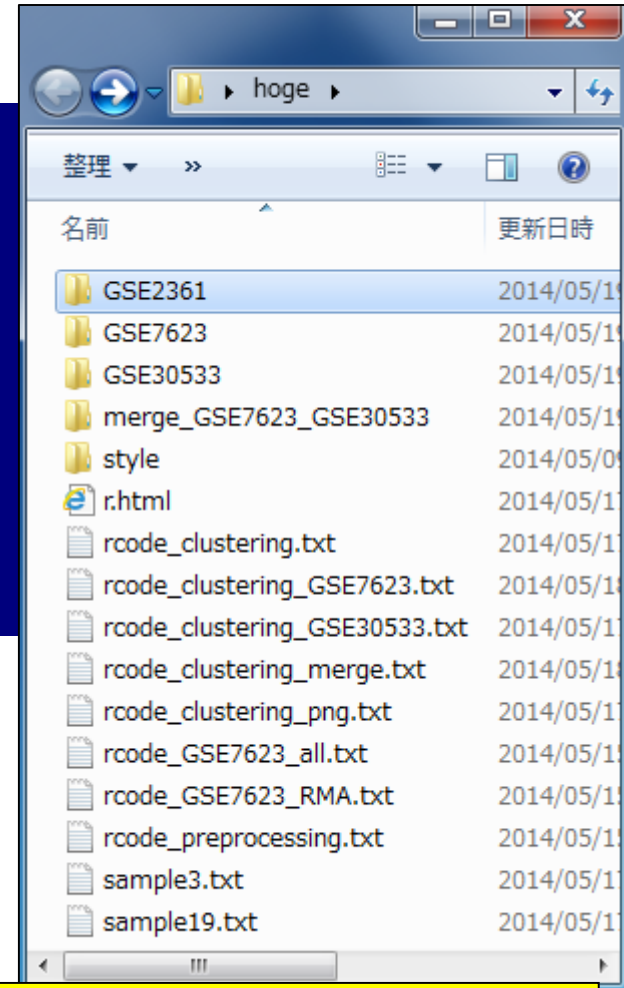
講義室後ろにあるUSBメモリ
中のhogeフォルダをデスクトッ
プにコピーしておいてください。

機能ゲノム学 第2回

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

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



前回 (5/14) 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。

講義予定

- 第1回(2014年5月14日)
 - 原理、各種データベース、生データ取得、遺伝子発現行列作成(データ正規化)
 - 教科書の1.2節、2.2節周辺
- 第2回(2014年5月21日)
 - クラスタリング(データ変換や距離の定義など)、実験デザイン、分布
 - 教科書の3.2節周辺
- 第3回(2014年5月28日)
 - 発現変動解析(多重比較問題)、各種プロット(M-A plotや平均-分散プロット)
 - 教科書の3.2節と4.2節周辺
- 第4回(2014年6月4日)
 - 機能解析(Gene Ontology解析やパスウェイ解析)、分類など

授業の目標・概要

細胞中で発現している全転写物(トランスクリプトーム)の解析技術は、マイクロアレイから次世代シーケンサ(RNA-seq)に移行しつつあります。RNA-seqデータ解析の多くは、マイクロアレイの知識を前提としています。また、ニュートリゲノミクス(食品系)分野では、マイクロアレイは現在でも主流派です。マイクロアレイデータを主な例として、各種トランスクリプトーム解析手法について解説します。



Contents (第2回)

- プロブレベルデータ取得および前処理法適用(発展形)
 - Rパッケージ経由でダウンロードしたデータをそのまま前処理法にかける
 - サンプル名の順番が変わることもあるので注意
 - 何度も同じ作業をしなくていいように有効活用
 - 3つの論文の遺伝子発現行列をコピーで作成
 - ヒトの36サンプルからなるデータ (Ge et al., 2005)
 - ラットの24サンプルからなるデータ (Nakai et al., 2008)
 - ラットの10サンプルからなるデータ (Kamei et al., 2013)
- サンプル間クラスタリング
 - 階層的クラスタリング、距離の定義、クラスターをまとめる方法など
 - 前処理法 (MAS5, RMA, RobLoxBioC) と相関係数 (Spearman and Pearson) の違い
 - 同一アレイ由来データセットのマージ (ラット24サンプル+ラット10サンプル)
 - 3' 発現アレイの長所 (教科書p7)
 - ラット10サンプル (通常 対 鉄欠乏) クラスタリング結果の印象は、外群 (ラット24サンプル) の有無でずいぶん異なる (教科書p106-107)
- 実験デザイン

CELファイル取得 → 前処理法の実行

一気にやる方法もあります

- ① CELファイル取得
- ② 前処理法の実行

- サンプルデータ (last modified 2013/11/25)
- 書籍 | [CELについて](#) (last modified 2014/04/17) **NEW**
- 書籍 | [トランスクリプトーム解析 | 1.1 はじめに](#) (last modified 2014/05/09) **NEW**
- 書籍 | [トランスクリプトーム解析 | 1.2.1 原理\(Affymetrix 3'発現アレイ\)](#) (last modified 2014/05/12) **NEW**
- 書籍 | [トランスクリプトーム解析 | 1.2.2 最近の知見](#) (last modified 2014/05/09) **NEW**
- 書籍 | [トランスクリプトーム解析 | 2.2.1 生データ\(プローブレベルデータ\)取得](#) (last modified 2014/04/18) **NEW**
- 書籍 | [トランスクリプトーム解析 | 2.2.2 データの正規化\(基礎\)](#) (last modified 2014/04/17) **NEW**
- 書籍 | [トランスクリプトーム解析 | 2.2.3 データの正規化\(計算例\)](#) (last modified 2014/04/18) **NEW**
- 書籍 | [トランスクリプトーム解析 | 2.2.4 データの正規化\(その他\)](#) (last modified 2014/04/18) **NEW**
- 書籍 | [トランスクリプトーム解析 | 2.2.5 アンテーション情報](#) (last modified 2014/04/18) **NEW**
- 書籍 | [トランスクリプトーム解析 | 3.2.1 クラスターリング\(データ変換や距離の定義など\)](#) (last modified 2014/04/18) **NEW**
- 書籍 | [トランスクリプトーム解析 | 3.2.2 実験デザイン, データ分布, 統計解析との関係](#) (last modified 2014/04/18) **NEW**
- 書籍 | [トランスクリプトーム解析 | 3.2.3 多重比較問題](#) (last modified 2014/04/19) **NEW**
- 書籍 | [トランスクリプトーム解析 | 3.2.4 各種プロット \(M-A plotや平均-分散プロットなど\)](#) (last modified 2014/04/19) **NEW**
- 書籍 | [トランスクリプトーム解析 | 4.2.1 2群間比較](#) (last modified 2014/04/19) **NEW**
- 書籍 | [トランスクリプトーム解析 | 4.2.2 他の実験デザイン \(paired, multi-factor, 3群間\)](#) (last modified 2014/04/19) **NEW**
- 書籍 | [トランスクリプトーム解析 | 4.2.3 多群間比較\(特異的発現パターン\)](#) (last modified 2014/04/19) **NEW**
- インタロ | [発現データ取得 | 公共DBから](#) (last modified 2014/05/11) **NEW**
- インタロ | [発現データ取得 | inSilicoDb\(Taminau 2011\)](#) (last modified 2014/05/11) **NEW**
- インタロ | [発現データ取得 | ArrayExpress\(Kauffmann 2005\)](#) (last modified 2014/05/11) **NEW**
- インタロ | [発現データ取得 | GEOQuery\(Davis 2007\)](#) (last modified 2014/05/11) **NEW**
- インタロ | [アンテーション情報取得 | 公共DB\(GEO\)から](#) (last modified 2014/05/11) **NEW**
- インタロ | [アンテーション情報取得 | GEOQuery\(Davis 2007\)](#) (last modified 2014/05/11) **NEW**
- インタロ | [アンテーション情報取得 | Rのパッケージ*.dbから](#) (last modified 2014/05/11) **NEW**

- 正規化 | [Affymetrix GeneChip | CELについて](#) (last modified 2013/09/02)
- 正規化 | [Affymetrix GeneChip | frma \(McCall 2010\)](#) (last modified 2013/08/21)
- 正規化 | [Affymetrix GeneChip | rmx \(Kohl 2010\)](#) (last modified 2013/11/19) 推奨
- 正規化 | [Affymetrix GeneChip | GRSN \(Pelz 2008\)](#) (last modified 2013/05/27)
- 正規化 | [Affymetrix GeneChip | Hook \(Binder 2008\)](#) (last modified 2013/05/30)
- 正規化 | [Affymetrix GeneChip | DFW \(Chen 2007\)](#) (last modified 2013/08/20)
- 正規化 | [Affymetrix GeneChip | FARMS \(Hochreiter 2006\)](#) (last modified 2013/08/20)
- 正規化 | [Affymetrix GeneChip | multi-mgMOS \(Liu 2005\)](#) (last modified 2013/08/21)
- 正規化 | [Affymetrix GeneChip | GCRMA \(Wu 2004\)](#) (last modified 2013/08/21)
- 正規化 | [Affymetrix GeneChip | PLIER \(Affymetrix 2004\)](#) (last modified 2013/08/21)
- 正規化 | [Affymetrix GeneChip | VSN \(Huber 2002\)](#) (last modified 2013/08/21)
- 正規化 | [Affymetrix GeneChip | RMA \(Irizarry 2003\)](#) (last modified 2013/08/21)
- 正規化 | [Affymetrix GeneChip | MAS5.0 \(Hubbell 2002\)](#) (last modified 2013/11/25)
- 正規化 | [Affymetrix GeneChip | MBEI \(Li 2001\)](#) (last modified 2013/08/21)

CELファイル取得 → 前処理法の実行

- 書籍 | トランスクリプトーム解析 | [4.2.2 他の実験デザイン\(paired, multi-factor, 3群間\)](#) (last modified 2014/04/19)
- 書籍 | トランスクリプトーム解析 | [4.2.3 多群間比較\(特異的発現パターン\)](#) (last modified 2014/04/20) **NEW**
- イントロ | 発現データ取得 | [公共DBから](#) (last modified 2014/05/11) **NEW**
- イントロ | 発現データ取得 | [inSilicoDb\(Taminau 2011\)](#) (last modified 2013/08/20)
- イントロ | 発現データ取得 | [ArrayExpress\(Kauffmann 2009\)](#) (last modified 2013/08/29) 推奨
- イントロ | 発現データ取得 | [GEOQuery\(Davis 2007\)](#) (last modified 2013/08/20)
- イントロ | アンテーション情報取得 | [公共DB\(GEO\)から](#) (last modified 2013/08/18)
- イントロ | アンテーション情報取得 | **イントロ | 発現データ取得 | [ArrayExpress\(Kauffmann_2009\)](#) **NEW****
- イントロ | アンテーション情報取得
- イントロ | プローブ配列情報取得
- イントロ | トランスクリプトーム配列
- イントロ | トランスクリプトーム配列
- イントロ | [Affymetrix CELファイル](#)

マイクロアレイデータベース [ArrayExpress](#) に登録されているデータを [ArrayExpress](#) というRパッケージで取得するやり方を示します。GEO IDでも検索可能であり、CELファイルデータも取得可能、任意の preprocessing法を適用可能、などの利点からこのパッケージ経由での利用をお勧めします。

「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. AffymetrixデータE-MEXP-1422 (Bourgon et al., PNAS, 2010)のCELファイルを取得し、RMA法(Irizarry et al., Biostatistics, 2003)を実行して得られた発現情報を取得したい場合:

以下の [ArrayExpress](#)関数のオプションを `save=F`から`save=T`に変更すると、CELファイルなどを含む全データのダウンロードも同時に行ってくれます。が、そんなことをいちいちやらなくても `readAffy`関数を用いて読み込んだ状態と同じなので直接RMA (Irizarry et al., Biostatistics, 2003)などの任意の正規化法を適用可能です。

```

out_f <- "data_rma.txt" #出力ファイル名を指定してout_fに格納
param <- "E-MEXP-1422" #入手したいIDを指定

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

#前処理(データ取得)
hoge <- ArrayExpress(param, save=F) #paramで指定したIDのデータを取得した結果をhogeに格納

#本番
eset <- rma(hoge) #RMAを実行し、結果をesetに保存

#ファイルに保存
write.exprs(eset, file=out_f) #結果を指定したファイル名で保存
    
```

これが一気にやるための
テンプレートスクリプトです

R経由で生データ取得 (教科書の § 2.2.1)

■ Affymetrix GeneChip

- Ge et al., *Genomics*, **86**: 127–141, 2005
 - GSE2361、ヒト36サンプル、GPL96を利用
- Nakai et al., *Biosci Biotechnol Biochem.*, **72**: 139–148, 2008
 - GSE7623、ラット24サンプル、GPL1355を利用
- Kamei et al., *PLoS One*, **8**: e65732, 2013
 - GSE30533、ラット10サンプル、GPL1355を利用

GSE7623のプローブレベル
データ取得からRMA前処理
法の実行までを一気にやる

■ Illumina BeadChip

- Sharma et al., *Cancer Cell*, **23**: 35–47, 2013
 - GSE28680、ヒト24サンプル、GPL10558を利用

■ NGSデータも…

- Neyret-Kahn et al., *Genome Res.*, **23**: 1563–1579, 2013
 - GSE42213、ヒト26サンプル、GPL10999とGPL11154を利用
 - GSE42211、ヒト20サンプル、GPL10999とGPL11154を利用 (ChIP-seq)
 - GSE42212、ヒト6サンプル、GPL10999を利用 (RNA-seq)
- Huang et al., *Development*, **139**: 2161–2169, 2012
 - GSE36469、シロイヌナズナ8サンプル、GPL13222を利用

マイクロアレイデータベース [ArrayExpress](#) に登録されているデータを [ArrayExpress](#) というRパッケージで取得するやり方を示します。
 GEO IDでも検索可能であり、CELファイルデータも取得可能、任意の preprocessing法を適用可能、などの利点からこのパッケージ
 経由での利用をお勧めします。
 「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. Affymetrixデータ [E-MEXP-1422 \(Bourgon et al., PNAS, 2010\)](#) のCELファイルを取得し、RMA法 ([Irizarry et al., Biostatistics, 2003](#)) を実行して得られた発現情報を取得したい場合:

以下の [ArrayExpress](#) 関数のオプションを `save=F` から `save=T` に変更すると、CELファイルなどを含む全データのダウンロードも同時に行ってくれます。が、そんなことをいちいちやらなくても `ReadAffy` 関数を用いて読み込んだ状態と同じなので直接RMA ([Irizarry et al., Biostatistics, 2003](#)) などの任意の正規化法を適用可能です。

```
out_f <- "data_rma.txt"           #出力ファイル名を指定してout_fに格納
param <- "E-MEXP-1422"          #入手したいIDを指定

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

```
#前処理(データ取得)
hoge <- ArrayExpress(param, save=T)

#####
### GSE7623のCELファイルダウンロードとRMA前処理実行 ###
#####
out_f <- "data_rma.txt"         #出力ファイル名を指定してout_fに格納↓
param <- "GSE7623"             #入手したいIDを指定↓

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

↓
#前処理(データ取得)↓
hoge <- ArrayExpress(param, save=F) #paramで指定したIDのデータを取得した結果をhogeに格納↓

↓
#本番↓
eset <- rma(hoge)                #RMAを実行し、結果をesetに保存↓

↓
#ファイルに保存↓
write.exprs(eset, file=out_f)   #結果を指定したファイル名で保存↓

←
```

必要最小限の変更後

```
> getwd()
[1] "C:/Users/kadota/Desktop/hoge"
> #####
> ### GSE7623のCELファイルダウンロードとRMA前処理実行 ###
> #####
> out_f <- "data_rma.txt" #出力ファイル名を指定してout_fに格納
> param <- "GSE7623" #入手したいIDを指定
>
> #必要なライブラリをロード
> library(AffyBatch)
ArrayExpress: Reading pheno data from SDRF
ArrayExpress: Reading data files
Read 49 items
GSE7623 was successfully loaded into AffyBatch
警告メッセージ:
1: In strsplit(g, "\t") : 入力文字列 1 はこのロケールでは不適切です
2: In strsplit(g, "\t") : 入力文字列 1 はこのロケールでは不適切です
>
> #本番実行
> eset <- rma(ho
Background corre
Normalizing
Calculating Expr
警告メッセージ:
1: replacing pre
2: replacing pre
>
> #ファイルに保存
> write.expr(es
> |
```

```
R Console
ArrayExpress: Reading pheno data from SDRF
ArrayExpress: Reading data files
Read 49 items
GSE7623 was successfully loaded into AffyBatch
警告メッセージ:
1: In strsplit(g, "\t") : 入力文字列 1 はこのロケールでは不適切です
2: In strsplit(g, "\t") : 入力文字列 1 はこのロケールでは不適切です
>
> #本番実行
> eset <- rma(ho
```

```
Background corre
Normalizing
Calculating Expr
警告メッセージ:
1: replacing pre
2: replacing pre
>
> #ファイルに保存
> write.expr(es
> |
```

data_rma.txt - Excel

ファイル ホーム 挿入 ページレイアウト 数式 データ 校閲 表示 アドイン

L16 : X ✓ fx 9.81808731077457

	A	B	C	D	E	F	G	H	I
1		GSM184434	GSM184414	GSM184418	GSM184427	GSM184437	GSM184420	GSM184421	GSM184422
2	1367452_at	9.65660613	10.524571	10.195674	10.3803459	9.45655687	9.74547642	10.0793083	10.3261
3	1367453_at	9.56071583	9.66424171	9.65657585	9.61752258	9.24437968	9.86596543	9.9082944	9.72840
4	1367454_at	9.60223688	9.65369911	9.79645967	9.61747771	9.85195553	9.727593	9.78282946	9.47288
5	1367455_at	11.4411168	10.7645706	10.5480622	10.9425437	11.2800063	10.7885641	10.7542919	11.1411
6	1367456_at	11.5355052	11.7127252	11.3259473	11.6187254	11.6465211	11.3108027	11.4757781	11.5824
7	1367457_at	9.07839762	8.96244355	8.76459798	9.26592163	9.09357188	8.38811185	8.45840226	9.11513
8	1367458_at	8.50569408	8.28183868	7.5182244	7.59895175	7.8007724	7.85110536	7.98179758	8.08825
9	1367459_at	11.8201046	11.8071	11.6800452	11.7443567	11.6770477	11.5597179	11.5977639	11.914
10	1367460_at	11.1040311	11.6317936	11.5323152	12.0131555	10.996116	11.6244301	11.6467453	11.8150

準備完了

列の順番は異なっているものの、中の数値は同じ

hogeオブジェクトの中身は実質的に同じ

正規化 | Affymetrix GeneChip | RMA (Irizarry_2003)

Affymetrix chip (GeneChip™)を用いて得られた*.CELファイルを元に、RMA(Irizarry et al., [Biostatistics, 2003](#))アルゴリズムを用いてSummary scoreを算出。

「ファイル」-「ディレクトリの変更」で適切なディレクトリに移動し以下をコピー。

1. (CELファイルがあるディレクトリ上で)手元にあるCELファイルの読み込みから行う場合:

作業ディレクトリ中のCEL
ファイルを読み込んだ後の
hogeオブジェクトの中身

```
out_f <- "hoge1.txt"

#必要なパッケージをロード
library(affy)

#データファイルの読み込み(*.CELファイル)
hoge <- ReadAffy()

#本番
eset <- rma(hoge)

#ファイルに保存
write.exprs(eset, file=out_f)
```

#出力ファイル名を指定してout_fに格納

```
R Console
> library(affy) #パッケージの読み込み
> #データファイルの読み込み(*.CELファイル)
> hoge <- ReadAffy() #*.CELファイルの読み込み
> hoge

AffyBatch object
size of arrays=834x834 features (24 kb)
cdf=Rat230_2 (31099 affyids)
number of samples=24
number of genes=31099
annotation=rat2302
notes=
警告メッセージ:
1: replacing previous import by 'utils::head' when loading 'rat2302cdf'
2: replacing previous import by 'utils::tail' when loading 'rat2302cdf'
> sampleNames(hoge)
 [1] "GSM184414.CEL" "GSM184415.CEL" "GSM184416.CEL" "GSM184417.CEL"
 [5] "GSM184418.CEL" "GSM184419.CEL" "GSM184420.CEL" "GSM184421.CEL"
 [9] "GSM184422.CEL" "GSM184423.CEL" "GSM184424.CEL" "GSM184425.CEL"
[13] "GSM184426.CEL" "GSM184427.CEL" "GSM184428.CEL" "GSM184429.CEL"
[17] "GSM184430.CEL" "GSM184431.CEL" "GSM184432.CEL" "GSM184433.CEL"
[21] "GSM184434.CEL" "GSM184435.CEL" "GSM184436.CEL" "GSM184437.CEL"
> |
```

hogeオブジェクトの中身は実質的に同じ

イントロ | 発現データ取得 | ArrayExpress(Kauffmann_2009) NEW

```

#####
GEO IDでも検索
経由での利用を
「ファイル」-「デ
1. Affymetrixデ
Biostatistics, 200
以下のArrayEx
時に行ってくれ
(Irizarry et al.,
out_f <- "data_rma.txt"
param <- "GSE7623"
#必要なパッケージをロード↓
library(ArrayExpress)
library(affy)
#前処理(データ取得)↓
hoge <- ArrayExpress(param)
#本番↓
library(affy)
eset <- rma(hoge)
#ファイルに保存↓
write.exprs(eset, file=out_f)
hoge <- Arr

```

```

R Console
> eset <- rma(hoge) #RMAを実
Background correcting
Normalizing
Calculating Expression
警告メッセージ:
1: replacing previous import by 'utils::head' when loading 'rat2302cdf'
2: replacing previous import by 'utils::tail' when loading 'rat2302cdf'
> #ファイルに保存
> write.exprs(eset, file=out_f) #結果を指定したファイル名で保存
> hoge
AffyBatch object
size of arrays=834x834 features (44 kb)
cdf=Rat230_2 (31099 affyids)
number of samples=24
number of genes=31099
annotation=rat2302
notes=E-GEOD-7623
      E-GEOD-7623
      NA
      c("unknown_experiment_design_type", "transcription profiling by arra$
> sampleNames(hoge)
[1] "GSM184434" "GSM184414" "GSM184418" "GSM184427" "GSM184437" "GSM184420"
[7] "GSM184421" "GSM184423" "GSM184431" "GSM184422" "GSM184417" "GSM184433"
[13] "GSM184426" "GSM184430" "GSM184424" "GSM184428" "GSM184425" "GSM184419"
[19] "GSM184415" "GSM184429" "GSM184416" "GSM184435" "GSM184436" "GSM184432"
> |

```

ウェブサイトからダウンロード後のhogeオブジェクトの中身

こういうこともある、ということを知っておこう

hogeオブジェクトを有効利用

CELファイルのダウンロードと前処理法実行を切り分けることもできます

```
#####↓
### GSE7623のCELファイルダウンロードとRMA前処理実行 ###↓
#####↓
out_f <- "data_rma.txt" #出力ファイル名を指定してout_fに格納↓
param <- "GSE7623" #入手したいIDを指定↓
↓
#必要なパッケージをロード↓
library(ArrayExpress) #パッケージの読み込み↓
library(affy) #パッケージの読み込み↓
↓
#前処理(データ取得)↓
hoge <- ArrayExpress(param, save=F) #paramで指定したIDのデータを取得した結果をhogeに格納↓
```

```
↓
#本番↓
eset <- rma(hoge)
↓
#ファイルに保存↓
write.exprs(eset, file=out_f)
←
```

```
#####↓
### GSE7623のCELファイルダウンロード ###↓
#####↓
param <- "GSE7623" #入手したいIDを指定↓
library(ArrayExpress) #パッケージの読み込み↓
hoge <- ArrayExpress(param, save=F) #paramで指定したIDのデータを取得した結果をhogeに格納↓
↓
#####↓
### RMA前処理法実行 ###↓
#####↓
out_f <- "data_rma.txt" #出力ファイル名を指定してout_fに格納↓
library(affy) #パッケージの読み込み↓
eset <- rma(hoge) #RMAを実行し、結果をesetに保存↓
write.exprs(eset, file=out_f) #結果を指定したファイル名で保存↓
↓
```

```
#####↓
### GSE7623のCELファイルダウンロード ###↓
#####↓
param <- "GSE7623" #入手したいIDを指定↓
library(ArrayExpress) #パッケージの読み込み↓
hoge <- ArrayExpress(param, save=F) #paramで指定したIDのデータを取得した結果をhogeに格納↓
↓
#####↓
### RMA前処理法実行 ###↓
#####↓
out_f <- "data_rma.txt" #出力ファイル名を指定してout_fに格納↓
library(affy) #パッケージの読み込み↓
eset <- rma(hoge) #RMAを実行し、結果をesetに保存↓
write.exprs(eset, file=out_f) #結果を指定したファイル名で保存↓
↓
#####↓
### MAS5前処理法実行 ###↓
#####↓
out_f <- "data_mas.txt" #出力ファイル名を指定してout_fに格納↓
library(affy) #パッケージの読み込み↓
eset <- mas5(hoge) #MASを実行し、結果をesetに保存↓
exprs(eset)[exprs(eset) < 1] <- 1 #対数変換 (log2) できるようにシグナル強度が1未満のもの
exprs(eset) <- log(exprs(eset), 2) #底を2として対数変換↓
write.exprs(eset, file=out_f) #結果を指定したファイル名で保存↓
↓
#####↓
### RMX (RobLoxBioC)前処理法実行 ###↓
#####↓
out_f <- "data_rob.txt" #出力ファイル名を指定してout_fに格納↓
library(RobLoxBioC) #パッケージの読み込み↓
eset <- robloxbioc(hoge) #rmxを実行し、結果をesetに保存↓
exprs(eset)[exprs(eset) < 1] <- 1 #対数変換 (log2) できるようにシグナル強度が1未満のもの
exprs(eset) <- log(exprs(eset), 2) #底を2として対数変換↓
write.exprs(eset, file=out_f) #結果を指定したファイル名で保存↓
↓
```

rcode_GSE7623_all.txt

CELファイル情報を含むhogeオブジェクトを有効利用して、複数の前処理法を一気に実行

利点: 作業ディレクトリがどこであろうが、任意のGSE IDを指定して前処理法実行結果ファイルを得ることができる

門田のやり方

```
#####↓  
### MAS5 ###↓  
#####↓  
out_f <- "data_mas.txt"  
library(affy)  
hoge <- ReadAffy() ←  
eset <- mas5(hoge) ←  
summary(exprs(eset))  
exprs(eset)[exprs(eset) < 1] <- 1  
summary(exprs(eset))  
exprs(eset) <- log(exprs(eset), 2)  
write.exprs(eset, file=out_f)
```

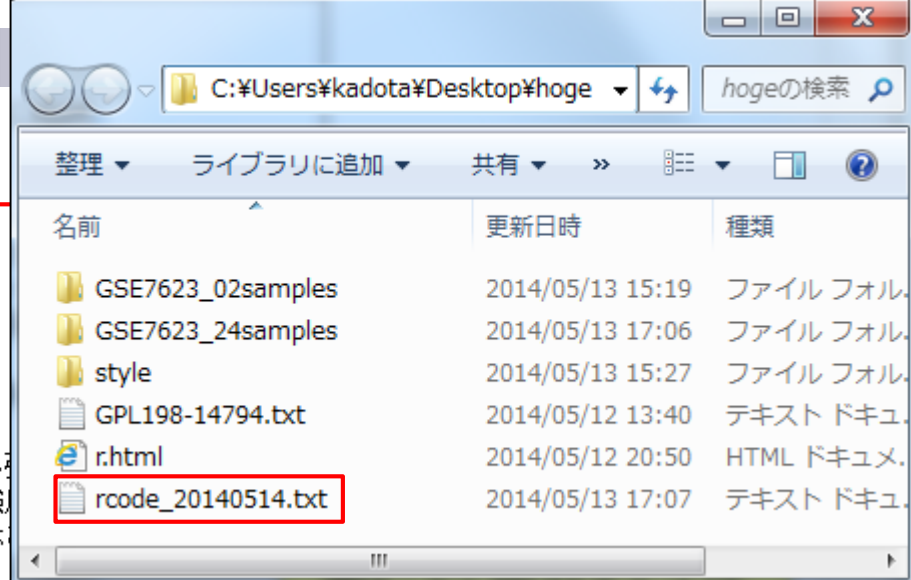
```
#出力ファイル名を指定してout_fに格納↓  
#パッケージの読み込み↓  
#*.CELファイルの読み込み↓  
#MASを実行し、結果をesetに保存↓  
#得られたesetの遺伝子発現行列のシグナル強  
#対数変換 (log2) できるようにシグナル強  
#上記処理後のシグナル強度分布を再び表示  
#底を2として対数変換↓  
#結果を指定したファイル名で保存↓
```

```
#####↓  
### RMA ###↓  
#####↓  
out_f <- "data_rma.txt"  
library(affy)  
hoge <- ReadAffy() ←  
eset <- rma(hoge) ←  
write.exprs(eset, file=out_f)
```

```
#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↓  
#*.CELファイルの読み込み↓  
#RMAを実行し、結果をesetに保存↓  
#結果を指定したファイル名で保存↓
```

```
#####↓  
### RMX (RobLoxBioC) ###↓  
#####↓  
out_f <- "data_rob.txt"  
library(RobLoxBioC)  
hoge <- ReadAffy() ←  
eset <- robloxbioc(hoge) ←  
summary(exprs(eset))  
exprs(eset)[exprs(eset) < 1] <- 1  
summary(exprs(eset))  
exprs(eset) <- log(exprs(eset), 2)  
write.exprs(eset, file=out_f)
```

```
#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↓  
#*.CELファイルの読み込み↓  
#rmxを実行し、結果をesetに保存↓  
#得られたesetの遺伝子発現行列のシグナル強度分布を表  
#対数変換 (log2) できるようにシグナル強度が1未満のも  
#上記処理後のシグナル強度分布を再び表示させて確認↓  
#底を2として対数変換↓  
#結果を指定したファイル名で保存↓
```



rcode_20140514.txt

2014/05/14のスライド75
のコードはただの基本形。
これは毎回作業ディレク
トリ中のCELファイルの読み
込みを行っており非効率

- 正規化 | Affymetrix GeneChip | [RMA \(Irizarry 2003\)](#)
- 正規化 | Affymetrix GeneChip | [MAS5.0 \(Hubbell 2002\)](#)
- 正規化 | Affymetrix GeneChip | [rmx \(Kohl 2010\)](#)

```
#####↓
### 作業ディレクトリ中のCELファイルの読み込み ###↓
#####↓
library(affy) #パッケージの読み込み↓
hoge <- ReadAffy() #*.CELファイルの読み込み↓
↓
↓
#####↓
### RMA前処理法実行 ###↓
#####↓
out_f <- "data_rma.txt" #出力ファイル名を指定してout_fに格納↓
library(affy) #パッケージの読み込み↓
eset <- rma(hoge) #RMAを実行し、結果をesetに保存↓
write.exprs(eset, file=out_f) #結果を指定したファイル名で保存↓
↓
#####↓
### MAS5前処理法実行 ###↓
#####↓
out_f <- "data_mas.txt" #出力ファイル名を指定してout_fに格納↓
library(affy) #パッケージの読み込み↓
eset <- mas5(hoge) #MASを実行し、結果をesetに保存↓
exprs(eset)[exprs(eset) < 1] <- 1 #対数変換 (log2) できるようにシグナル強度が1未満のもの
exprs(eset) <- log(exprs(eset), 2) #底を2として対数変換↓
write.exprs(eset, file=out_f) #結果を指定したファイル名で保存↓
↓
#####↓
### RMX (RobLoxBioC)前処理法実行 ###↓
#####↓
out_f <- "data_rob.txt" #出力ファイル名を指定してout_fに格納↓
library(RobLoxBioC) #パッケージの読み込み↓
eset <- robloxbioc(hoge) #rmxを実行し、結果をesetに保存↓
exprs(eset)[exprs(eset) < 1] <- 1 #対数変換 (log2) できるようにシグナル強度が1未満のもの
exprs(eset) <- log(exprs(eset), 2) #底を2として対数変換↓
write.exprs(eset, file=out_f) #結果を指定したファイル名で保存↓
```

rcode_preprocessing.txt

CELファイル情報を含むhogeオブジェクトを有効利用して、複数の前処理法を一気に実行

利点: CELファイルを含むディレクトリの指定さえすれば、プログラムのどこも変更する必要がない

遺伝子発現行列データは作成済み

■ Affymetrix GeneChip

- Ge et al., *Genomics*, **86**: 127–141, 2005
 - GSE2361、GPL96 (Affymetrix Human Genome U133A Array)、22,283 probesets
 - ヒト36サンプル: Heart (心臓)、Thymus (胸腺)、Spleen (脾臓)、Ovary (卵巣)、Kidney (腎臓)、Skeletal Muscle (骨格筋)、Pancreas (膵臓)、Prostate (前立腺)、…
- Nakai et al., *Biosci Biotechnol Biochem.*, **72**: 139–148, 2008
 - GSE7623、GPL1355 (Affymetrix Rat Genome 230 2.0 Array)、31,099 probesets
 - ラット24サンプル: Brown adipose tissue (褐色脂肪組織; BAT) 8サンプル、White adipose tissue (白色脂肪組織; WAT) 8サンプル、Liver (肝臓; LIV) 8サンプル
 - BAT 8サンプル: 通常 (BAT_fed) 4サンプル 対 24時間絶食 (BAT_fas) 4サンプル
 - WAT 8サンプル: 通常 (WAT_fed) 4サンプル 対 24時間絶食 (WAT_fas) 4サンプル
 - LIV 8サンプル: 通常 (LIV_fed) 4サンプル 対 24時間絶食 (LIV_fas) 4サンプル
- Kamei et al., *PLoS One*, **8**: e65732, 2013
 - GSE30533、GPL1355 (Affymetrix Rat Genome 230 2.0 Array)、31,099 probesets
 - ラット10サンプル: 全てLiver (肝臓) サンプル
 - iron-deficient diet (Iron_def) 5サンプル 対 control diet (Control) 5サンプル

hogeフォルダ中に3つの前処理法の実行結果ファイルがあります。
 MAS5 (data_mas.txt)、RMA (data_rma.txt)、RMX (data_rob.txt)

遺伝子発現行列データは作成済み

Affymetrix GeneChip

□ Ge et al., *Genomics*, **86**: 127-141, 2005

■ GSE2361、GPL96 (Affymetrix Human Genome U133A Array) 22,283 probesets

■ ヒト36サンプル: Heart (臓)、Skeletal Muscle

□ Nakai et al., *Biosci Biotech*

■ GSE7623、GPL1355 (

■ ラット24サンプル: Brown adipose tissue (白色脂肪組織)

□ BAT 8サンプル: 3

□ WAT 8サンプル: 3

□ LIV 8サンプル: 通

□ Kamei et al., *PLoS One*,

■ GSE30533、GPL1355

■ ラット10サンプル: 全て

■ iron-deficient diet (Ird

イントロ | 発現データ取得 | 公共DBから **NEW**

遺伝子発現(主にマイクロアレイ)データベースをリストアップします。

一次データベース

• [GEO: Barrett et al., Nucleic Acids Res., 2013](#)

→ [GSE7623](#)(ラット24サンプル, 62MB): [Nakai et al., BBB, 2008](#)

→ [GSE30533](#)(ラット10サンプル, 25MB): [Kamei et al., PLoS One, 2013](#)

◦ [GSE2361](#)(ヒト36サンプル, 130MB): [Ge et al., Genomics, 2005](#)

◦ [GSE10246](#)(マウス182サンプル, 1.1GB): [Lattin et al., Immunome Res., 2008](#)

◦ [GSE1133](#)(ヒトとマウス438サンプル, 1.7GB): [Su et al., Proc Natl Acad Sci U S A, 2004](#)

◦ [GSE5364](#)(ヒト341サンプル, 生データなし): [Yu et al., PLoS Genet., 2008](#)

◦ [GSE15998](#)(マウス106サンプル, 4.0GB): [原著論文はなし?!エクソアレイ](#)

• [ArrayExpress: Rustici et al., Nucleic Acids Res., 2013](#)

◦ [GSE7623](#)(ラット24サンプル, 62MB): [Nakai et al., BBB, 2008](#)

◦ [GSE30533](#)(ラット10サンプル, 25MB): [Kamei et al., PLoS One, 2013](#)

◦ [GSE2361](#)(ヒト36サンプル, 130MB): [Ge et al., Genomics, 2005](#)

◦ [GSE10246](#)(マウス182サンプル, 1.1GB): [Lattin et al., Immunome Res., 2008](#)

◦ [GSE1133](#)(リンク先なし): [Su et al., Proc Natl Acad Sci U S A, 2004](#)

◦ [GSE5364](#)(ヒト341サンプル, 生データなし): [Yu et al., PLoS Genet., 2008](#)

◦ [GSE15998](#)(マウス106サンプル, 4.0GB): [原著論文はなし?!エクソアレイ](#)

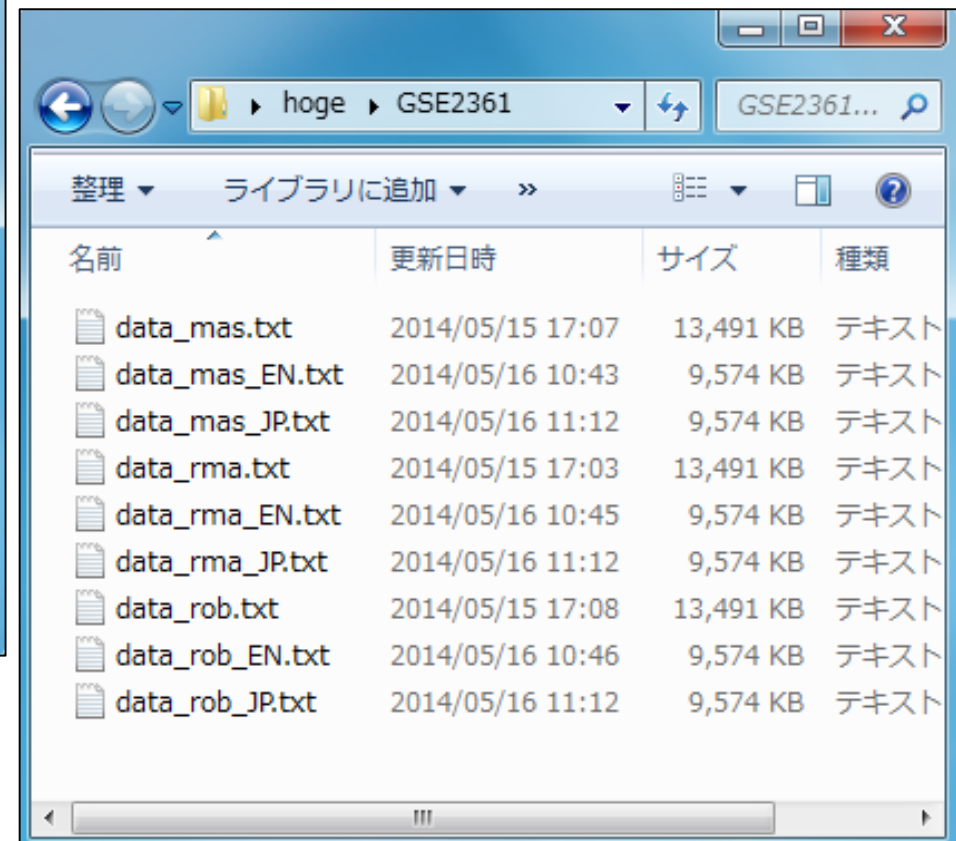
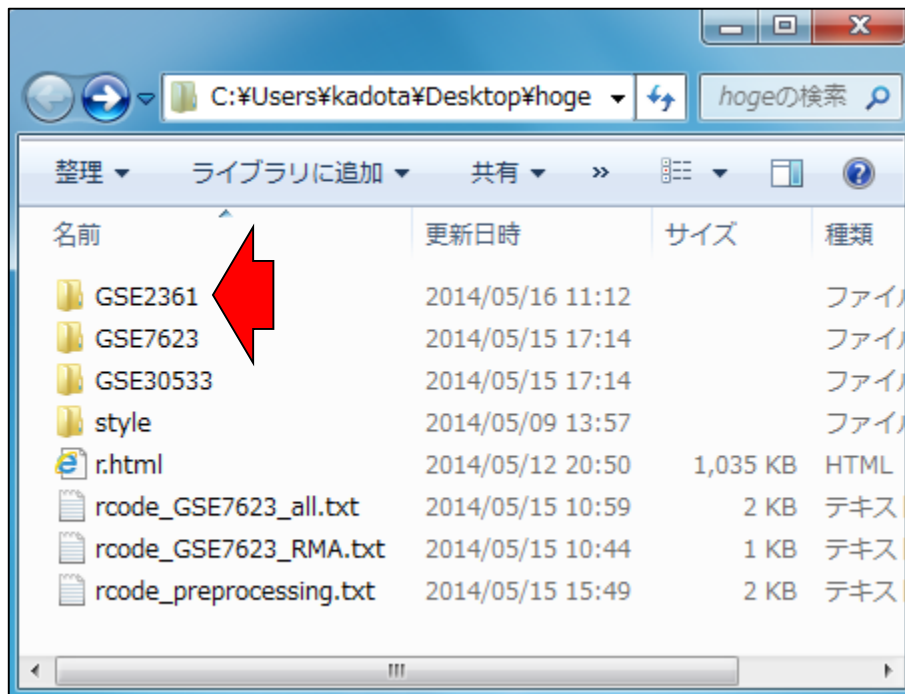
これらのファイルを用いて各種データ解析を行っていきます

遺伝子発現行列データは作成済み

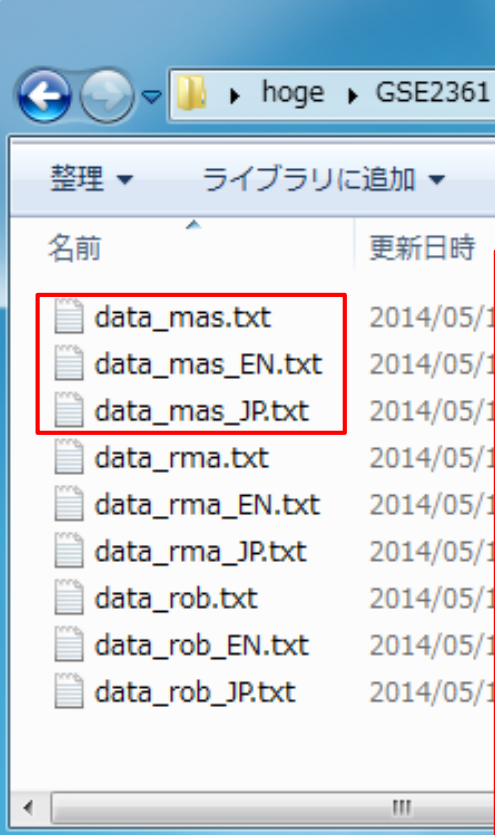
Affymetrix GeneChip

□ Ge et al., *Genomics*, **86**: 127–141, 2005

- GSE2361、GPL96 (Affymetrix Human Genome U133A Array)、22,283 probesets
- ヒト36サンプル: Heart (心臓)、Thymus (胸腺)、Spleen (脾臓)、Ovary (卵巣)、Kidney (腎臓)、Skeletal Muscle (骨格筋)、Pancreas (膵臓)、Prostate (前立腺)、…



*_EN.txtや*_JP.txtのように入力ファイルの段階で(手作業で)解析結果を見やすくするのが一般的。好きなのをご利用ください。いずれも対数変換後のデータです。



data_mas.txt

	A	B	C	D	E	F	G	H	I	J
1		GSM44671	GSM44672	GSM44673	GSM44674	GSM44675	GSM44676	GSM44677	GSM44678	GSM44679
2	1007_s_at	10.82435	11.21261	9.898072	10.8948	11.75531	10.92032	11.61243	11.22101	11.48
3	1053_at	6.621657	6.47488	7.371465	4.168622	7.350244	7.209918	8.362777	7.176571	7.75
4	117_at	8.259603	8.647364	8.689724	7.875649	5.083001	8.165044	7.96707	8.418082	5.719
5	121_at	11.26699	11.04186	11.39999	11.02855	13.13267	11.39138	12.32899	11.0559	11.28
6	1255_g_at	7.1757	6.477278	6.781766	7.048799	7.30767	6.600267	6.42359	6.694412	7.30
7	1294_at	9.137586	9.718507	9.083742	9.014997	8.813377	8.402562	9.146946	8.421351	10.12
8	1316_at	8.235789	7.418995	8.07462	7.280095	8.238176	7.600147	7.422262	7.399894	7.67

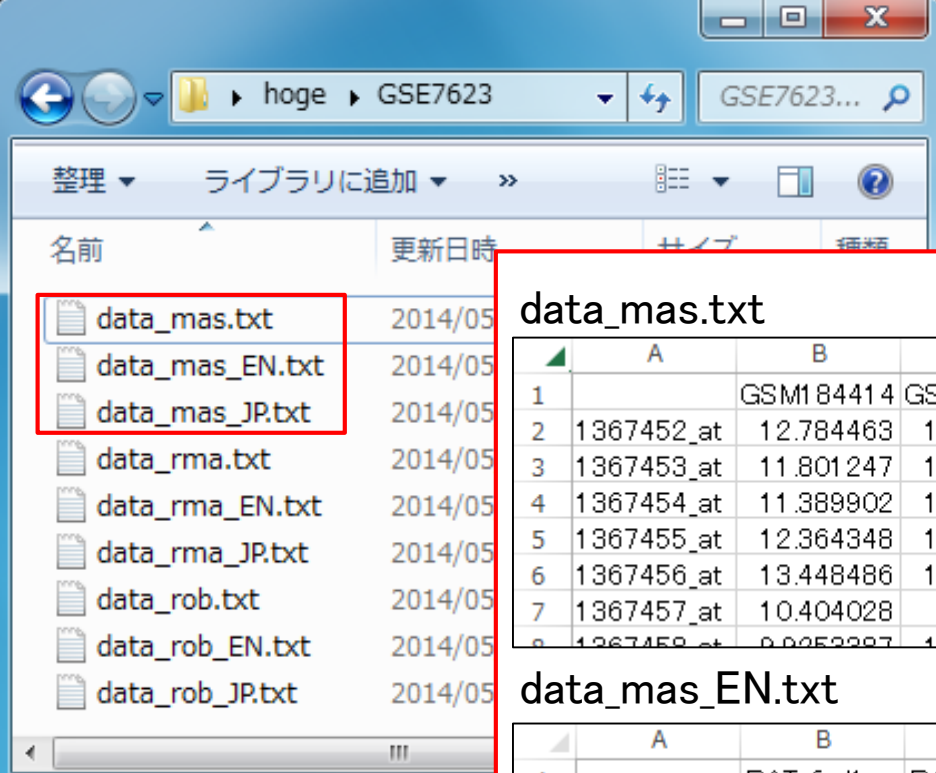
data_mas_EN.txt

	A	B	C	D	E	F	G	H	I	J
1		Heart	Thymus	Spleen	Ovary	Kidney	Skeletal_Mu	Pancreas	Prostate	Small_I
2	1007_s_at	10.82435	11.21261	9.898072	10.8948	11.75531	10.92032	11.61243	11.22101	11.48
3	1053_at	6.621657	6.47488	7.371465	4.168622	7.350244	7.209918	8.362777	7.176571	7.75
4	117_at	8.259603	8.647364	8.689724	7.875649	5.083001	8.165044	7.96707	8.418082	5.719
5	121_at	11.26699	11.04186	11.39999	11.02855	13.13267	11.39138	12.32899	11.0559	11.28
6	1255_g_at	7.1757	6.477278	6.781766	7.048799	7.30767	6.600267	6.42359	6.694412	7.30
7	1294_at	9.137586	9.718507	9.083742	9.014997	8.813377	8.402562	9.146946	8.421351	10.12
8	1316_at	8.235789	7.418995	8.07462	7.280095	8.238176	7.600147	7.422262	7.399894	7.67

data_mas_JP.txt

	A	B	C	D	E	F	G	H	I	J
1		心臓	胸腺	脾臓	卵巣	腎臓	骨格筋	膵臓	前立腺	小腸
2	1007_s_at	10.82435	11.21261	9.898072	10.8948	11.75531	10.92032	11.61243	11.22101	11.48
3	1053_at	6.621657	6.47488	7.371465	4.168622	7.350244	7.209918	8.362777	7.176571	7.75
4	117_at	8.259603	8.647364	8.689724	7.875649	5.083001	8.165044	7.96707	8.418082	5.719
5	121_at	11.26699	11.04186	11.39999	11.02855	13.13267	11.39138	12.32899	11.0559	11.28
6	1255_g_at	7.1757	6.477278	6.781766	7.048799	7.30767	6.600267	6.42359	6.694412	7.30
7	1294_at	9.137586	9.718507	9.083742	9.014997	8.813377	8.402562	9.146946	8.421351	10.12
8	1316_at	8.235789	7.418995	8.07462	7.280095	8.238176	7.600147	7.422262	7.399894	7.67

GSE7623 (Nakai et al., 2008)の対数変換後のデータ



data_mas.txt

	A	B	C	D	E	F	G	H	I
1		GSM184414	GSM184415	GSM184416	GSM184417	GSM184418	GSM184419	GSM184420	GSM184421
2	1367452_at	12.784463	12.447082	12.805908	12.304718	12.589425	12.607532	11.815378	12.439008
3	1367453_at	11.801247	12.152935	11.942227	11.968477	11.845375	11.681727	12.078672	12.048008
4	1367454_at	11.389902	11.160757	11.145987	11.212088	11.540652	11.308877	11.49885	11.402008
5	1367455_at	12.364348	12.529744	12.432574	12.604011	12.441991	12.249935	12.281827	12.190008
6	1367456_at	13.448486	13.543046	13.552794	13.629799	13.36913	13.244278	13.424371	13.329008
7	1367457_at	10.404028	10.69632	10.475078	10.45579	10.141921	10.290666	10.146529	10.260008
8	1367458_at	9.9253387	10.244544	9.9720008	9.9576072	8.702884	9.3578792	9.2134367	9.499008

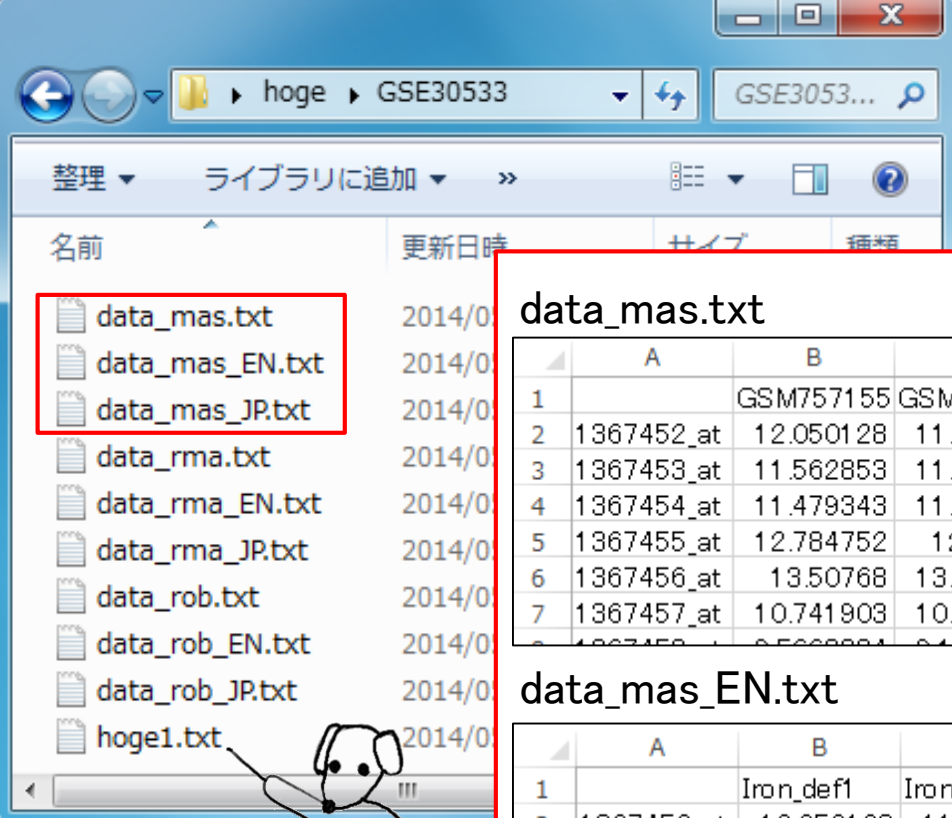
data_mas_EN.txt

	A	B	C	D	E	F	G	H	I
1		BAT_fed1	BAT_fed2	BAT_fed3	BAT_fed4	BAT_fas1	BAT_fas2	BAT_fas3	BAT_fas4
2	1367452_at	12.784463	12.447082	12.805908	12.304718	12.589425	12.607532	11.815378	12.439008
3	1367453_at	11.801247	12.152935	11.942227	11.968477	11.845375	11.681727	12.078672	12.048008
4	1367454_at	11.389902	11.160757	11.145987	11.212088	11.540652	11.308877	11.49885	11.402008
5	1367455_at	12.364348	12.529744	12.432574	12.604011	12.441991	12.249935	12.281827	12.190008
6	1367456_at	13.448486	13.543046	13.552794	13.629799	13.36913	13.244278	13.424371	13.329008
7	1367457_at	10.404028	10.69632	10.475078	10.45579	10.141921	10.290666	10.146529	10.260008
8	1367458_at	9.9253387	10.244544	9.9720008	9.9576072	8.702884	9.3578792	9.2134367	9.499008

data_mas_JP.txt

	A	B	C	D	E	F	G
1		褐色脂肪_満腹1	褐色脂肪_満腹2	褐色脂肪_満腹3	褐色脂肪_満腹4	褐色脂肪_空腹1	褐色脂肪_空腹2
2	1367452_at	12.7844634	12.44708219	12.80590758	12.30471769	12.58942538	12.6075319
3	1367453_at	11.80124704	12.15293493	11.94222741	11.96847729	11.84537542	11.6817274
4	1367454_at	11.38990178	11.16075717	11.14598707	11.21208786	11.54065185	11.3088766
5	1367455_at	12.36434768	12.52974368	12.43257392	12.60401124	12.44199125	12.2499348
6	1367456_at	13.44848649	13.54304603	13.55279359	13.62979898	13.36912977	13.2442783
7	1367457_at	10.40402803	10.69631952	10.47507777	10.4557902	10.14192076	10.2906657
8	1367458_at	9.925338749	10.24454259	9.972000815	9.957607169	8.70288404	9.35787919

GSE30533 (Kamei et al., 2013)の対数変換後のデータ。教科書中で用いているデータセットです。



data_mas.txt

	A	B	C	D	E	F	G	H	I
1		GSM757155	GSM757156	GSM757157	GSM757158	GSM757159	GSM757160	GSM757161	GSM757162
2	1367452_at	12.050128	11.923417	11.991814	11.921645	11.728346	12.079674	12.058299	11.979419
3	1367453_at	11.562853	11.587924	11.62175	11.747779	11.776342	11.633797	11.50868	11.479419
4	1367454_at	11.479343	11.676545	11.608875	11.652704	11.86008	11.70589	11.975747	12.009519
5	1367455_at	12.784752	12.58787	12.696588	12.789029	13.002298	12.678645	12.780148	12.552219
6	1367456_at	13.50768	13.533852	13.483843	13.524296	13.452217	13.472369	13.594191	13.603819
7	1367457_at	10.741903	10.136369	10.609696	10.256467	10.312411	10.49749	10.304808	10.427419

data_mas_EN.txt

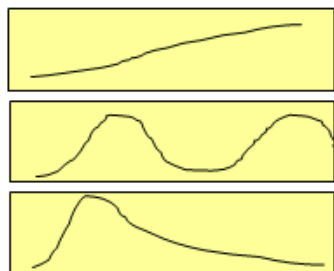
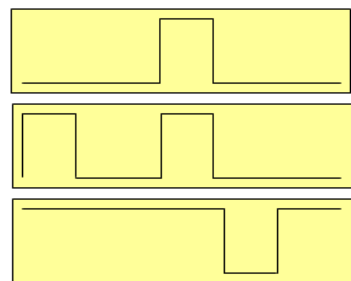
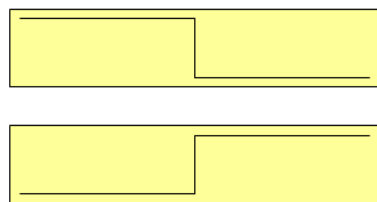
	A	B	C	D	E	F	G	H	I
1		Iron_def1	Iron_def2	Iron_def3	Iron_def4	Iron_def5	Control1	Control2	Control3
2	1367452_at	12.050128	11.923417	11.991814	11.921645	11.728346	12.079674	12.058299	11.979419
3	1367453_at	11.562853	11.587924	11.62175	11.747779	11.776342	11.633797	11.50868	11.479419
4	1367454_at	11.479343	11.676545	11.608875	11.652704	11.86008	11.70589	11.975747	12.009519
5	1367455_at	12.784752	12.58787	12.696588	12.789029	13.002298	12.678645	12.780148	12.552219
6	1367456_at	13.50768	13.533852	13.483843	13.524296	13.452217	13.472369	13.594191	13.603819
7	1367457_at	10.741903	10.136369	10.609696	10.256467	10.312411	10.49749	10.304808	10.427419

data_mas_JP.txt

	A	B	C	D	E	F	G	H	I
1		鉄欠乏1	鉄欠乏2	鉄欠乏3	鉄欠乏4	鉄欠乏5	通常1	通常2	通常3
2	1367452_at	12.050128	11.923417	11.991814	11.921645	11.728346	12.079674	12.058299	11.979419
3	1367453_at	11.562853	11.587924	11.62175	11.747779	11.776342	11.633797	11.50868	11.479419
4	1367454_at	11.479343	11.676545	11.608875	11.652704	11.86008	11.70589	11.975747	12.009519
5	1367455_at	12.784752	12.58787	12.696588	12.789029	13.002298	12.678645	12.780148	12.552219
6	1367456_at	13.50768	13.533852	13.483843	13.524296	13.452217	13.472369	13.594191	13.603819
7	1367457_at	10.741903	10.136369	10.609696	10.256467	10.312411	10.49749	10.304808	10.427419

データ解析もいろいろ

発現変動遺伝子同定



遺伝子発現行列

二群間比較用

	A群		B群	
	A1	A2	B1	B2
gene 1	$x_{1,1}^A$	$x_{1,2}^A$	$x_{1,1}^B$	$x_{1,2}^B$
gene 2	$x_{2,1}^A$	$x_{2,2}^A$	$x_{2,1}^B$	$x_{2,2}^B$
...
gene i	$x_{i,1}^A$	$x_{i,2}^A$	$x_{i,1}^B$	$x_{i,2}^B$
...
gene n	$x_{n,1}^A$	$x_{n,2}^A$	$x_{n,1}^B$	$x_{n,2}^B$

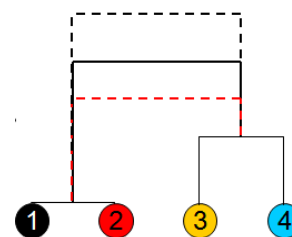
様々な組織(条件)

	S1	S2	S3	S4	...
gene 1	$x_{1,1}$	$x_{1,2}$	$x_{1,3}$	$x_{1,4}$...
gene 2	$x_{2,1}$	$x_{2,2}$	$x_{2,3}$	$x_{2,4}$...
...
gene i	$x_{i,1}$	$x_{i,2}$	$x_{i,3}$	$x_{i,4}$...
...
gene n	$x_{n,1}$	$x_{n,2}$	$x_{n,3}$	$x_{n,4}$...

時系列データ

	T1	T2	T3	T4	...
gene 1	$x_{1,1}$	$x_{1,2}$	$x_{1,3}$	$x_{1,4}$...
gene 2	$x_{2,1}$	$x_{2,2}$	$x_{2,3}$	$x_{2,4}$...
...
gene i	$x_{i,1}$	$x_{i,2}$	$x_{i,3}$	$x_{i,4}$...
...
gene n	$x_{n,1}$	$x_{n,2}$	$x_{n,3}$	$x_{n,4}$...

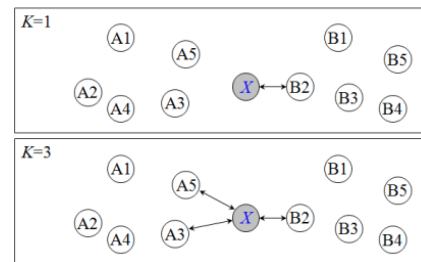
クラスタリング



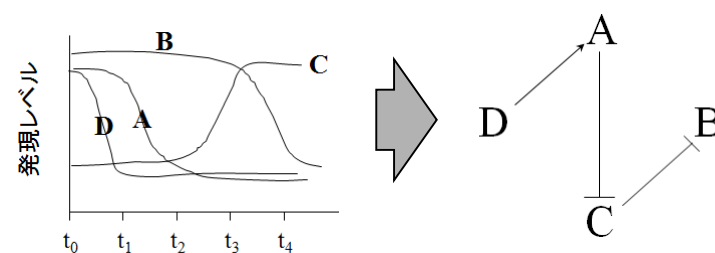
機能解析

- Gene Ontology (GO)
- パスウェイ解析

分類(診断)

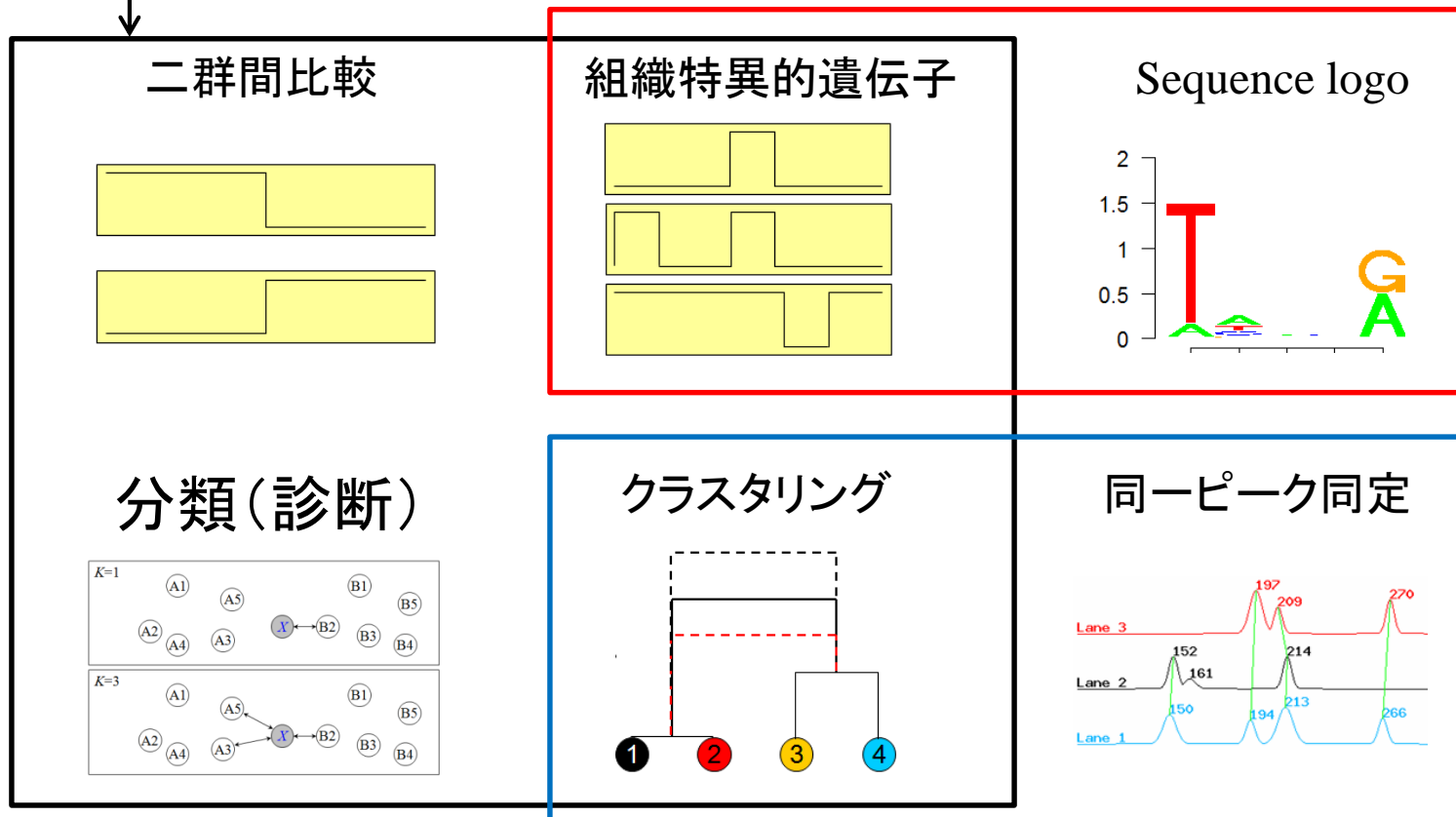


遺伝子ネットワーク推定



バイオインフォマティクス要素技術

■ 相関係数や**エントロピー**などの応用例を紹介



基本スキルのみでいろいろなことができます

第3章 データ解析(基礎)

§ 3.2.1 クラスタリング(データ変換や距離の定義など)

- 書籍 | トランスクリプトーム解析 | [1.1 はじめに](#) (last modified 2014/05/09) NEW
- 書籍 | トランスクリプトーム解析 | [1.2.1 原理\(Affymetrix 3'発現アレイ\)](#) (last modified 2014/05/12) NEW
- 書籍 | トランスクリプトーム解析 | [1.2.2 最近の知見](#) (last modified 2014/05/12) NEW
- 書籍 | トランスクリプトーム解析 | [2.2.1 生データ\(プローブレベルデータ\)取得](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [2.2.2 データの正規化\(基礎\)](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [2.2.3 データの正規化\(計算例\)](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [2.2.4 データの正規化\(その他\)](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [2.2.5 アノテーション情報](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [3.2.1 クラスタリング\(データ変換や距離の定義など\)](#) (last modified 2014/05/16) NEW
- 書籍 | トランスクリプトーム解析 | [3.2.2 実験デザイン, データ分布, 統計解析との関係](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [3.2.3 多重比較問題](#) (last modified 2014/04/19) NEW

教科書中で用いている `hoge1.txt` は、GSE30533 (Kamei et al., 2013) の対数変換 (\log_2 変換) 前の MAS5 データ

書籍 | トランスクリプトーム解析 | 3.2.1 クラスタリング(データ変換や距離の定義など) NEW

シリーズ Useful R 第7巻トランスクリプトーム解析のp99-107のRコードです。

「ファイル」-「ディレクトリの変更」でデスクトップ上の "`E-GEOD-30533.raw.1`" など任意のディレクトリに移動し以下をコピー。

p99の網掛け部分:

書籍中では作業ディレクトリがデスクトップ上の "`E-GEOD-30533.raw.1`" という前提になっていますが、p40で作成した MAS5 データファイル (`hoge1.txt`) を置いてあるディレクトリであればどこでも構いません。

```
in_f <- "hoge1.txt"
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))
data.dist <- as.dist(1 - cor(data, method = "spearman"))
out <- hclust(data.dist, method = "average")
plot(out) # 図3-1作成部分
```

第3章 データ解析(基礎)

- 書籍 | トランスクリプトーム解析 | [1.1 はじめに](#) (last modified 2014/05/09) NEW
- 書籍 | トランスクリプトーム解析 | [1.2.1 原理\(Affymetrix 3'発現アレイ\)](#) (last modified 2014/05/09) NEW
- 書籍 | トランスクリプトーム解析 | [1.2.2 最近の知見](#) (last modified 2014/05/09) NEW
- 書籍 | トランスクリプトーム解析 | [2.2.1 生データ\(プローブレベル\)の取得](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [2.2.2 データの正規化\(基礎\)](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [2.2.3 データの正規化\(計算例\)](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [2.2.4 データの正規化\(その他\)](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [2.2.5 アンテーション情報](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [3.2.1 クラスターリング\(データ変換や距離の定義など\)](#) (last modified 2014/04/18) NEW
- 書籍 | トランスクリプトーム解析 | [3.2.2 実験デザイン, データ分布, 統計解析との関係](#) (last modified 2014/04/18) NEW

教科書中で用いている `hoge1.txt` は、GSE30533 (Kamei et al., 2013) の対数変換 (\log_2 変換) 前の MAS5 データ

p40の網掛け部分(上):

`hoge1.txt` と同じものができていると思います。

```
out_f <- "hoge1.txt"
library(affy)
hoge <- ReadAffy()
eset <- mas5(hoge)
write.exprs(eset, file=out_f)
```

```
#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
#*.CELファイルの読み込み
#MAS5を実行し、結果をesetに保存
#結果をout_fで指定したファイル名で保存
```

	A	B	C	D	E	F	G	H	I
1		GSM757155	GSM757156	GSM757157	GSM757158	GSM757159	GSM757160	GSM757161	GSM757162
2	1367452_at	4240.8	3884.2	4072.8	3879.5	3393.0	4328.6	4264.9	4039
3	1367453_at	3025.3	3078.3	3151.3	3439.0	3507.8	3177.8	2913.8	2855
4	1367454_at	2855.1	3273.3	3123.3	3219.7	3717.4	3340.6	4027.7	4123
5	1367455_at	7056.6	6156.4	6638.3	7077.5	8205.1	6556.2	7034.1	6006
6	1367456_at	11647.1	11860.3	11456.2	11782.0	11207.8	11365.5	12366.9	12449
7	1367457_at	1712.5	1125.5	1562.6	1223.2	1271.6	1445.6	1264.9	1377
8	1367458_at	758.4	575.5	568.6	494.5	681.0	610.6	442.0	565

File Explorer window showing the directory structure for GSE30533. The files listed include:

- data_mas.txt
- data_mas_EN.txt
- data_mas_JP.txt
- data_rma.txt
- data_rma_EN.txt
- data_rma_JP.txt
- data_rob.txt
- data_rob_EN.txt
- data_rob_JP.txt
- hoge1.txt

	A	B	C	D	E	F	G	H	I
1		GSM757155	GSM757156	GSM757157	GSM757158	GSM757159	GSM757160	GSM757161	GSM757162
2	1367452_at	12.050128	11.923417	11.991814	11.921645	11.728346	12.079674	12.058299	11.979412
3	1367453_at	11.562853	11.587924	11.62175	11.747779	11.776242	11.633797	11.50868	11.479412
4	1367454_at	11.479343	11.671924	11.62175	11.747779	11.776242	11.633797	11.50868	11.479412
5	1367455_at	12.784752	12.50128	12.050128	11.923417	11.991814	11.921645	11.728346	12.079674
6	1367456_at	13.50768	13.53128	13.50768	13.53128	13.50768	13.53128	13.50768	13.53128
7	1367457_at	10.741903	10.136369	10.609696	10.256467	10.312411	10.49749	10.304808	10.427412

GSE30533 (Kamei et al., 2013) の対数変換後のMAS5データ

```
#####↓
### CELファイルの読み込みとMAS5前処理法実行 ###↓
#####↓
out_f <- "data_mas.txt" #出力ファイル名を指定してout_fに格納↓
library(affy) #パッケージの読み込み↓
hoge <- ReadAffy() #*.CELファイルの読み込み↓
eset <- mas5(hoge) #MASを実行し、結果をesetに保存↓
exprs(eset)[exprs(eset) < 1] <- 1 #対数変換 (log2) できるようにシグナル強度が1未満のものを1にしておく
exprs(eset) <- log(exprs(eset), 2) #底を2として対数変換↓
write.exprs(eset, file=out_f) #結果を指定したファイル名で保存↓
↓
```

対数変換部分

§ 3.2.1 クラスタリング

- ・書籍 | トランスクリプトーム解析 | [2.2.2 データの正規化\(基礎\)](#) (last modified 2014/04/17) NEW
- ・書籍 | トランスクリプトーム解析 | [2.2.3 データの正規化\(計算例\)](#) (last modified 2014/04/18) NEW
- ・書籍 | トランスクリプトーム解析 | [2.2.4 データの正規化\(その他\)](#) (last modified 2014/04/18) NEW
- ・書籍 | トランスクリプトーム解析 | [2.2.5 アノテーション情報](#) (last modified 2014/04/18) NEW
- ・書籍 | トランスクリプトーム解析 | [3.2.1 クラスタリング\(データ変換や距離の定義など\)](#) (last modified 2014/05/16) NEW
- ・書籍 | トランスクリプトーム解析 | [3.2.2 実験デザイン, データ分布, 統計解析との関係](#) (last modified 2014/04/18) NEW
- ・書籍 | トランスクリプトーム解析 | [3.2.3 多重比較問題](#) (last modified 2014/04/19) NEW

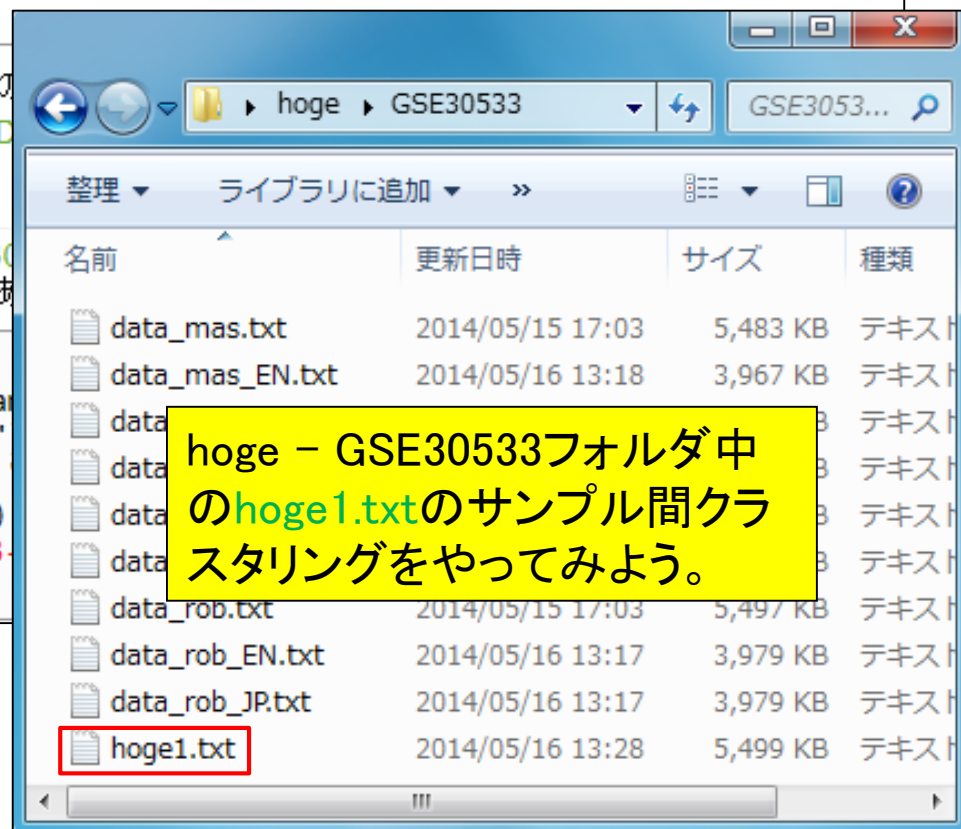
書籍 | トランスクリプトーム解析 | 3.2.1 クラスタリング(データ変換や距離の定義など) NEW

シリーズ Useful R 第7巻 トランスクリプトーム解析のp99-1070
 「ファイル」-「ディレクトリの変更」でデスクトップ上の"E-GEOD

p99の網掛け部分:

書籍中では作業ディレクトリがデスクトップ上の"E-GEOD-30
 MAS5データファイル([hoge1.txt](#))を置いてあるディレクトリであ

```
in_f <- "hoge1.txt"
data <- read.table(in_f, header=TRUE, row.names=1)
colnames(data) <- c(paste("G1_", 1:5, sep=""))
data.dist <- as.dist(1 - cor(data, method = "spearmanr"))
out <- hclust(data.dist, method = "average")
plot(out) # 図3.2.1
```



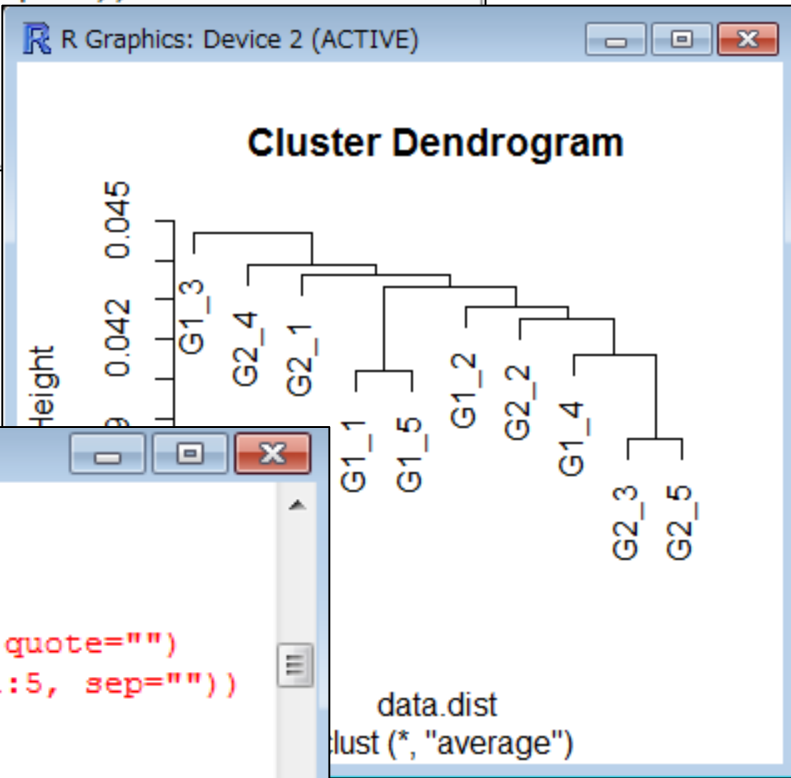
シリーズ Useful R 第
「ファイル」-「ディレク
p99の網掛け部分:
書籍中では作業デ
MAS5データファイル

	A	B	C	D	E	F	G	H	I
1		GSM757155	GSM757156	GSM757157	GSM757158	GSM757159	GSM757160	GSM757161	GSM7571
2	1367452_at	4240.8	3884.2	4072.8	3879.5	3393.0	4328.6	4264.9	4039
3	1367453_at	3025.3	3078.3	3151.3	3439.0	3507.8	3177.8	2913.8	2855
4	1367454_at	2855.1	3273.3	3123.3					
5	1367455_at	7056.6	6156.4	6638.3					
6	1367456_at	11647.1	11860.3	11456.2					
7	1367457_at	1712.5	1125.5	1562.6	1223.2	1271.6	1445.6	1264.9	1377
8	1367458_at	759.1	575.5	568.6	404.5	681.0	610.6	442.0	565

GSE30533 (Kamei et al., 2013)
の対数変換前のMAS5データ

```
in_f <- "hoge1.txt"
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))
data.dist <- as.dist(1 - cor(data, method = "spearman"))
out <- hclust(data.dist, method = "average")
plot(out)
# 図3-1作成部分
```

なぜサンプル名のところが変わっているのかを説明します



```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/GSE30533"
> in_f <- "hoge1.txt"
> data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
> colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))
> data.dist <- as.dist(1 - cor(data, method = "spearman"))
> out <- hclust(data.dist, method = "average")
> plot(out)
# 図3-1作成部分
```

書籍中では作業ディレクトリがデスクトップ上の"E-GEOD-30533.raw.1"という前提になっていますが、p40で作成したMAS5データファイル([hoge1.txt](#))を置いてあるディレクトリであればどこでも構いません。

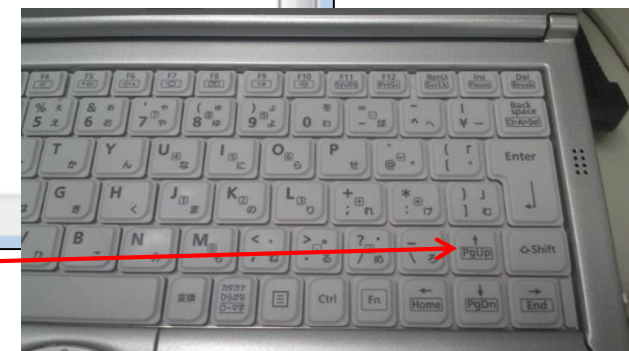
```
in_f <- "hoge1.txt"
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))
data.dist <- as.dist(1 - cor(data, method = "spearman"))
out <- hclust(data.dist, method = "average")
plot(out) # 図3-1作成部分
```

サンプル名部分に相当する
colnames(data)に任意の文字
列ベクトルを代入しています

R Console

```
> data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
> colnames(data)
[1] "GSM757155_.Fe_short_27.CEL" "GSM757156_.Fe_short_31.CEL"
[3] "GSM757157_.Fe_short_33.CEL" "GSM757158_.Fe_short_35.CEL"
[5] "GSM757159_.Fe_short_37.CEL" "GSM757160_control_28.CEL"
[7] "GSM757161_control_30.CEL"   "GSM757162_control_32.CEL"
[9] "GSM757163_control_34.CEL"   "GSM757164_control_36.CEL"
> colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))
> colnames(data)
[1] "G1_1" "G1_2" "G1_3" "G1_4" "G1_5" "G2_1" "G2_2" "G2_3" "G2_4" "G2_5"
> paste("G1_", 1:5, sep="")
[1] "G1_1" "G1_2" "G1_3" "G1_4" "G1_5"
> paste("G1", 1:5, sep="_")
[1] "G1_1" "G1_2" "G1_3" "G1_4" "G1_5"
> paste("Iron_def", 1:5, sep="")
[1] "Iron_def1" "Iron_def2" "Iron_def3" "Iron_def4" "Iron_def5"
> |
```

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。



§ 3.2.1 クラスターリング

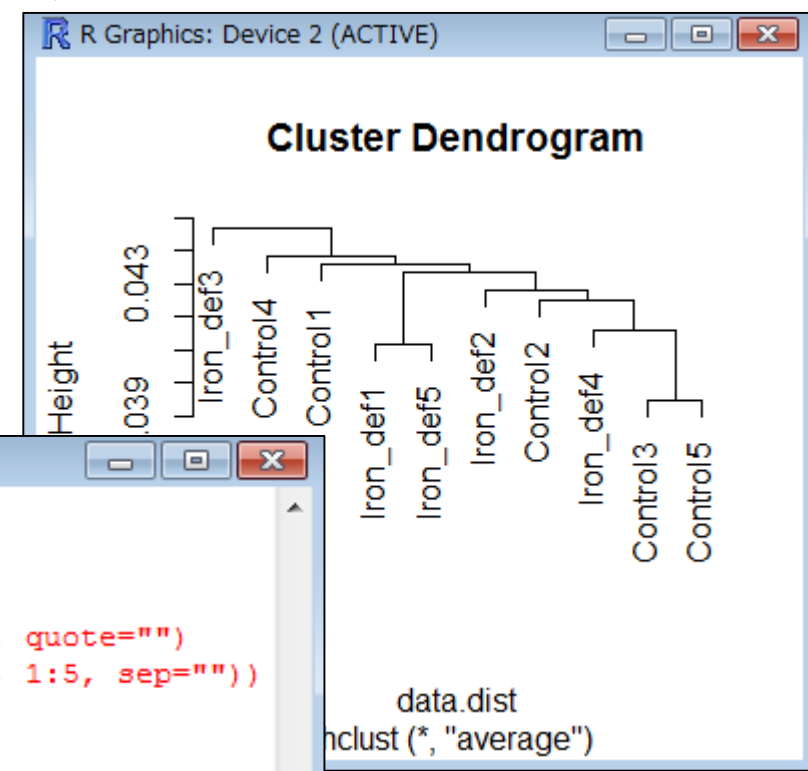
	A	B	C	D	E	F	G	H	I
1		Iron_def1	Iron_def2	Iron_def3	Iron_def4	Iron_def5	Control1	Control2	Control3
2	1367452_at	12.050128	11.923417	11.991814	11.921645	11.728346	12.079674	12.058299	11.979
3	1367453_at	11.562853	11.587924	11.69475	11.747779	11.776249	11.622797	11.50868	11.4794
4	1367454_at	11.479343	11.676545	11.9	11.9	11.9	11.9	11.9	11.9
5	1367455_at	12.784752	12.58787	12.9	12.9	12.9	12.9	12.9	12.9
6	1367456_at	13.50768	13.533852	13.9	13.9	13.9	13.9	13.9	13.9
7	1367457_at	10.741903	10.136369	10.609696	10.256467	10.312411	10.49749	10.304808	10.4274

GSE30533 (Kamei et al., 2013) の対数変換後のMAS5データ

rcode_clustering.txt

```
#####
### MAS5データのクラスターリング ###
#####
in_f <- "data_mas_EN.txt"
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
#colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))
data.dist <- as.dist(1 - cor(data, method = "spearman"))
out <- hclust(data.dist, method = "average")
plot(out)
<
```

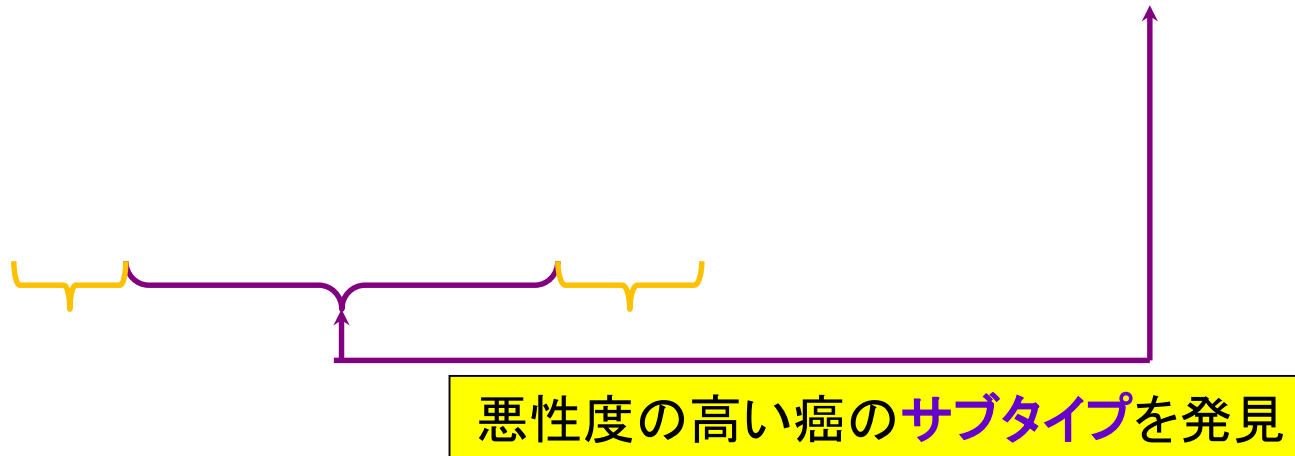
対数変換の有無にかかわらずクラスターリング結果(樹形図)のトポロジーは不変。理由は、Spearman相関係数を採用しているから。



```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/GSE30533"
> in_f <- "data_mas_EN.txt"
> data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
> #colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))
> data.dist <- as.dist(1 - cor(data, method = "spearman"))
> out <- hclust(data.dist, method = "average")
> plot(out)
> |
```

(サンプル間)クラスタリングの実例

- 悪性黒色腫(メラノーマ)31サンプルのデータ



クラスタリング (教師なし学習)

■ 階層的クラスタリング

- 発現パターンの類似した遺伝子を集めて系統樹を作成

■ 非階層的 (分割最適化) クラスタリング

□ K-meansクラスタリング

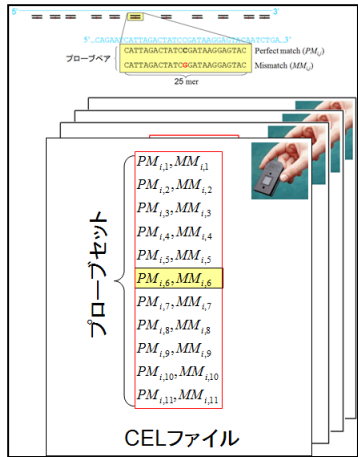
- 「K個のクラスターに分割 (Kの数は主観的に決定) する」と予め指定し、各クラスター内の遺伝子 (サンプル) 間の距離の総和が最小になるようなK個のクラスターを作成

□ 自己組織化マップ (SOM)

□ 主成分分析 (PCA)

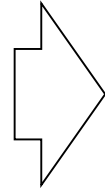
どれを使うかはほぼ趣味の問題です

得られる樹形図の可能性は無数



生データファイル

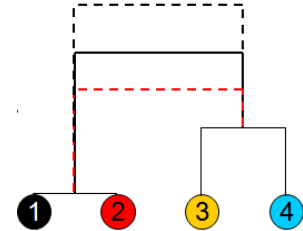
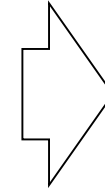
前処理法



	S1	S2	S3	S4	...
gene 1	$x_{1,1}$	$x_{1,2}$	$x_{1,3}$	$x_{1,4}$...
gene 2	$x_{2,1}$	$x_{2,2}$	$x_{2,3}$	$x_{2,4}$...
...
gene i	$x_{i,1}$	$x_{i,2}$	$x_{i,3}$	$x_{i,4}$...
...
gene n	$x_{n,1}$	$x_{n,2}$			

遺伝子発現行列

クラスタリング



<p>①前処理法</p> <ul style="list-style-type: none"> ・MAS ・RMA ・RobLoxBioC 	×	<p>②スケーリング</p> <ul style="list-style-type: none"> ・対数変換 ・相対値(0~1) ・Z-score化 	×	<p>③距離</p> <ul style="list-style-type: none"> ・1-相関係数 ・ユークリッド ・マンハッタン ・キャンベラ ... 	×	<p>④群の併合</p> <ul style="list-style-type: none"> ・単連結法 ・完全連結法 ・平均連結法 ・ワード法 ...
--	---	--	---	---	---	--

クラスタリング (教師なし学習)

■ 決めておくべき二つの基準 (事柄)

□ 距離 (類似度) の定義

- ユークリッド距離、マンハッタン距離など

□ クラスタをまとめる (併合する) 方法

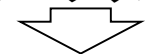
- クラスタ間の距離を定義する方法、とほぼ同じ
- 最短距離法、平均連結法、ワード法など

得られた結果の妥当性を何らかの知見に基づいて評価するため、結果の正当性を主張する視点が複数存在しうる。
→ 見当をつける、程度に利用すべし

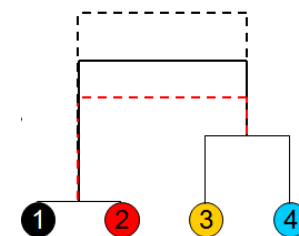
入力例

	x^1	x^2	x^3	x^4
Sample 1	38	42	204	204
Sample 2	0	3	182	182
Sample 3	6	6	19	168
Sample 4	141	106	11	11
Sample 5	8	5	79	179
Sample 6	16	20	96	126
Sample 7	2	5	164	164
Sample 8	132	101	222	0
Sample 9	7	9	164	164
Sample 10	2	8	16	116
Sample 11	66	79	195	195
Sample 12	8	9	241	241
Sample 13	6	8	177	177

クラスタリング



出力例



距離(類似度)の定義

■ ベクトル x と y の発現パターンの距離 $D(x,y)$

$$\text{相関係数 } r_{xy} = \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (-1 \leq r_{xy} \leq 1)$$

i	x	y
1	x_1	y_1
2	x_2	y_2
3	x_3	y_3
4	x_4	y_4
5	x_5	y_5
...
n	x_n	y_n

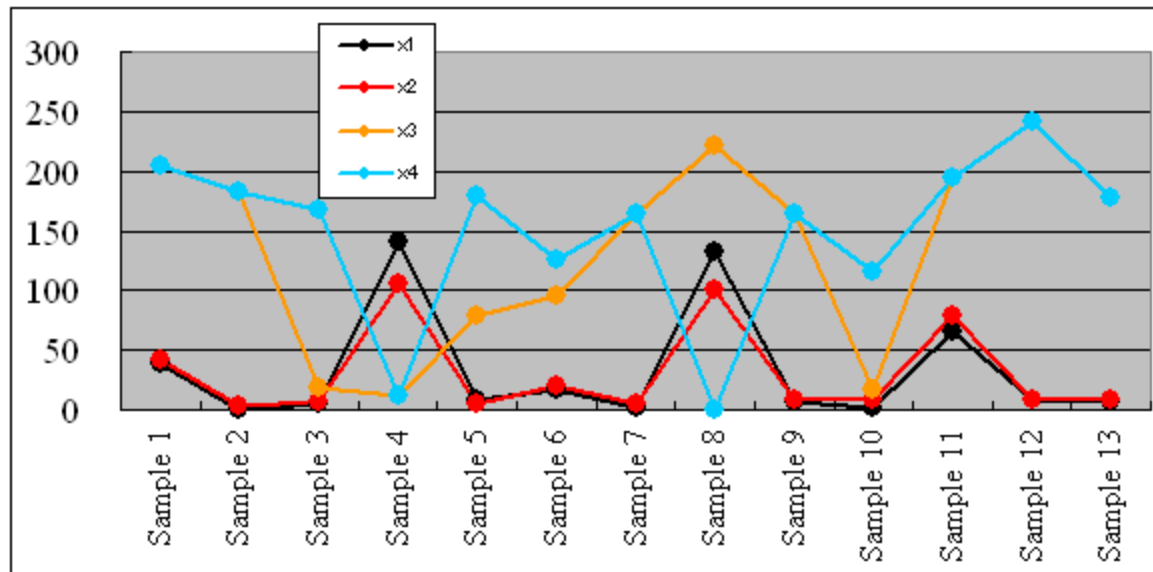
$$\left\{ \begin{array}{l} x \text{ と } y \text{ の発現パターンが酷似} \rightarrow r \approx 1 \\ x \text{ と } y \text{ の発現パターンがばらばら} \rightarrow r \approx 0 \\ x \text{ と } y \text{ の発現パターンがほぼ正反対} \rightarrow r \approx -1 \end{array} \right.$$

$$\text{距離 } D(x,y) = 1 - r \quad (0 \leq D \leq 2) \quad \left\{ \begin{array}{l} r = 1 \rightarrow D = 1 - 1 = 0 \\ r = 0 \rightarrow D = 1 - 0 = 1 \\ r = -1 \rightarrow D = 1 - (-1) = 2 \end{array} \right.$$

「1－相関係数」を距離として定義することができます

相関係数 → 距離(計算例)

	x^1	x^2	x^3	x^4
Sample 1	38	42	204	204
Sample 2	0	3	182	182
Sample 3	6	6	19	168
Sample 4	141	106	11	11
Sample 5	8	5	79	179
Sample 6	16	20	96	126
Sample 7	2	5	164	164
Sample 8	132	101	222	0
Sample 9	7	9	164	164
Sample 10	2	8	16	116
Sample 11	66	79	195	195
Sample 12	8	9	241	241
Sample 13	6	8	177	177



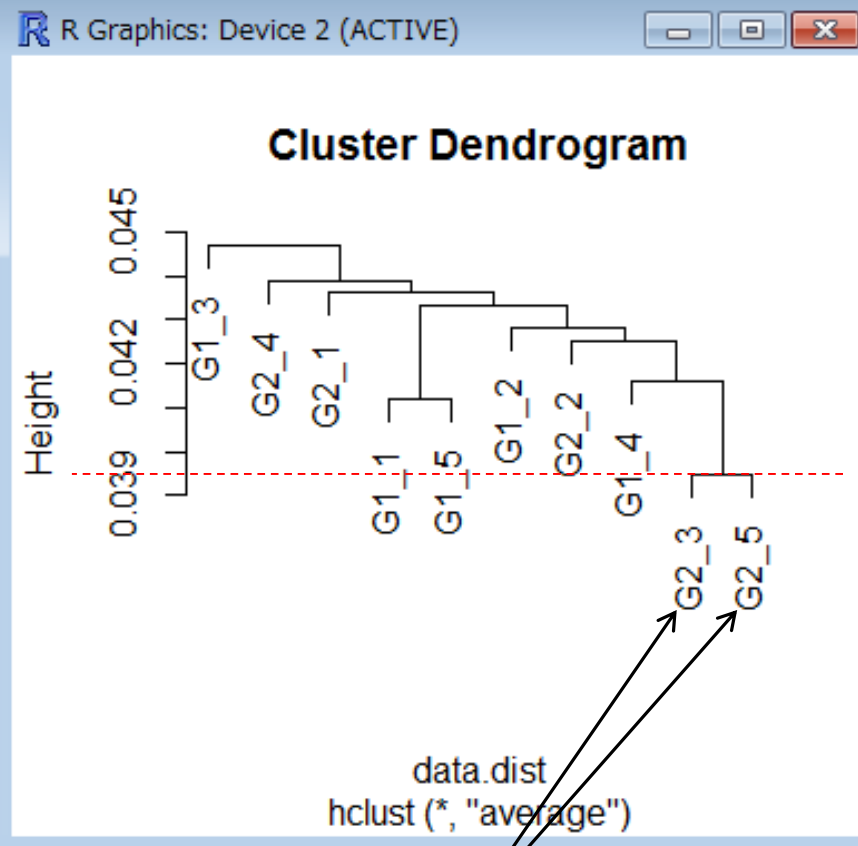
相関係数 $r_{x^1x^2} = 0.98 \rightarrow$ 距離 $D_{x^1x^2} = 1 - 0.98 = 0.02$

相関係数 $r_{x^1x^3} = -0.01 \rightarrow$ 距離 $D_{x^1x^3} = 1 - (-0.01) = 1.01$

相関係数 $r_{x^1x^4} = -0.78 \rightarrow$ 距離 $D_{x^1x^4} = 1 - (-0.78) = 1.78$

書籍中では作業ディレクトリがデスクトップ上の"E-GEOD-30533.raw.1"という前提になっていますが、p40で作成したMAS5データファイル([hoge1.txt](#))を置いてあるディレクトリであればどこでも構いません

```
in_f <- "hoge1.txt"
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t",
  colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))
data.dist <- as.dist(1 - cor(data, method = "spearman"))
out <- hclust(data.dist, method = "average")
plot(out) # 図3-1作成部分
```



G2_3とG2_5の発現ベクトル間の距離

R Console

```
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/GSE30533"
> in_f <- "hoge1.txt"
> data <- read.table(in_f, header=TRUE, row.names=1, sep="\t",
  colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))
> data.dist <- as.dist(1 - cor(data, method = "spearman"))
> out <- hclust(data.dist, method = "average")
> plot(out) # 図3-1作成部分
> data.dist
```

	G1_1	G1_2	G1_3	G1_4	G1_5	G2_1	G2_2	G2_3	G2_4	G2_5
G1_2	0.04533075									
G1_3	0.04530167	0.04403145								
G1_4	0.04298677	0.04437527	0.04491689							
G1_5	0.04119100	0.04173156	0.04377032	0.04466745						
G2_1	0.04336350	0.04423320	0.04547613	0.04395383	0.04331610					
G2_2	0.04550683	0.04330957	0.04650139	0.04415664	0.04429918	0.04386371				
G2_3	0.04243402	0.04060421	0.04257708	0.04034358	0.04158203	0.04243729	0.04110153			
G2_4	0.04462469	0.04483793	0.04526199	0.04423116	0.04244358	0.04430399	0.04405626	0.04210250		
G2_5	0.04244397	0.04304633	0.04445892	0.04287979	0.04232089	0.04416588	0.04230360	0.03949559	0.04432153	

相関係数

```

R Console
> cor(data[,8], data[,10], method="spearman")
[1] 0.9605044
> 1 - cor(data[,8], data[,10], method="spearman")
[1] 0.03949559
>
> cor(data[,8], data[,10], method="pearson")
[1] 0.9928171
> cor(rank(data[,8]), rank(data[,10]), method="pearson")
[1] 0.9605044
>
> cor(data[, "G2_3"], data[, "G2_5"], method="spearman")
[1] 0.9605044
> cor(rank(data[,8]), rank(data[,10]), method="spearman")
[1] 0.9605044
> |
    
```

対数変換前の数値で Spearman相関係数

対数変換前の1 - Spearman相関係数

対数変換前の Pearson相関係数

対数変換前の順位で Pearson相関係数

対数変換前の順位で Spearman相関係数

Spearman相関係数を用いれば、対数変換の有無に関わらず、距離の値が変わらないようにすることもできる。しかし、ユークリッド距離などそれ以外の多くの場合には対数変換の有無によって値が変わるため、マイクロアレイデータは対数変換後の値で取り扱うのが一般的である。p100-106。

他の距離(類似度)を定義する手段

■ ベクトル x と y の発現パターンの距離 $D(x,y)$

□ ユークリッド距離 $D = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$

□ マンハッタン距離 $D = \sum_{i=1}^n |x_i - y_i|$

□ 最大距離 $D = \max(|x_1 - y_1|, \dots, |x_i - y_i|, \dots, |x_n - y_n|)$

□ キャンベラ距離 $D = \sum_{i=1}^n \frac{|x_i - y_i|}{|x_i + y_i|}$

□ ...

i	x	y
1	x_1	y_1
2	x_2	y_2
3	x_3	y_3
4	x_4	y_4
5	x_5	y_5
...
n	x_n	y_n

計算例 (サンプルxとy間の距離D)

	A	B	C	D	E	F
1		<i>x</i>	<i>y</i>		$ x_i - y_i $	$ x_i - y_i / x_i + y_i $
2	gene1	10.5	12.4		1.9	0.0830
3	gene2	6.4	7.1		0.7	0.0519
4	gene3	8	8.5		0.5	0.0303
5	gene4	10.8	11.4		0.6	0.0270
6	gene5	5.6	6.7		1.1	0.0894
7	gene6	8.4	8.9		0.5	0.0289
8	gene7	6.2	7		0.8	0.0606
9	gene8	6.1	6.8		0.7	0.0543
10	gene9	6.6	6.5		0.1	0.0076
11	gene10	5.1	5.8		0.7	0.0642

$$D = \sum_{i=1}^n |x_i - y_i| \quad \text{マンハッタン距離} = 1.9+0.7+0.5+0.6+1.1+0.5+0.8+0.7+0.1+0.7 = 7.6$$

$$D = \max(|x_i - y_i|) \quad \text{最大距離} = \max(1.9, 0.7, 0.5, 0.6, 1.1, 0.5, 0.8, 0.7, 0.1, 0.7) = 1.9$$

$$D = \sum_{i=1}^n \frac{|x_i - y_i|}{|x_i + y_i|} \quad \text{キャンベラ距離} = 0.0830+0.0519+0.0303+\dots+0.0642 = 0.4972$$

- 前処理 | フィルタリング | 分散が小さいものを除去 (last modified 2013/11/15)
- 前処理 | 同じ遺伝子名を持つものをまとめる (last modified 2013/7/31)
- 解析 | 基礎 | 共通遺伝子の抽出 (last modified 2013/6/2)
- 解析 | 基礎 | **ベクトル間の距離** (last modified 2013/11/22)
- 解析 | 基礎 | 遺伝子ごとの各種統計量の算出 (last modified 2013/6/2)

解析 | 基礎 | ベクトル間の距離

二つのベクトル間の距離を定義する方法は多数存在します。ここでは10 genes × 2 samplesのデータファイル ([sample19.txt](#))を読み込んで二つのサンプル間の距離をいくつかの方法で算出します。

「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. 10 genes × 2 samplesのデータファイル ([sample19.txt](#)) の場合:

```
in_f <- "sample19.txt"
#データファイルの読み込み
data <- read.table(in_f)
#本番
dist(t(data), method="euclidean")
dist(t(data), method="manhattan")
dist(t(data), method="maximum")
dist(t(data), method="canberra")
1 - cor(data, method="pearson")
dist(t(data), method="binary")
dist(t(data), method="minkowski")
1 - cor(data, method="spearman")
```

```
R Console
> #本番
> dist(t(data), method="euclidean") #ユークリッド (Euclidean) 距離
      sample1
sample2 2.792848
> dist(t(data), method="manhattan") #マンハッタン (Manhattan) 距離
      sample1
sample2    7.6
> dist(t(data), method="maximum") #チェビシエフ (Chebyshev) 距離
      sample1
sample2    1.9
> dist(t(data), method="canberra") #キャンベラ (Canberra) 距離
      sample1
sample2 0.4972074
> 1 - cor(data, method="pearson") #1 - Pearson相関係数
      sample1 sample2
sample1 0.0000000 0.02414407
sample2 0.02414407 0.00000000
>
> dist(t(data), method="binary") #ハミング (Hamming) 距離
      sample1
sample2    0
> dist(t(data), method="minkowski") #ミンコフスキー (Minkowski) 距離
      sample1
sample2 2.792848
> 1 - cor(data, method="spearman") #1 - Spearman相関係数
      sample1 sample2
sample1 0.0000000 0.1333333
sample2 0.1333333 0.0000000
```

	A	B	C
1	ID	sample1	sample2
2	gene1	10.5	12.4
3	gene2	6.4	7.1
4	gene3	8	8.5
5	gene4	10.8	11.4
6	gene5	5.6	6.7
7	gene6	8.4	8.9
8	gene7	6.2	7
9	gene8	6.1	6.8
10	gene9	6.6	6.5
11	gene10	5.1	5.8

他にどんな距離を利用可能か調べたい場合は...

?関数名で詳細な使用法を学ぶ

R Console

```
> ?dist
starting httpd help
>
> dist(t(data))
      sample1
sample2 2.792848
> |
```

Description

This function computes and returns the distance matrix computed by using the specified distance measure to compute the distances between the rows of a data matrix.

Usage

```
dist(x, method = "euclidean", diag = FALSE, upper = FALSE, p = 2)
```

```
as.dist(m, diag = FALSE, upper = FALSE)
## Default S3 method:
as.dist(m, diag = FALSE, upper = FALSE)
```

ユークリッド距離でよければ、「method="xxx"」のところに記述しなくてもいいようだ

```
## S3 method for class 'dist'
print(x, diag = NULL, upper = NULL,
      digits = getOption("digits"), justify = "none",
      right = TRUE, ...)
```

```
## S3 method for class 'dist'
as.matrix(x, ...)
```

Arguments

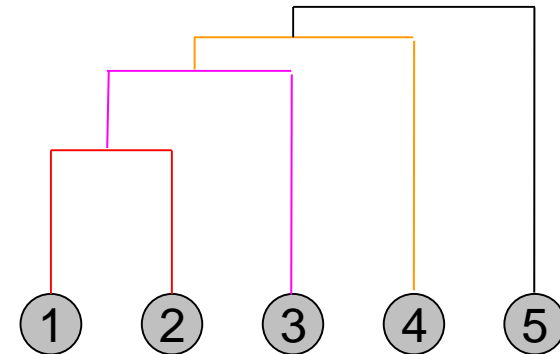
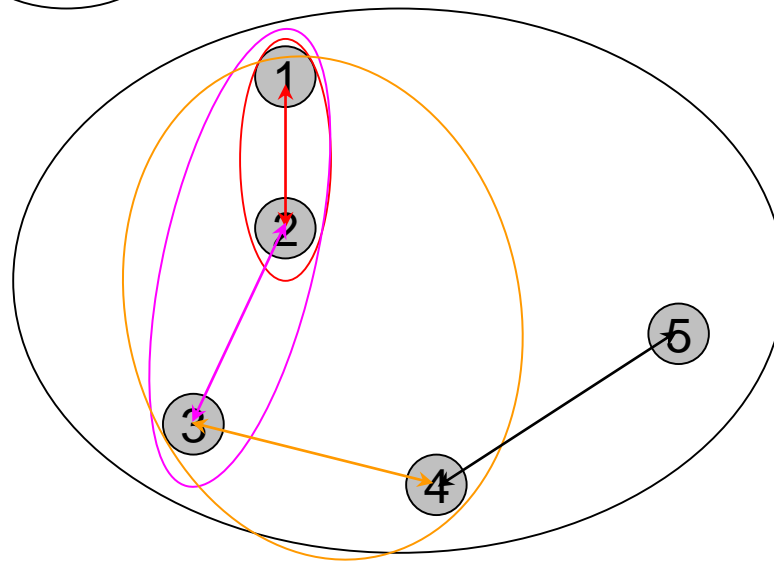
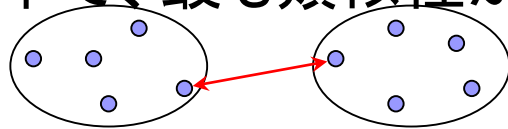
“binary”や“minkowski”というものも指定できるようだが「1-相関係数」を指定することはできないようだ...orz

x a numeric matrix, data frame or "dist" object.
method the distance measure to be used. This must be one of "euclidean", "maximum", "manhattan", "canberra", "binary" or "minkowski". Any unambiguous substring can be given.

クラスターをまとめる方法

■ 最短距離法（単連結法；single-linkage）

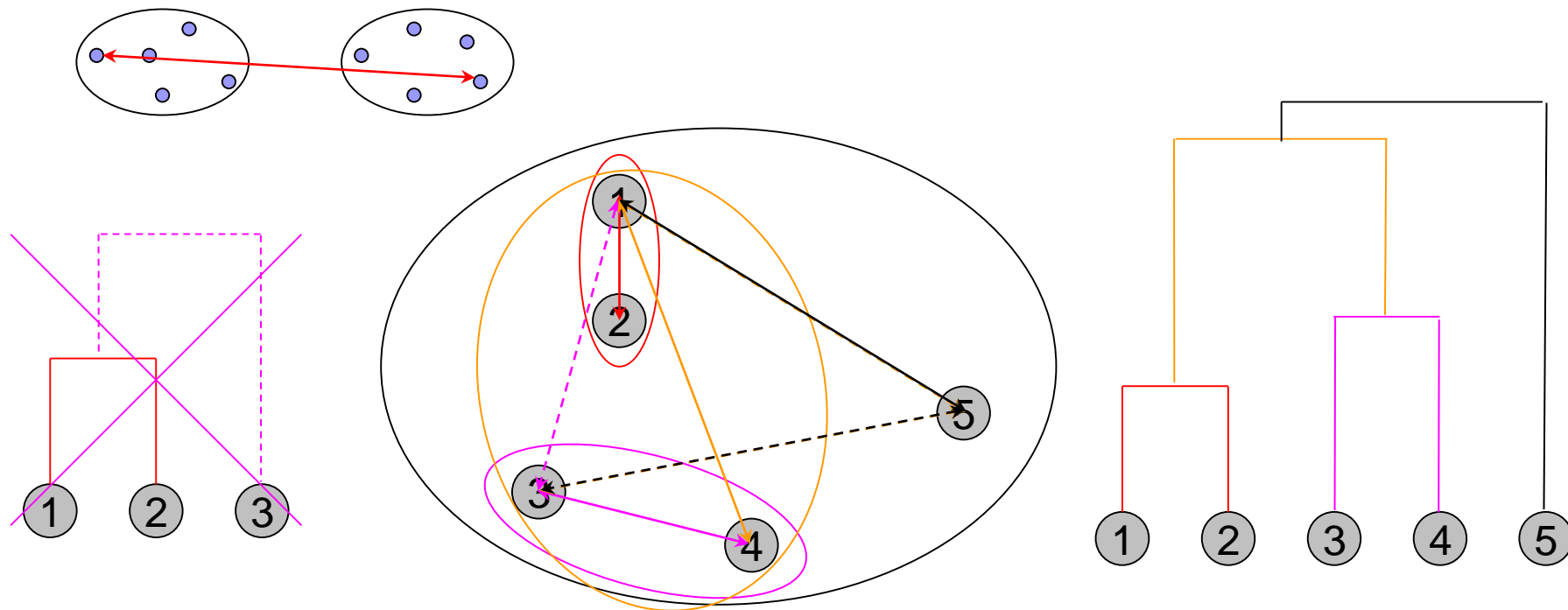
- 二つのクラスター間の類似度を、それぞれに含まれる要素対の中で、最も類似性が高い対の間の類似度で定義



クラスターをまとめる方法

■ 最長距離法（完全連結法；complete-linkage）

- クラスターをmergeするときの基準として、最遠距離を用いる

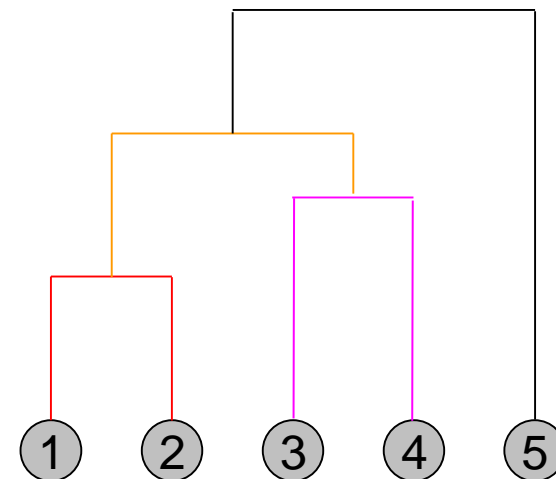
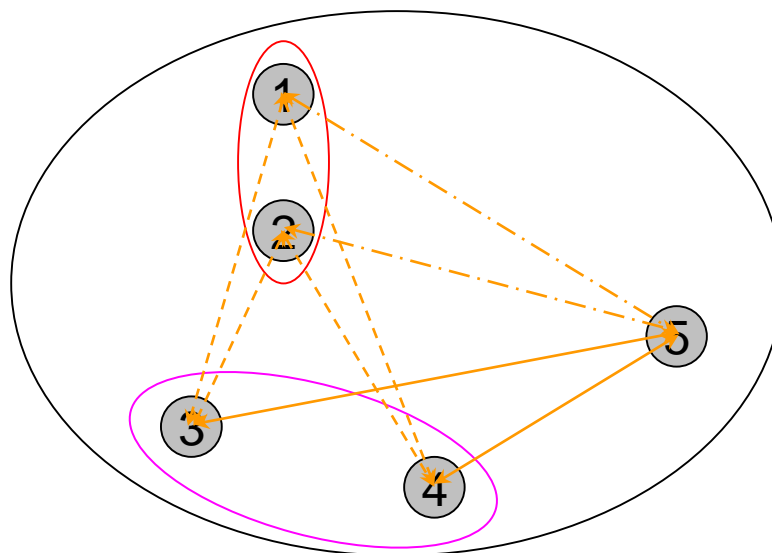
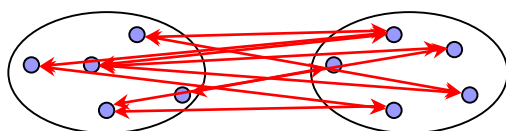


電気泳動波形解析（ピークアラインメント）にも応用可能

クラスターをまとめる方法

■ 群平均法（平均連結法；average-linkage）

- クラスターをmergeするときの基準として、群間平均距離を用いる



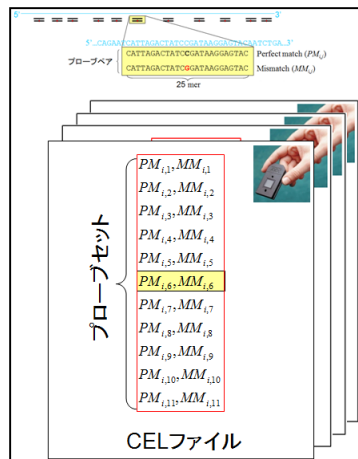
用いる方法によって得られるデンドログラム（樹形図）が異なる

他のクラスターをまとめる手段

- 重心法 (Centroid) : 重心間距離を利用
- ウォード法 : 群内平方和の増加量が最小となるクラスターと併合
- メディアン (Median) 法 : 群間距離の中央値を利用
- McQuitty法...
- 可変 (flexible) 法...

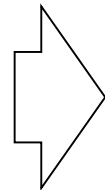
分野にもよるらしいが群平均法が最もよく利用されている?!(ウォード法も?!)...。
いろいろ試して総合的に判断することが重要

得られる樹形図の可能性は無数



生データファイル

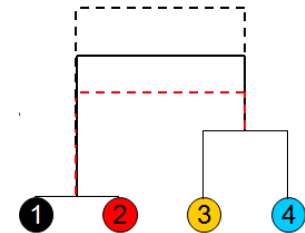
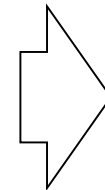
前処理法



	S1	S2	S3	S4	...
gene 1	$x_{1,1}$	$x_{1,2}$	$x_{1,3}$	$x_{1,4}$...
gene 2	$x_{2,1}$	$x_{2,2}$	$x_{2,3}$	$x_{2,4}$...
...
gene i	$x_{i,1}$	$x_{i,2}$	$x_{i,3}$	$x_{i,4}$...
...
gene n	$x_{n,1}$	$x_{n,2}$			

遺伝子発現行列

クラスタリング

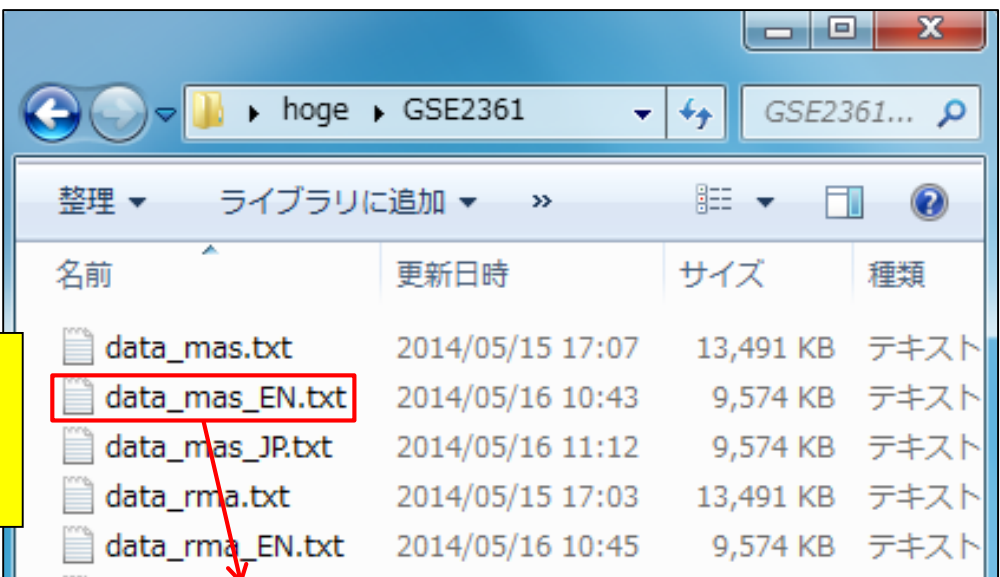
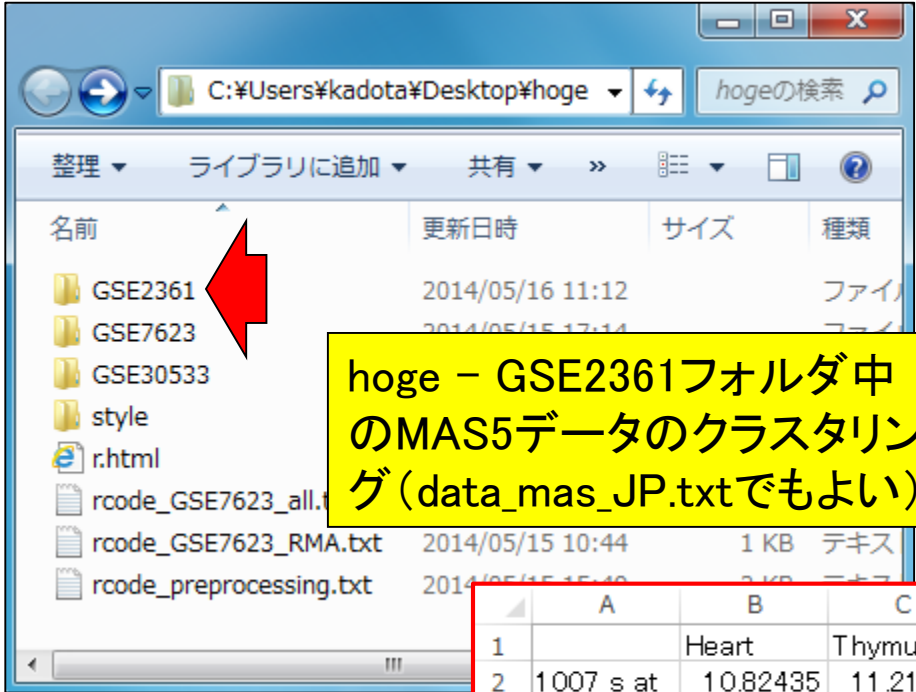


<p>①前処理法</p> <ul style="list-style-type: none"> ・MAS ・RMA ・RobLoxBioC 	<p>②スケールング</p> <ul style="list-style-type: none"> ・対数変換 ・相対値(0~1) ・Z-score化 	<p>③距離</p> <ul style="list-style-type: none"> ・1-相関係数 ・ユークリッド ・マンハッタン ・キャンベラ ... 	<p>④群の併合</p> <ul style="list-style-type: none"> ・単連結法 ・完全連結法 ・平均連結法 ・ワード法 ...
--	--	---	--

サンプル間クラスタリングをやってみよう

Affymetrix GeneChip

- Ge et al., *Genomics*, **86**: 127-141, 2005
 - GSE2361、GPL96 (Affymetrix Human Genome U133A Array)、22,283 probesets
 - ヒト36サンプル: Heart (心臓)、Thymus (胸腺)、Spleen (脾臓)、Ovary (卵巣)、Kidney (腎臓)、Skeletal Muscle (骨格筋)、Pancreas (膵臓)、Prostate (前立腺)、...



	A	B	C	D	E	F	G	H	I	J
1		Heart	Thymus	Spleen	Ovary	Kidney	Skeletal_Mu	Pancreas	Prostate	Small_I
2	1007_s_at	10.82435	11.21261	9.898072	10.8948	11.75531	10.92032	11.61243	11.22101	11.48
3	1053_at	6.621657	6.47488	7.371465	4.168622	7.350244	7.209918	8.362777	7.176571	7.75
4	117_at	8.259603	8.647364	8.689724	7.875649	5.083001	8.165044	7.96707	8.418082	5.719
5	121_at	11.26699	11.04186	11.39999	11.02855	13.13267	11.39138	12.32899	11.0559	11.28
6	1255_g_at	7.1757	6.477278	6.781766	7.048799	7.30767	6.600267	6.42359	6.694412	7.30
7	1294_at	9.137586	9.718507	9.083742	9.014997	8.813377	8.402562	9.146946	8.421351	10.12
8	1316_at	8.225789	7.418995	8.07469	7.980035	8.228176	7.600147	7.422269	7.209994	7.67

- 解析 | 基礎 | [ベクトル間の距離](#) (last modified 2013/11/22)
- 解析 | 基礎 | [遺伝子ごとの各種要約統計量の算出](#) (last modified 2013/6/2)
- 解析 | 基礎 | [最大発現量を示す組織の同定](#) (last modified 2013/6/2)
- 解析 | 基礎 | [似た発現パターンを持つ遺伝子の同定](#) (last modified 2013/6/2)
- 解析 | 基礎 | [平均-分散プロット](#) (last modified 2013/11/25)
- 解析 | [クラスタリング](#) | 階層的 | [について](#) (last modified 2009/8/12)
- 解析 | [クラスタリング](#) | 階層的 | [pvclust](#) (Suzuki 2006) (last modified 2010/8/5)
- 解析 | [クラスタリング](#) | 階層的 | [hclust](#) (last modified 2013/11/19)
- 解析 | [クラスタリング](#) | 階層的 | [hclust](#)後の詳細な解析 (last modified 2009/8/7)

解析 | クラスタリング | 階層的 | hclust NEW

階層的クラスタリングのやり方を示します。1.用いた前処理法(MAS5やRMAなど)、2.スケールリング方法(対数変換やZ-scoreなど)、3.距離(または非類似度)を定義する方法(ユークリッド距離など)、4.クラスターをまとめる方法(平均連結法やワード法など)でどの方法を採用するかで結果が変わってきます。

「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. サンプルデータ3の [sample3.txt](#) の場合:

サンプル間クラスタリング(距離: 1-Pearson相関係数、方法: 平均連結法(average))でR Graphics画面上に表示

```
in_f <- "sample3.txt" #入力ファイル名を指定してin_fに格納
param <- "average" #方法(method)を指定

#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.dist <- as.dist(1 - cor(data, method="pearson"))#サンプル間の距離を計算した結果をdata.distに格納
out <- hclust(data.dist, method=param) #階層的クラスタリングを実行した結果をoutに格納
plot(out) #樹形図(デンドログラム)の表示
```

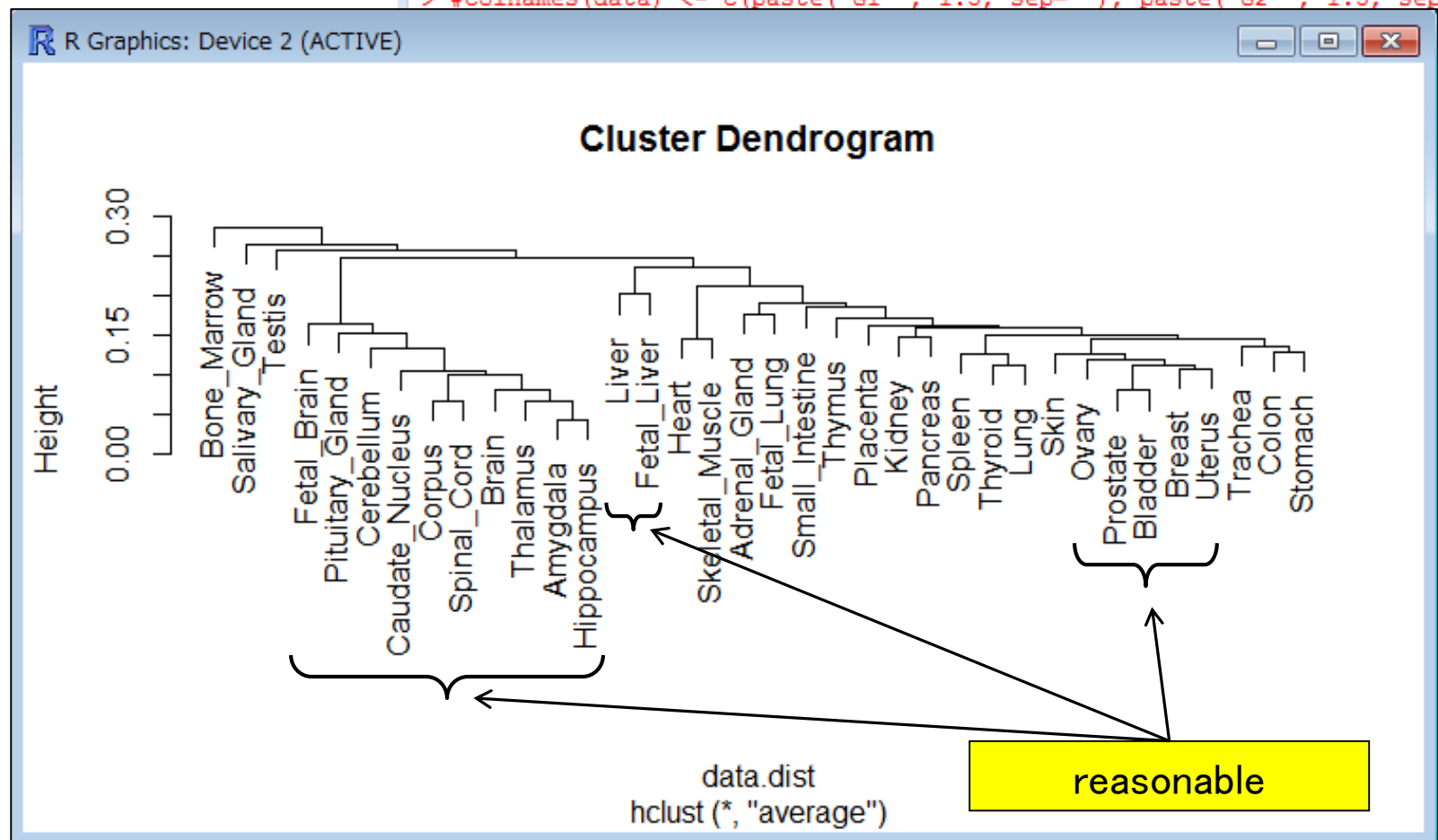
左記のテンプレートスクリプトとの違いは、入力ファイル名と Spearman相関係数の部分のみ

```
##### ↓
### MAS5データのクラスタリング ### ↓
##### ↓
in_f <- "data_mas_EN.txt" ↓
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="") ↓
#colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep="")) ↓
data.dist <- as.dist(1 - cor(data, method = "spearman")) ↓
out <- hclust(data.dist, method = "average") ↓
plot(out) ↓
<
```

rcode_clustering.txt

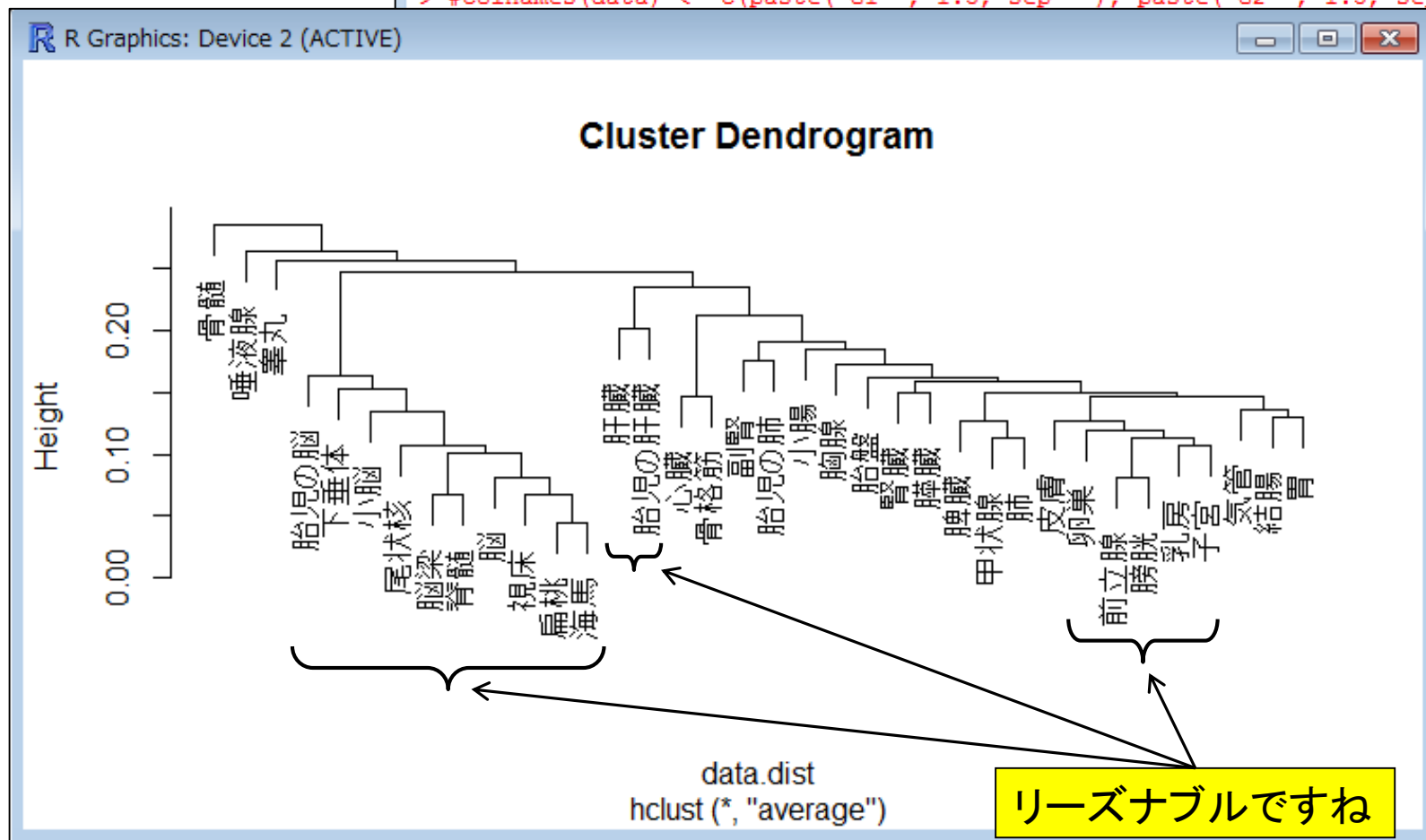

```
#####  
### MAS5データのクラスタリング ###  
#####  
in_f <- "data_mas_EN.txt"  
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")  
#colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))  
data.dist <- as.dist(1 - cor(data, method = "spearman"))  
out <- hclust(data.dist, method = "average")  
plot(out)
```

```
R Console  
> getwd()  
[1] "C:/Users/kadota/Desktop/hoge/GSE2361"  
> in_f <- "data_mas_EN.txt"  
> data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")  
> #colnames(data) <- c(paste("G1 ", 1:5, sep=""), paste("G2 ", 1:5, sep=""))
```



```
#####↓
### MAS5データのクラスタリング ###↓
#####↓
in_f <- "data_mas_JP.txt"↓
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")↓
#colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))↓
data.dist <- as.dist(1 - cor(data, method = "spearman"))↓
out <- hclust(data.dist, method = "average")↓
plot(out)↓
←
```

```
R Console
> in_f <- "data_mas_JP.txt"
> data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
> #colnames(data) <- c(paste("G1 ", 1:5, sep=""), paste("G2 ", 1:5, sep=""))
```



クラスタリング結果をファイルに保存

解析 | クラスタリング | 階層的 | hclust NEW

階層的クラスタリングのやり方を示します。1.用いた前処理法(MAS5やRMAなど)、2.スケール方法(対数変換やZ-scoreなど)、3.距離(または非類似度)を定義する方法(ユークリッド距離など)、4.クラスターをまとめる方法(平均連結法やワード法など)

3. サンプルデータ30の sample3.txt の場合:

サンプル間クラスタリング(距離: 1-Spearman相関係数、方法: 平均連結法(average))で図の大きさを指定してpng形式ファイルで保存するやり方です。

```
in_f <- "sample3.txt"      #入力ファイル名を指定してin_fに格納
out_f <- "hoge3.png"     #出力ファイル名を指定してout_fに格納
param <- "average"       #方法(method)を指定
param_fig <- c(500, 400)  #ファイル出力時の横幅と縦幅を指定(単位はピクセル)
```

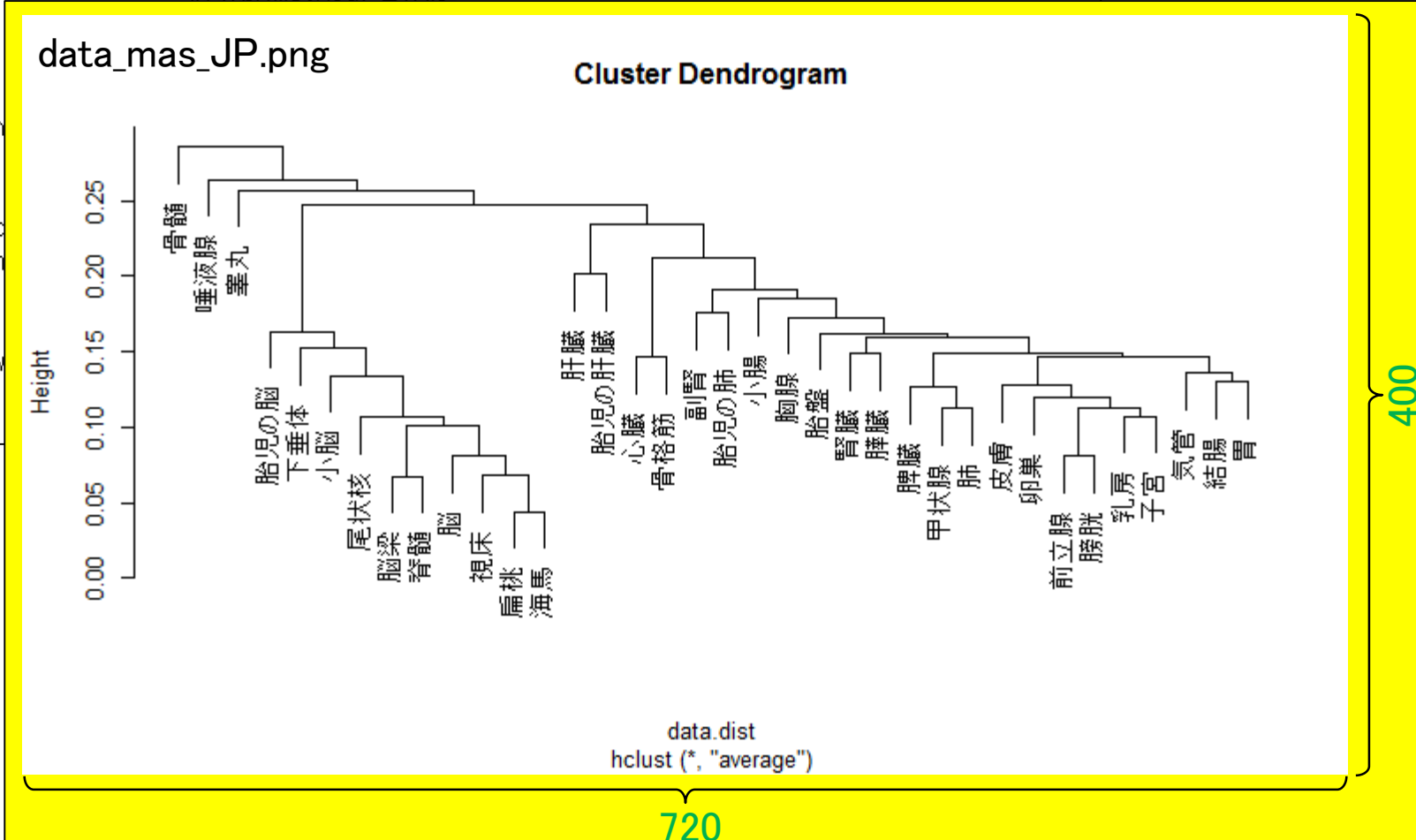
rancode_clustering_png.txt (変更点は赤矢印部分のみ)

```
#####↓
### MAS5データのクラスタリング, 距離: 1-Spearman相関係数, 方法: 平均連結法 ###↓
#####↓
in_f <- "data_mas_JP.txt"      #入力ファイル名を指定してin_fに格納↓
→ out_f <- "data_mas_JP.png"  #出力ファイル名を指定してout_fに格納↓
param <- "average"           #方法(method)を指定↓
→ param_fig <- c(720, 400)     #ファイル出力時の横幅と縦幅を指定(単位はピクセル)↓
↓
#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.dist <- as.dist(1 - cor(data, method="spearman"))#サンプル間の距離を計算した結果をdata.distに格納↓
out <- hclust(data.dist, method=param) #階層的クラスタリングを実行した結果をoutに格納↓
↓
#ファイルに保存↓
→ png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2])#出力ファイルの各種パラメータを指定↓
plot(out)                    #樹形図(デンドログラム)の表示↓
→ dev.off()                  #おまじない↓
```

クラスタリング結果をファイルに保存

rcode_clustering_png.txt (変更点は赤矢印部分のみ)

```
#####↓
### MAS5データのクラスタリング, 距離: 1-Spearman相関係数, 方法: 平均連結法 ###↓
#####↓
in_f <- "data_mas_JP.txt" #入力ファイル名を指定してin_fに格納↓
→ out_f <- "data_mas_JP.png" #出力ファイル名を指定してout_fに格納↓
param <- "average" #方法(method)を指定↓
→ param_fig <- c(720, 400)
↓
#入力ファイルの読み込み↓
data <- read.table(in_f, h
↓
#本番↓
data.dist <- as.dist(1 - c
out <- hclust(data.dist, m
↓
#ファイルに保存↓
→ png(out_f, pointsize=13, w
plot(out)
→ dev.off()
```

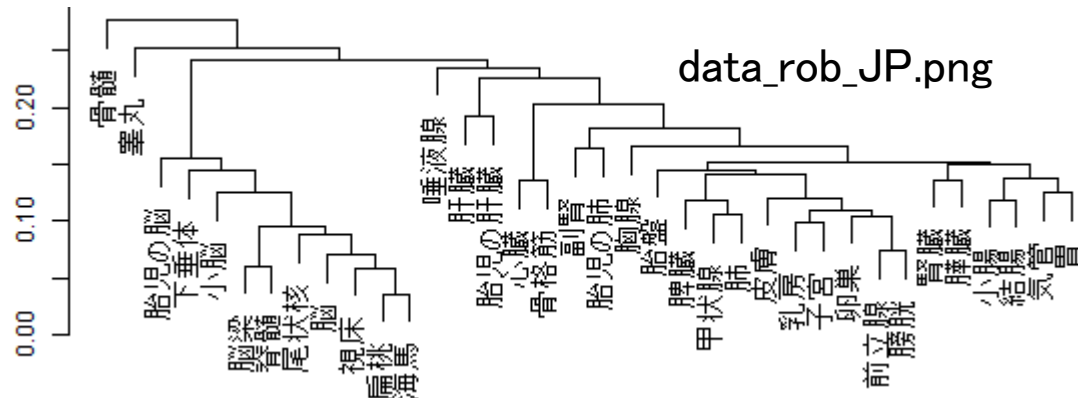
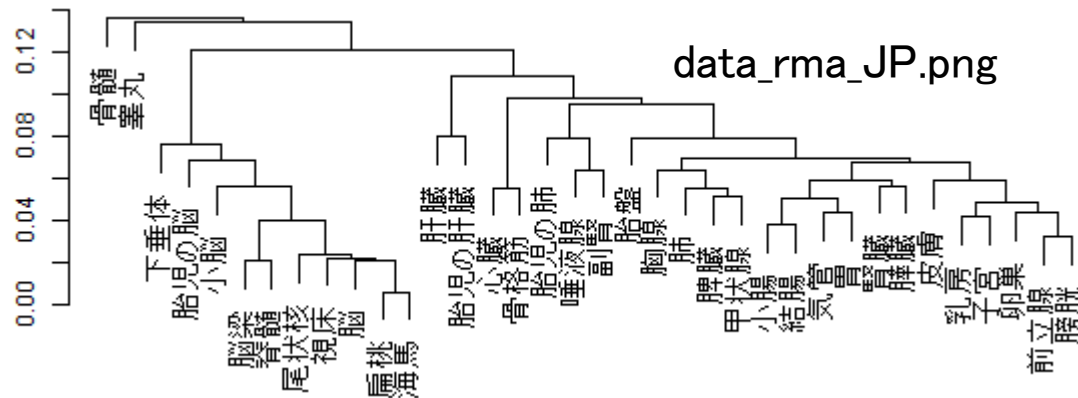
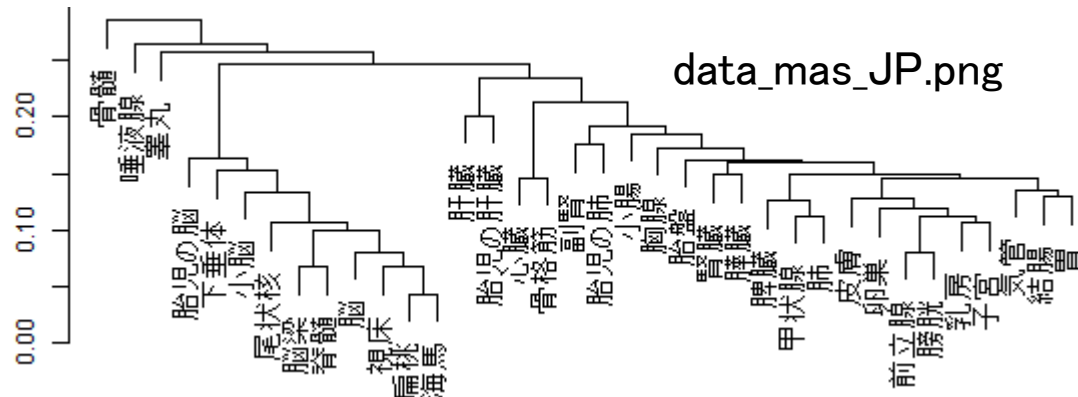


サンプル間クラスタリング (GSE2361)

GSE2361のサンプル間クラスタリングをMAS5, RMA, およびRobLoxBioC前処理法を適用したデータについても行ってみよう。

- 1 - Spearman相関係数
- 平均連結法

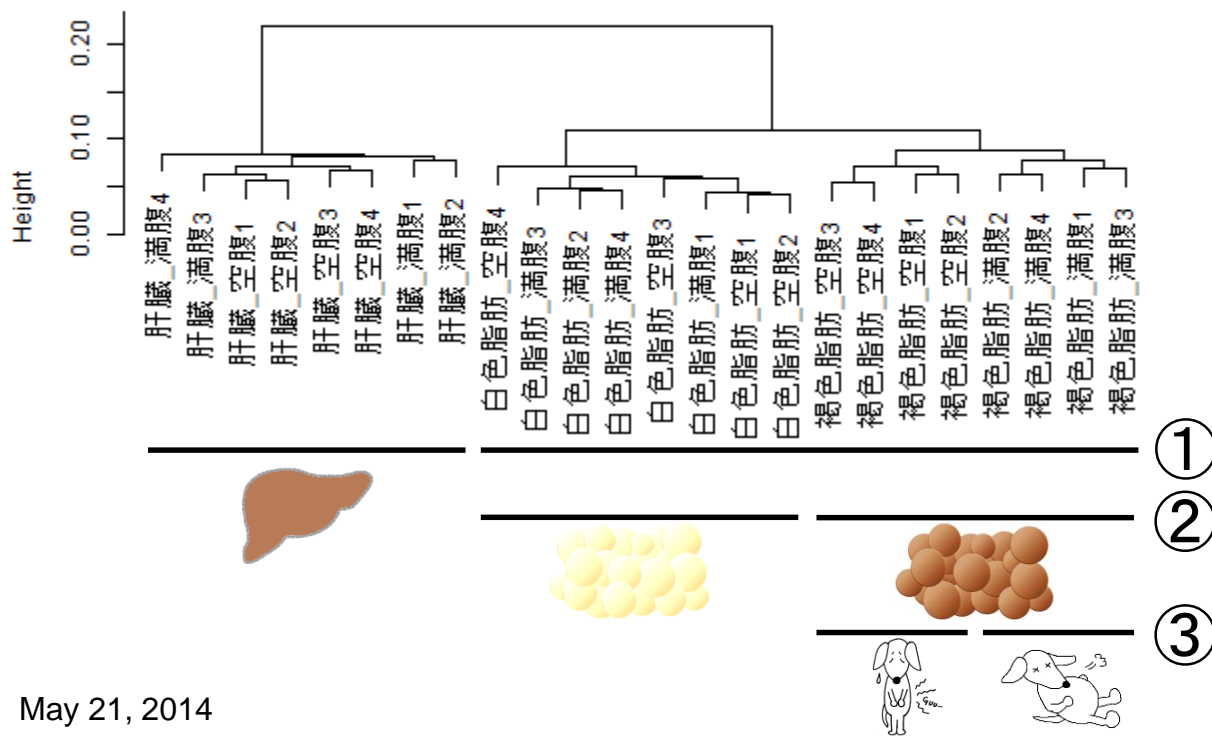
ヒント: rcode_clustering_png.txt
をテンプレートとすればよい。
正解: 右図の結果。どの前処理法を適用しても似たような樹形図が得られていることがわかる。



サンプル間クラスタリング (GSE7623)

- Nakai et al., *Biosci Biotechnol Biochem.*, 72: 139–148, 2008
- GSE7623、GPL1355 (Affymetrix Rat Genome 230 2.0 Array)、31,099 probesets
- ラット24サンプル: Brown adipose tissue (褐色脂肪組織; BAT) 8サンプル、White adipose tissue (白色脂肪組織; WAT) 8サンプル、Liver (肝臓; LIV) 8サンプル
 - BAT 8サンプル: 通常 (BAT_fed) 4サンプル 対 24時間絶食 (BAT_fas) 4サンプル
 - WAT 8サンプル: 通常 (WAT_fed) 4サンプル 対 24時間絶食 (WAT_fas) 4サンプル
 - LIV 8サンプル: 通常 (LIV_fed) 4サンプル 対 24時間絶食 (LIV_fas) 4サンプル

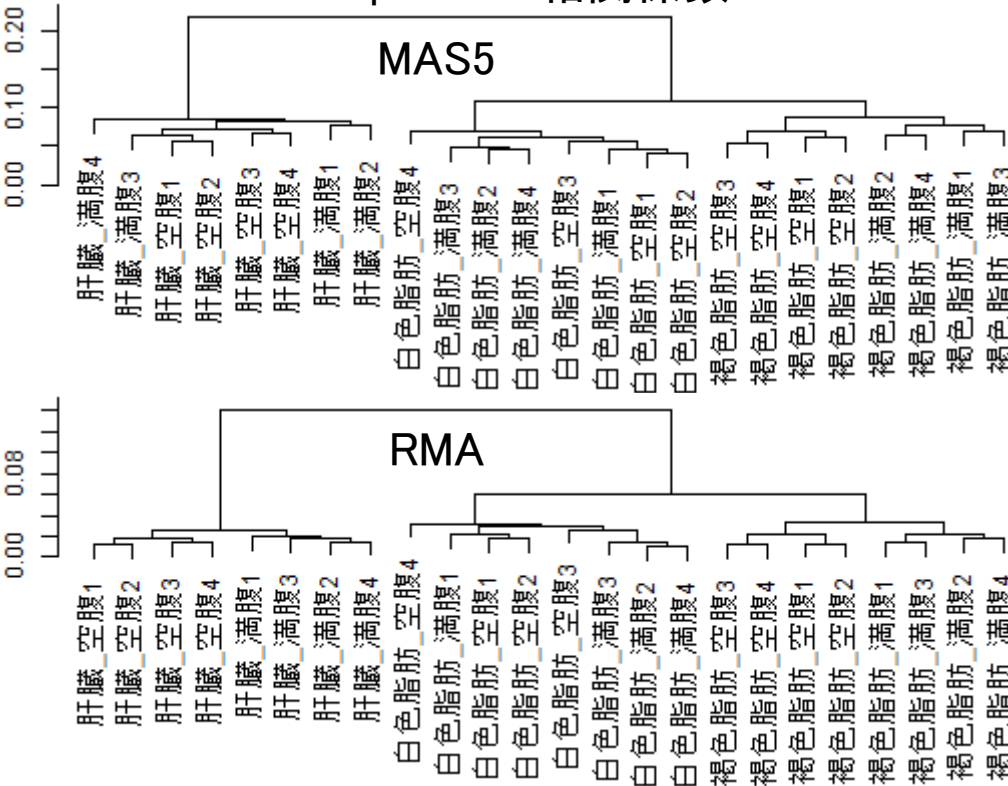
Cluster Dendrogram



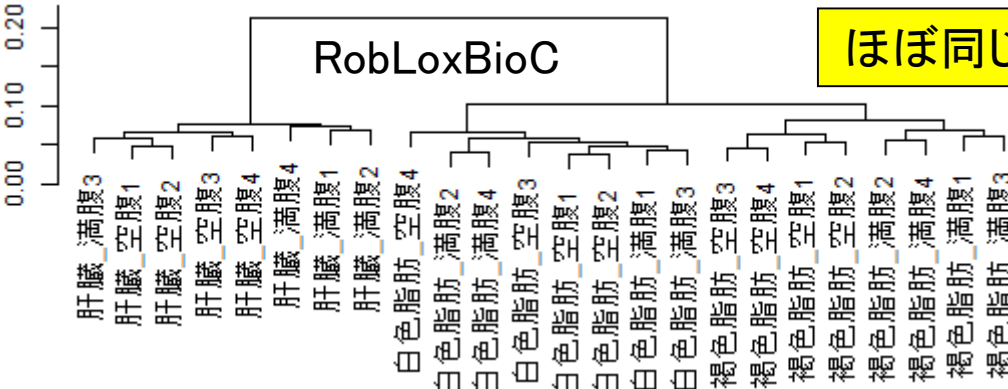
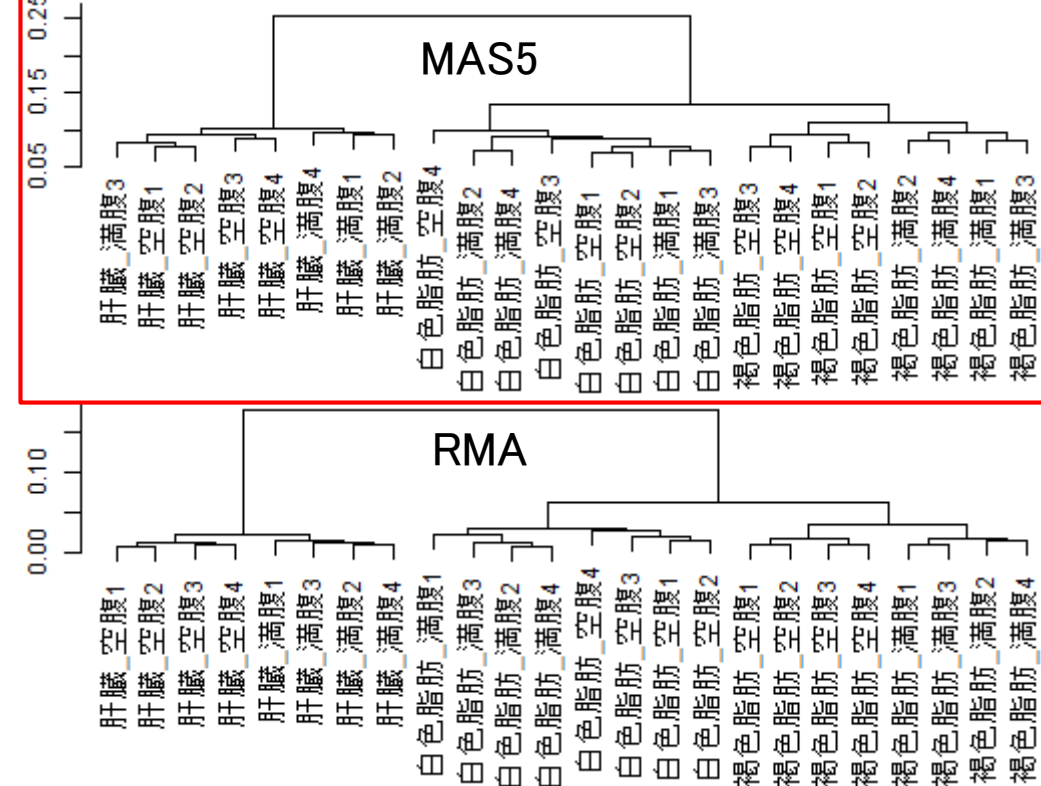
rancode_clustering_png.txtの実行結果。
 ①肝臓と脂肪間で大きく二つのクラスターに分かれている。
 ②脂肪の中でも白色脂肪と褐色脂肪に分かれている。
 ③褐色脂肪は空腹(24時間絶食)と満腹(通常)できれいに分かれている。

サンプル間クラスタリング (GSE7623)

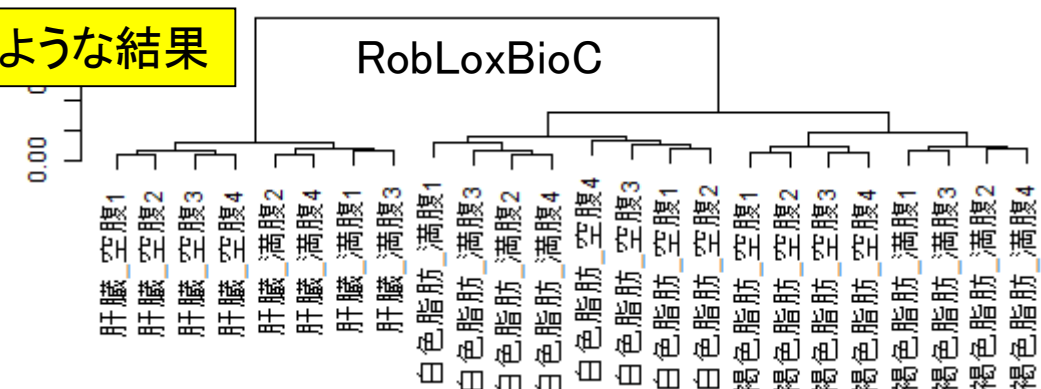
1 - Spearman相関係数



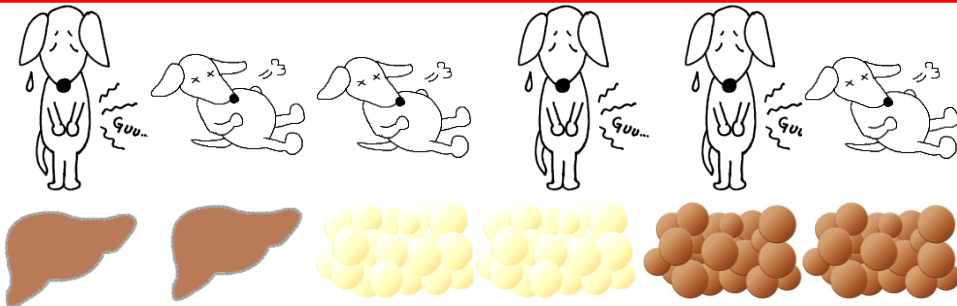
1 - Pearson相関係数



ほぼ同じような結果



原著論文と比較



ばっちりFig. 1 (の一部)を再現できました

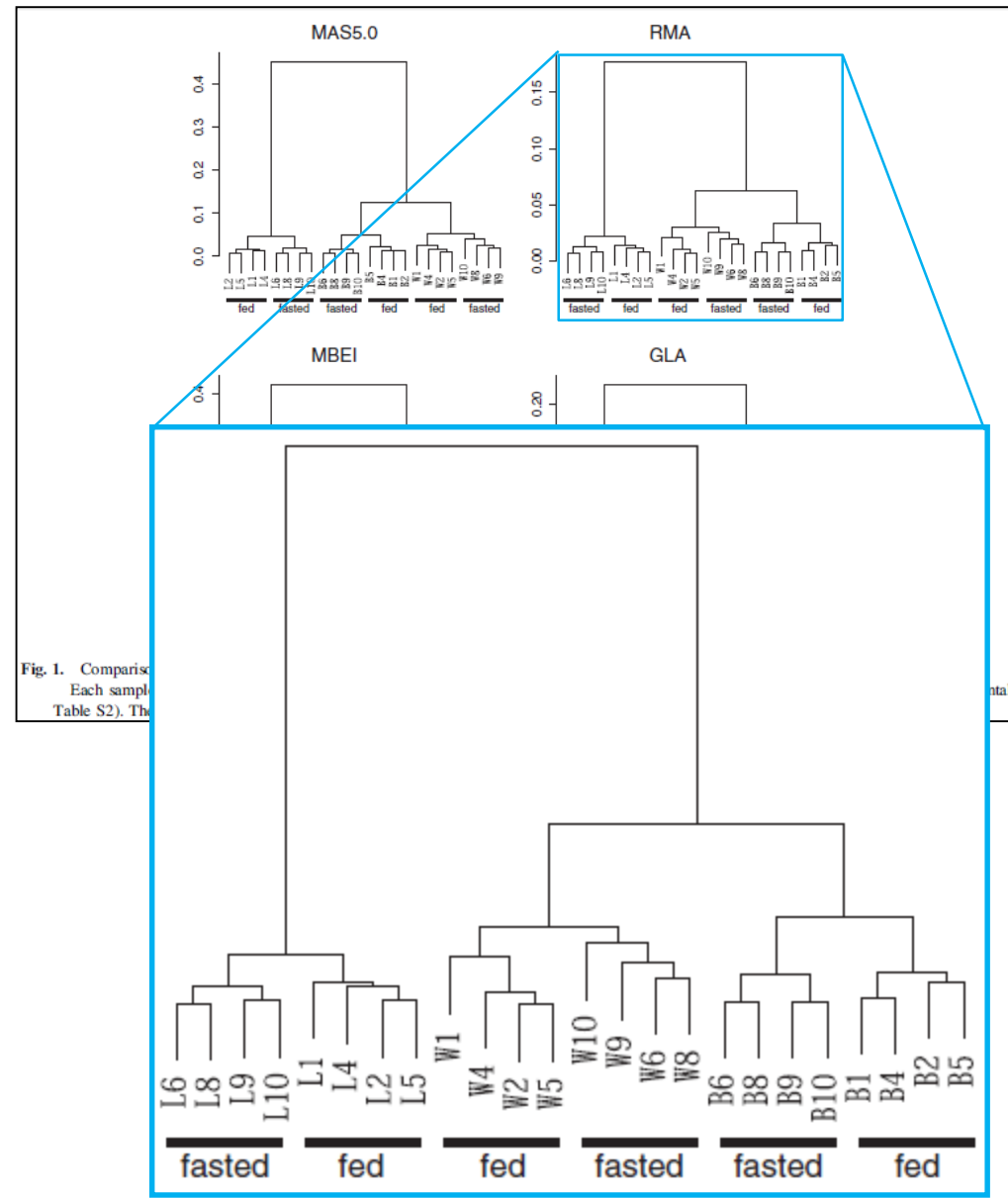
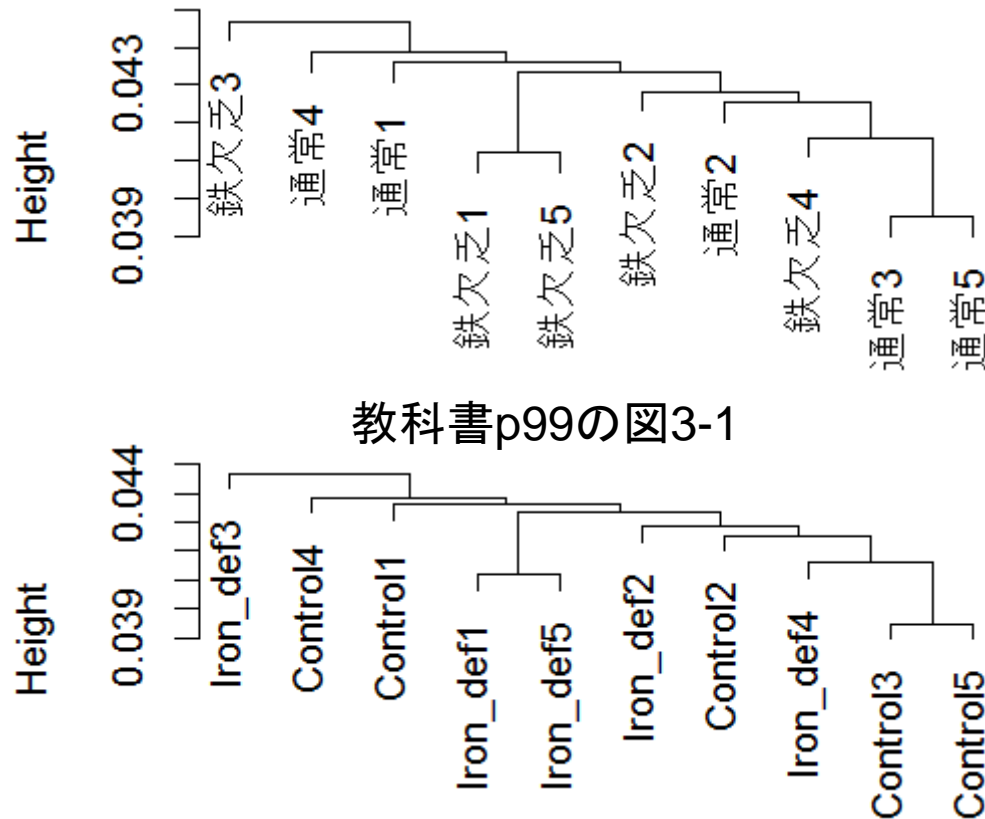


Fig. 1. Comparison of clustering methods. Each sample is labeled as in Table S2. The

サンプル間クラスタリング (GSE30533)

- Kamei et al., *PLoS One*, 8: e65732, 2013
 - GSE30533、GPL1355 (Affymetrix Rat Genome 230 2.0 Array)、31,099 probesets
 - ラット10サンプル: 全てLiver (肝臓) サンプル
 - iron-deficient diet (Iron_def) 5サンプル 対 control diet (Control) 5サンプル

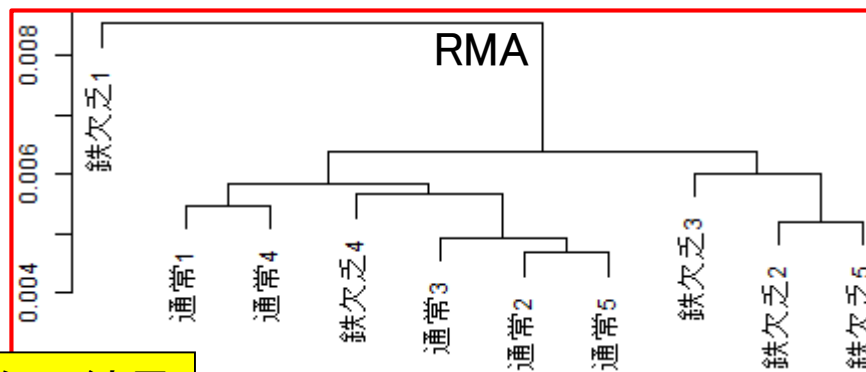
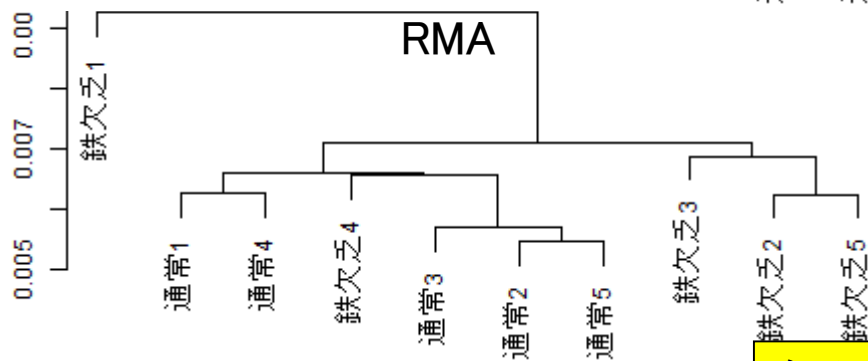
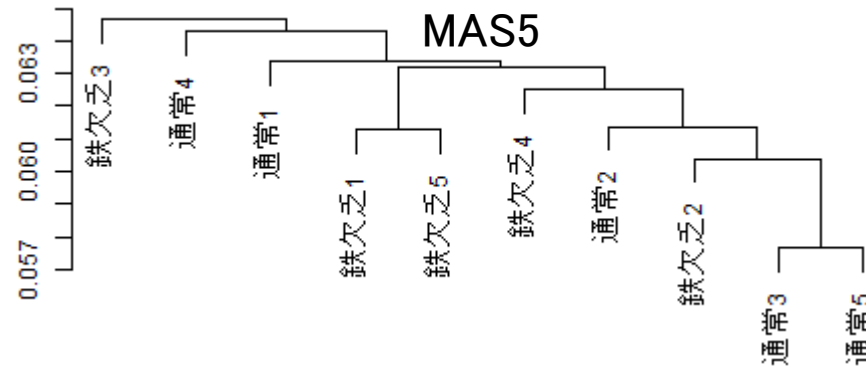
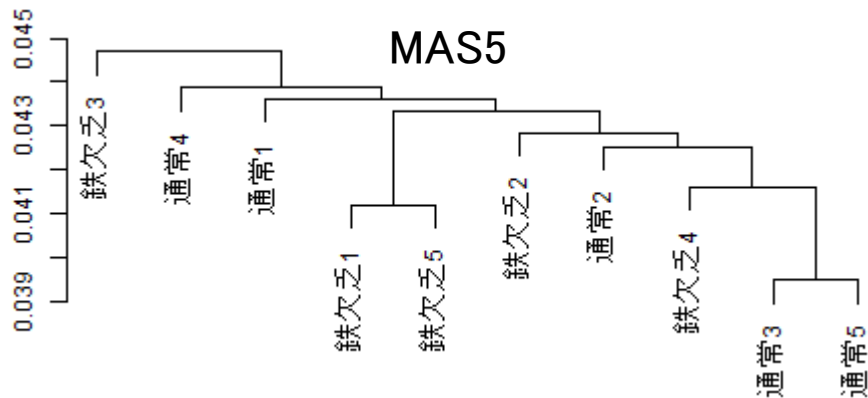


rancode_clustering_png.txtの実行結果。
肝臓全体の発現プロファイルが通常状態と鉄欠乏状態という違い程度では明確に区別できない、ということかもしれない…。

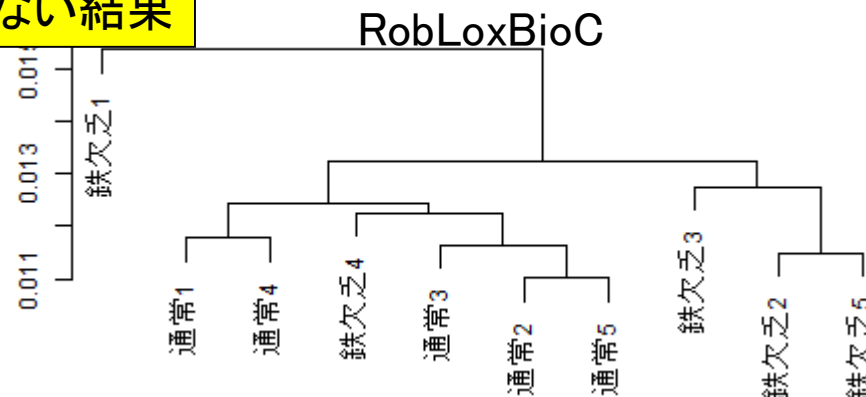
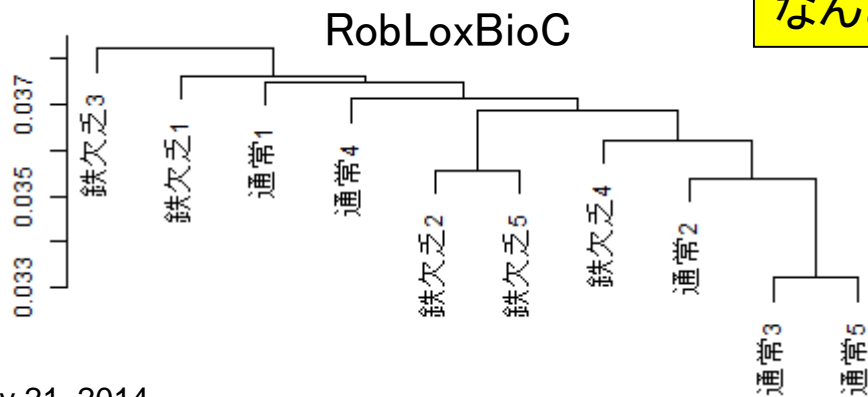
サンプル間クラスターリング (GSE7623)

1 - Spearman相関係数

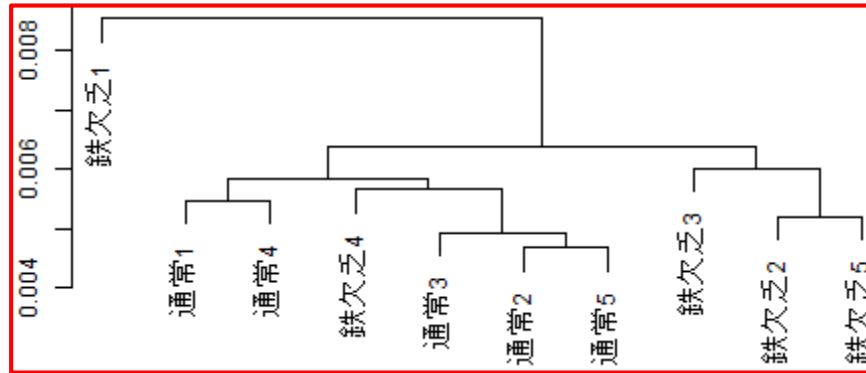
1 - Pearson相関係数



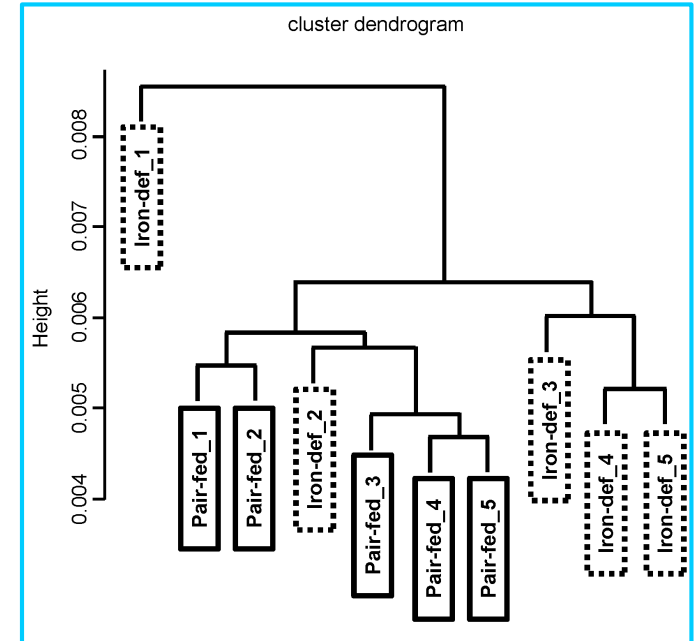
なんともいえない結果



原著論文と比較



原著論文のFigure S1



(サンプルのラベル番号が異なるだけで実質的には) 同じ結果

同一アレイデータはマージ可能

■ Affymetrix GeneChip

- Ge et al., *Genomics*, **86**: 127–141, 2005
 - GSE2361、GPL96 (Affymetrix Human Genome U133A Array)、22,283 probesets
 - ヒト36サンプル: Heart (心臓)、Thymus (胸腺)、Spleen (脾臓)、Ovary (卵巣)、Kidney (腎臓)、Skeletal Muscle (骨格筋)、Pancreas (膵臓)、Prostate (前立腺)、…
- Nakai et al., *Biosci Biotechnol Biochem.*, **72**: 139–148, 2008
 - GSE7623、GPL1355 (Affymetrix Rat Genome 230 2.0 Array)、31,099 probesets
 - ラット24サンプル: Brown adipose tissue (褐色脂肪組織; BAT) 8サンプル、White adipose tissue (白色脂肪組織; WAT) 8サンプル、Liver (肝臓; LIV) 8サンプル
 - BAT 8サンプル: 通常 (BAT_fed) 4サンプル 対 24時間絶食 (BAT_fas) 4サンプル
 - WAT 8サンプル: 通常 (WAT_fed) 4サンプル 対 24時間絶食 (WAT_fas) 4サンプル
 - LIV 8サンプル: 通常 (LIV_fed) 4サンプル 対 24時間絶食 (LIV_fas) 4サンプル
- Kamei et al., *PLoS One*, **8**: e65732, 2013
 - GSE30533、GPL1355 (Affymetrix Rat Genome 230 2.0 Array)、31,099 probesets
 - ラット10サンプル: 全てLiver (肝臓) サンプル
 - iron-deficient diet (Iron_def) 5サンプル 対 control diet (Control) 5サンプル

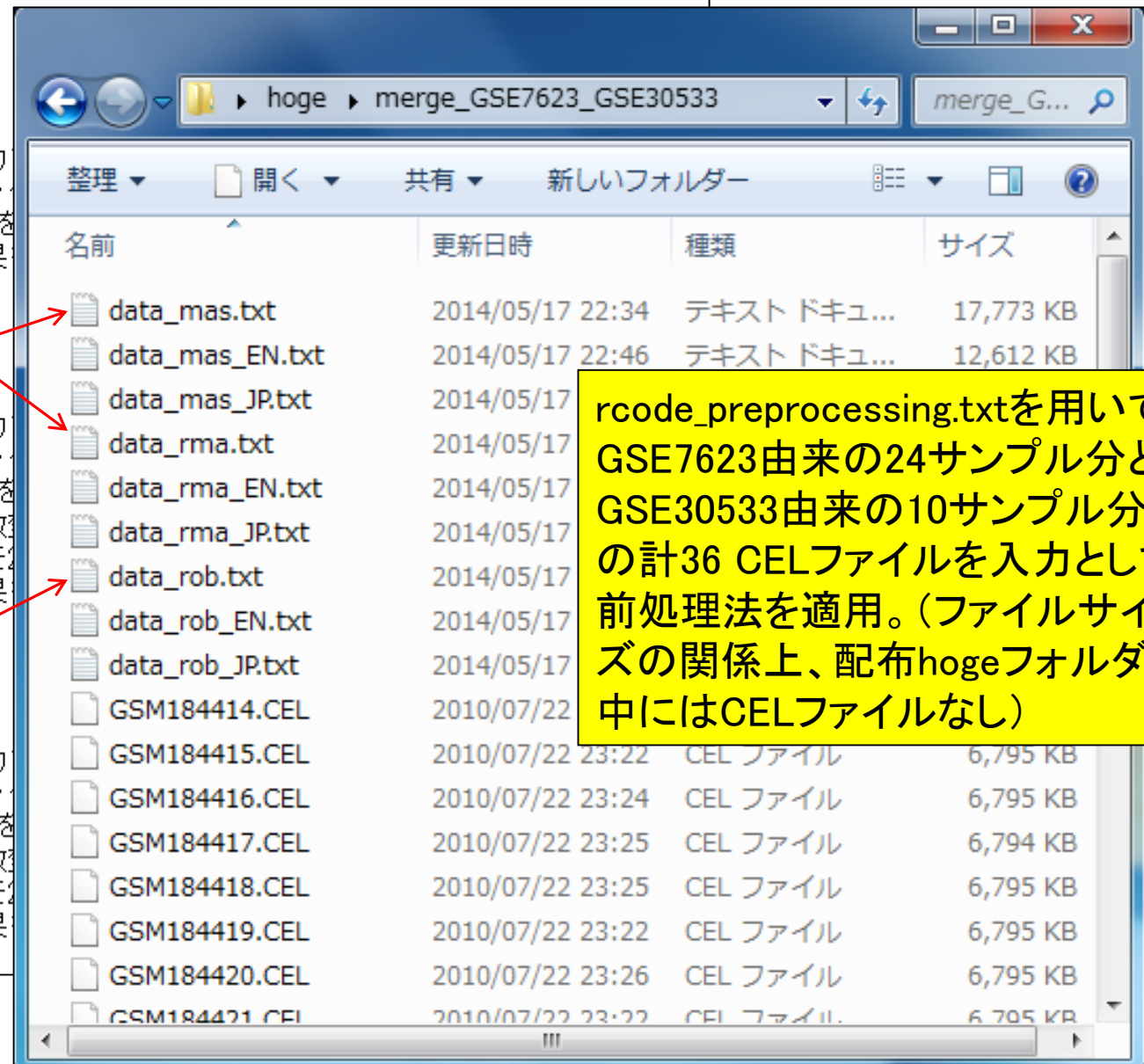
この2つの論文は同一プラットフォーム(同一アレイ)を利用している。
3' 発現アレイ用いることで、他の多くのデータセットとの比較が可能

```

##### ↓
### 作業ディレクトリ中のCELファイルの読み込み ### ↓
##### ↓
library(affy) #パッケージの読み込み↓
hoge <- ReadAffy() #*.CELファイルの読み込み↓
↓
↓
##### ↓
### RMA前処理法実行 ### ↓
##### ↓
out_f <- "data_rma.txt" #出力
library(affy) #パッケージの読み込み↓
eset <- rma(hoge) #RMAを適用↓
write.exprs(eset, file=out_f) #結果をファイルに保存↓
↓
##### ↓
### MAS5前処理法実行 ### ↓
##### ↓
out_f <- "data_mas.txt" #出力
library(affy) #パッケージの読み込み↓
eset <- mas5(hoge) #MAS5を適用↓
exprs(eset)[exprs(eset) < 1] <- 1 #対数変換の底を2に設定↓
exprs(eset) <- log(exprs(eset), 2)
write.exprs(eset, file=out_f) #結果をファイルに保存↓
↓
##### ↓
### RMX (RobLoxBioC)前処理法実行 ### ↓
##### ↓
out_f <- "data_rob.txt" #出力
library(RobLoxBioC) #パッケージの読み込み↓
eset <- robloxbioc(hoge) #rmxを適用↓
exprs(eset)[exprs(eset) < 1] <- 1 #対数変換の底を2に設定↓
exprs(eset) <- log(exprs(eset), 2)
write.exprs(eset, file=out_f) #結果をファイルに保存↓
↓

```

- 正規化 | Affymetrix GeneChip | [RMA \(Irizarry 2003\)](#)
- 正規化 | Affymetrix GeneChip | [MAS5.0 \(Hubbell 2002\)](#)
- 正規化 | Affymetrix GeneChip | [rmx \(Kohl 2010\)](#)

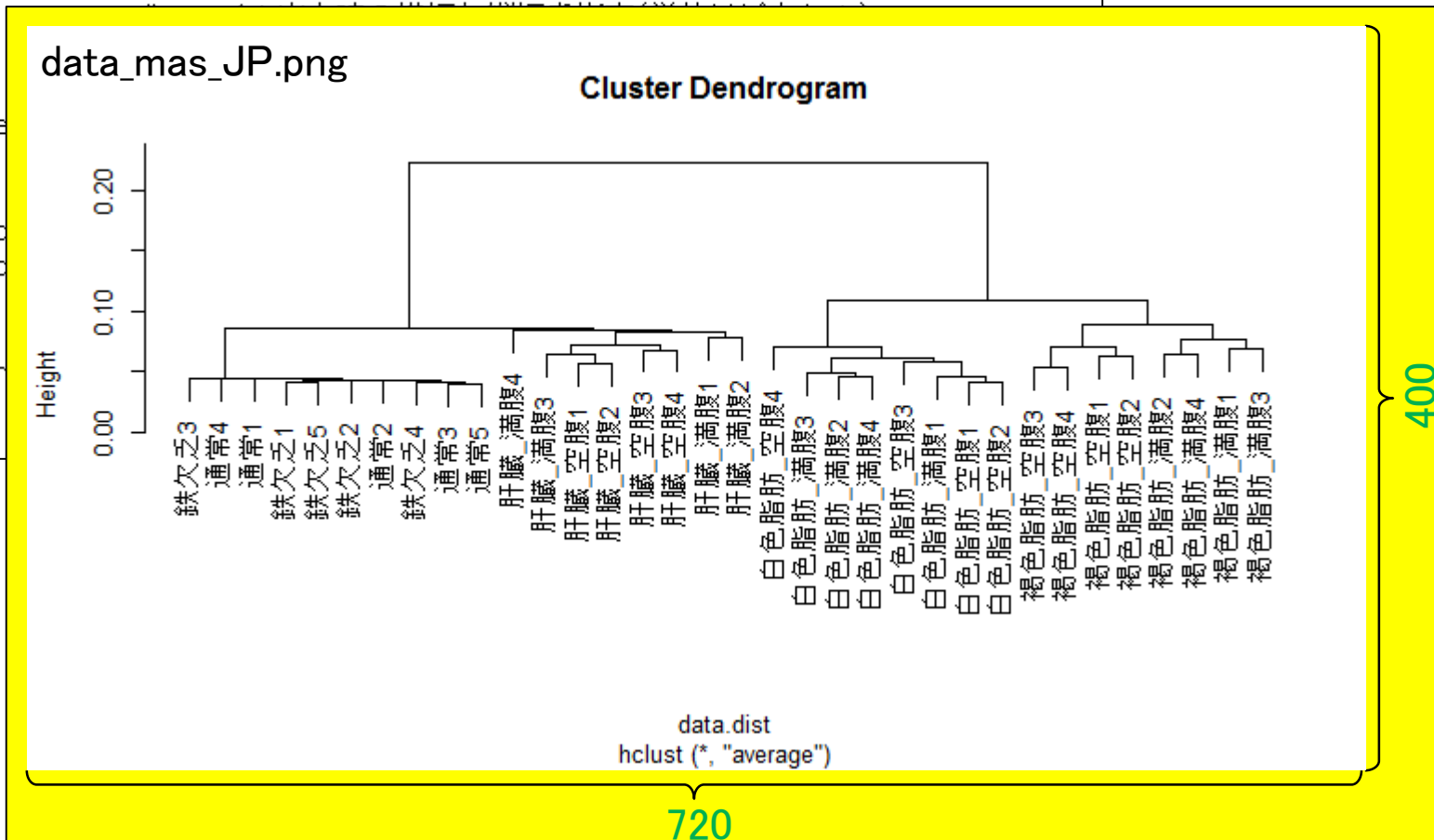


rcode_preprocessing.txtを用いて GSE7623由来の24サンプル分と GSE30533由来の10サンプル分の計36 CELファイルを入力として前処理法を適用。(ファイルサイズの関係上、配布hogeフォルダ中にはCELファイルなし)

クラスタリング結果をファイルに保存

rcode_clustering_png.txt

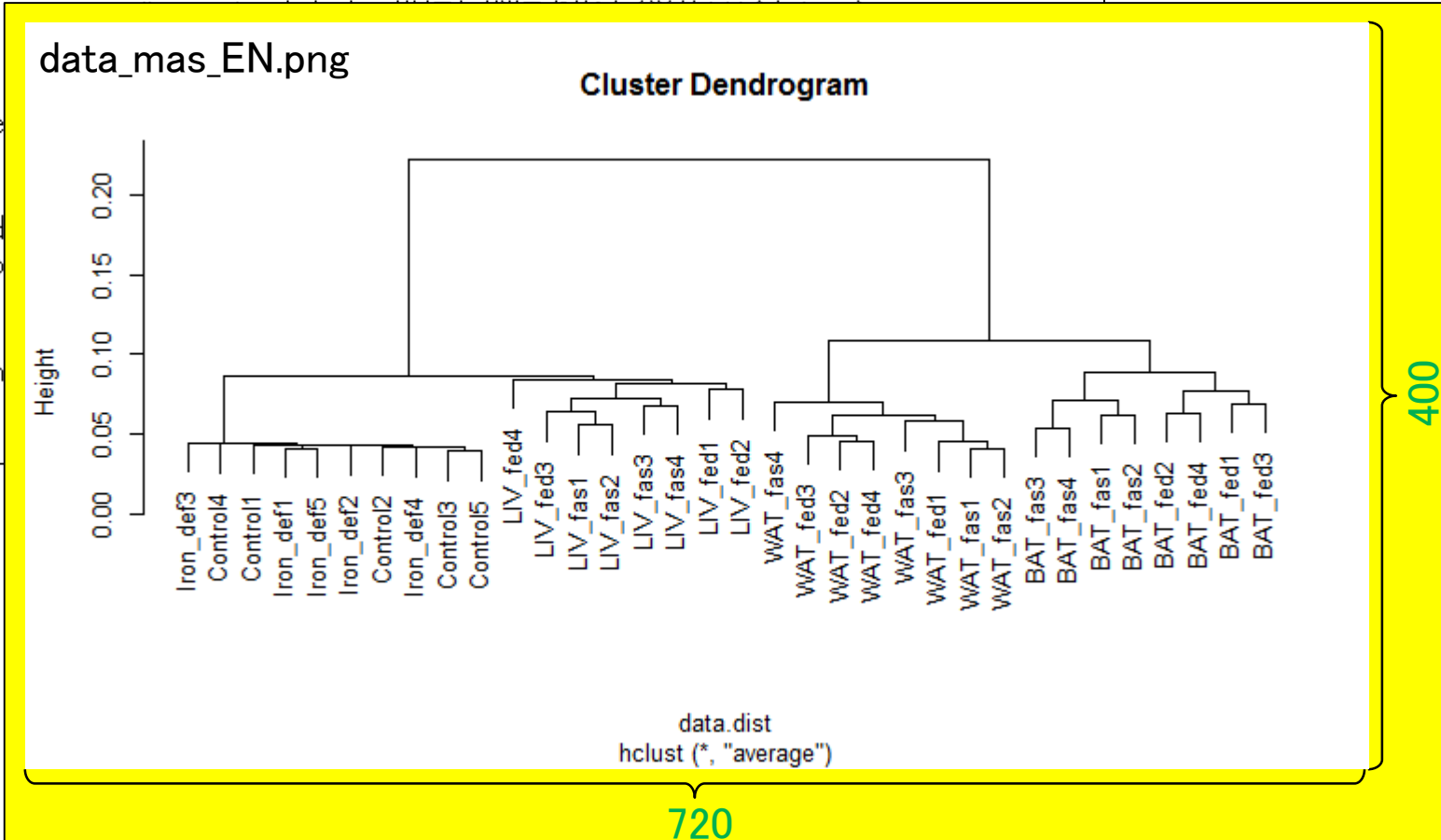
```
#####↓
### MAS5データのクラスタリング, 距離: 1-Spearman相関係数, 方法: 平均連結法 ###↓
#####↓
in_f <- "data_mas_JP.txt" #入力ファイル名を指定してin_fに格納↓
out_f <- "data_mas_JP.png" #出力ファイル名を指定してout_fに格納↓
param <- "average" #方法(method)を指定↓
param_fig <- c(720, 400)
↓
#入力ファイルの読み込み↓
data <- read.table(in_f, header=TRUE)
↓
#本番↓
data.dist <- as.dist(1 - cor(data))
out <- hclust(data.dist, method="average")
↓
#ファイルに保存↓
png(out_f, pointsize=13, width=720, height=400)
plot(out)
dev.off()
```



クラスタリング結果をファイルに保存

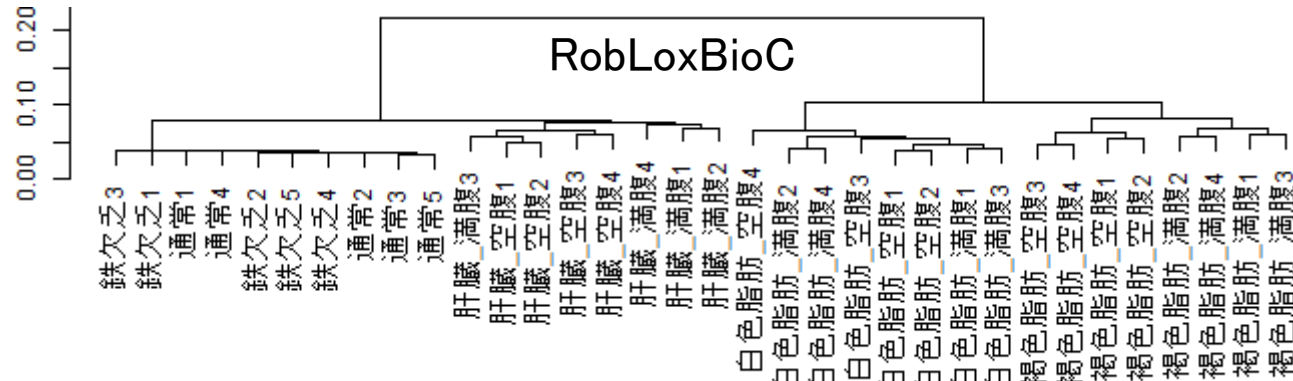
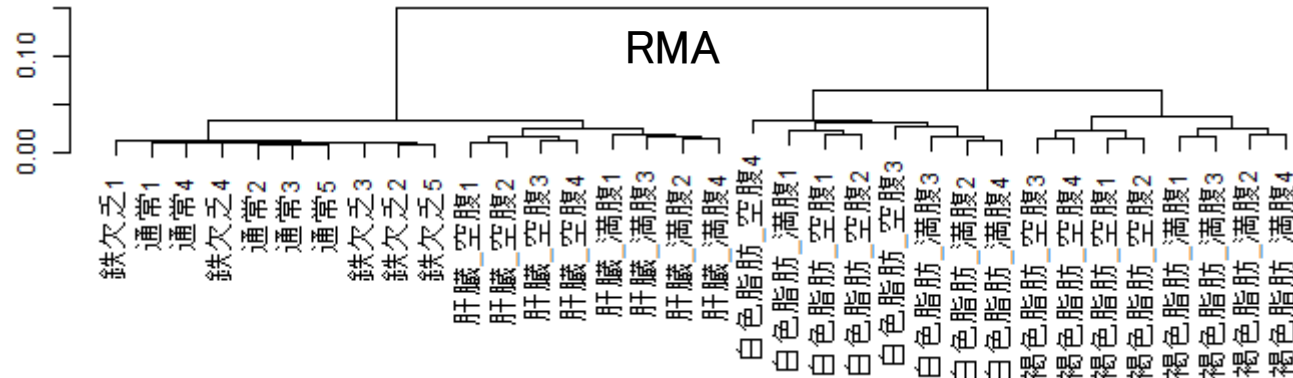
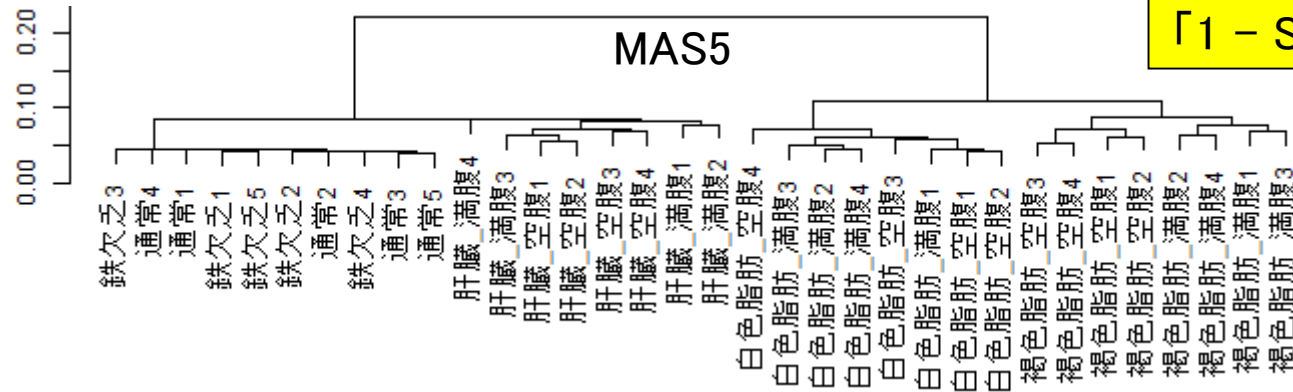
rcode_clustering_png.txt

```
#####↓
### MAS5データのクラスタリング, 距離: 1-Spearman相関係数, 方法: 平均連結法 ###↓
#####↓
in_f <- "data_mas_EN.txt"           #入力ファイル名を指定してin_fに格納↓
out_f <- "data_mas_EN.png"         #出力ファイル名を指定してout_fに格納↓
param <- "average"                 #方法(method)を指定↓
param_fig <- c(720, 400)
↓
#入力ファイルの読み込み↓
data <- read.table(in_f, header=)
↓
#本番↓
data.dist <- as.dist(1 - cor(d
out <- hclust(data.dist, method
↓
#ファイルに保存↓
png(out_f, pointsize=13, width
plot(out)
dev.off()
```



サンプル間クラスタリング (マージ後)

「1 - Spearman相関係数」の結果



サンプル間クラスタリング（マージ後）

0.20
0

MAS5

「1 - Spearman相関係数」の結果

```

param1 <- "average" #方法(method)を指定↓
param2 <- "spearman" #相関係数の種類を指定↓
param_fig <- c(720, 310) #ファイル出力時の横幅と縦幅を指定(単位はピクセル)↓
↓
#####↓
### MAS5データのクラスタリング, 距離: 1-Spearman相関係数, 方法: 平均連結法 ###↓
#####↓
in_f <- "data_mas_JP.txt" #入力ファイル名を指定してin_fに格納↓
out_f <- "data_mas_JP.png" #出力ファイル名を指定してout_fに格納↓
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
data.dist <- as.dist(1 - cor(data, method=param2))#サンプル間の距離を計算した結果をdata.distに格納↓
out <- hclust(data.dist, method=param1) #階層的クラスタリングを実行した結果をoutに格納↓
png(out_f, pointsize=12, width=param_fig[1], height=param_fig[2])#出力ファイルの各種パラメータを指定↓
plot(out) #樹形図 (デンドログラム) の表示↓
dev.off() #おまじない↓
↓
#####↓
### RMAデータのクラスタリング, 距離: 1-Spearman相関係数, 方法: 平均連結法 ###↓
#####↓
in_f <- "data_rma_JP.txt" #入力ファイル名を指定してin_fに格納↓
out_f <- "data_rma_JP.png" #出力ファイル名を指定してout_fに格納↓
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
data.dist <- as.dist(1 - cor(data, method=param2))#サンプル間の距離を計算した結果をdata.distに格納↓
out <- hclust(data.dist, method=param1) #階層的クラスタリングを実行した結果をoutに格納↓

```

サンプル間クラスタリング（マージ後）

「1 - Pearson相関係数」の結果にしたいとき

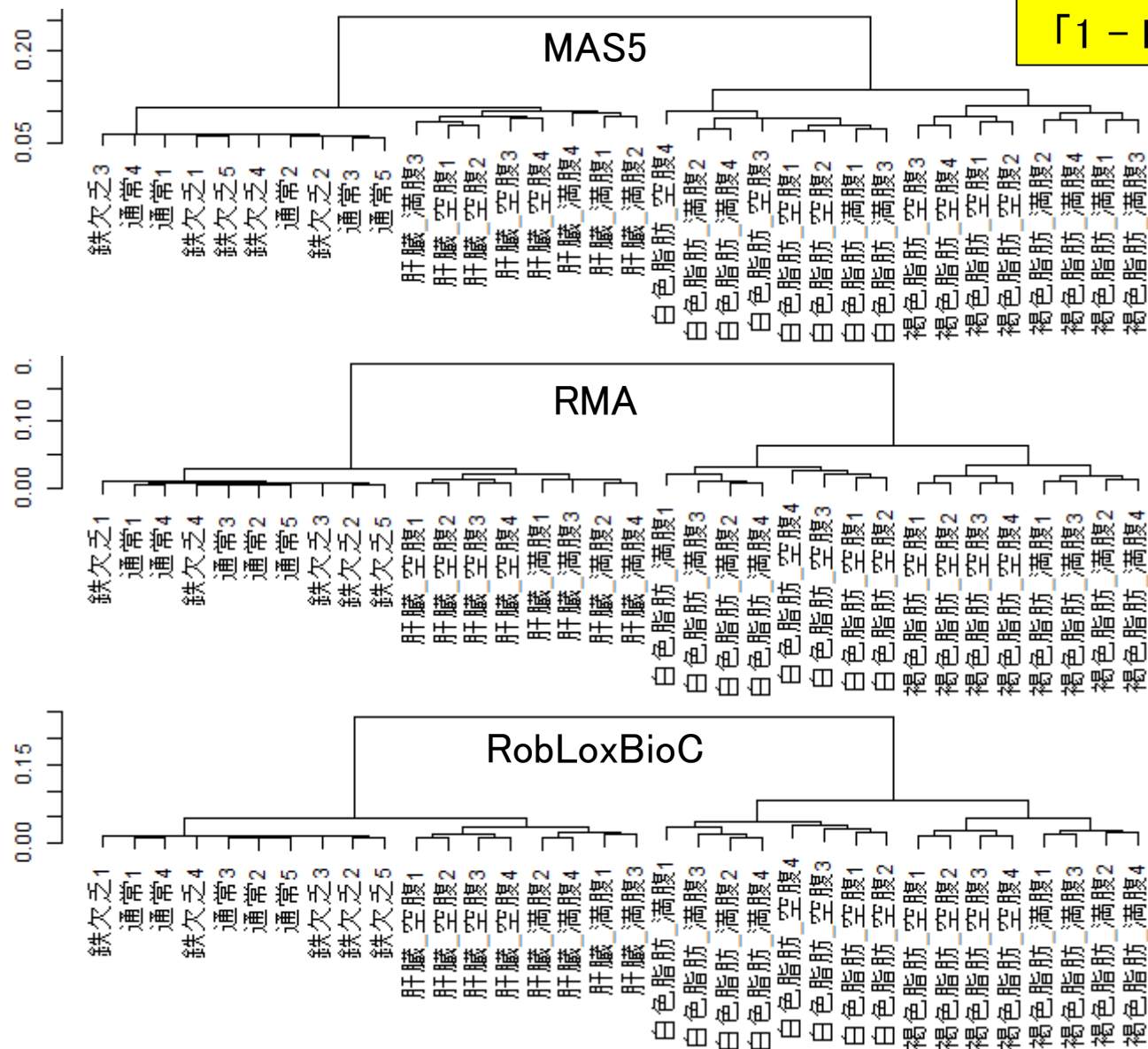
```

param1 <- "average" #方法(method)を指定↓
param2 <- "pearson" #相関係数の種類を指定↓
param_fig <- c(720, 310) #ファイル出力時の横幅と縦幅を指定(単位はピクセル)↓
↓
#####↓
### MAS5データのクラスタリング, 距離: 1-Spearman相関係数, 方法: 平均連結法 ###↓
#####↓
in_f <- "data_mas_JP.txt" #入力ファイル名を指定してin_fに格納↓
out_f <- "data_mas_JP.png" #出力ファイル名を指定してout_fに格納↓
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
data.dist <- as.dist(1 - cor(data, method=param2))#サンプル間の距離を計算した結果をdata.distに格納↓
out <- hclust(data.dist, method=param1) #階層的クラスタリングを実行した結果をoutに格納↓
png(out_f, pointsize=12, width=param_fig[1], height=param_fig[2])#出力ファイルの各種パラメータを指定↓
plot(out) #樹形図 (デンドログラム) の表示↓
dev.off() #おまじない↓
↓
#####↓
### RMAデータのクラスタリング, 距離: 1-Spearman相関係数, 方法: 平均連結法 ###↓
#####↓
in_f <- "data_rma_JP.txt" #入力ファイル名を指定してin_fに格納↓
out_f <- "data_rma_JP.png" #出力ファイル名を指定してout_fに格納↓
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
data.dist <- as.dist(1 - cor(data, method=param2))#サンプル間の距離を計算した結果をdata.distに格納↓
out <- hclust(data.dist, method=param1) #階層的クラスタリングを実行した結果をoutに格納↓

```

サンプル間クラスタリング (マージ後)

「1 - Pearson相関係数」の結果



サンプル間クラスタリング（マージ後）

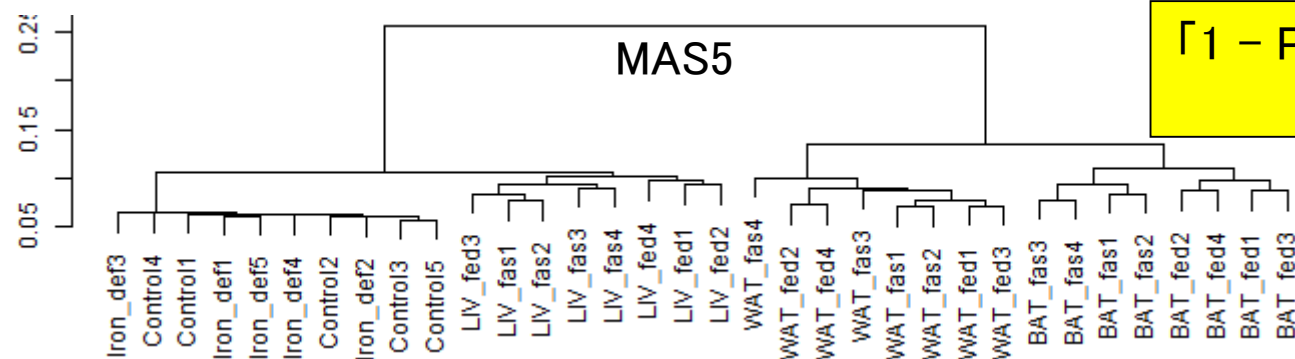
Japanese → English

```

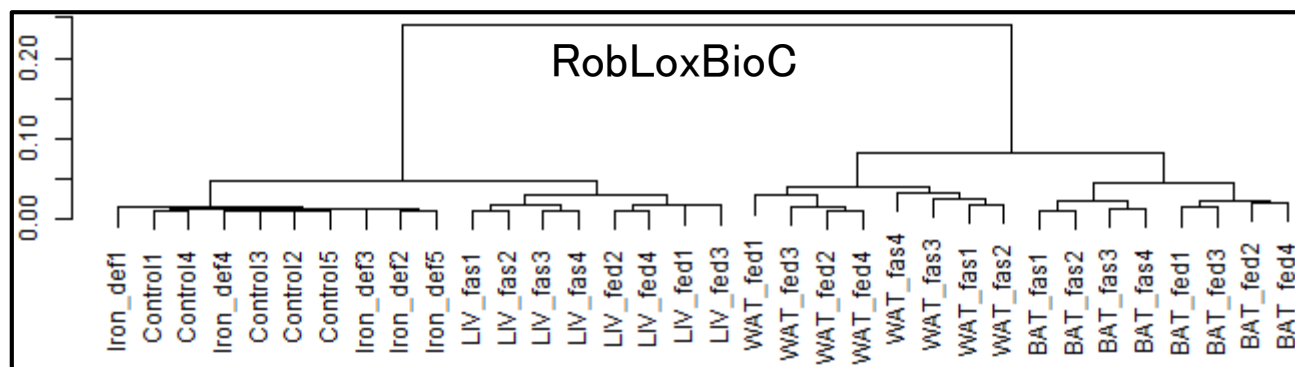
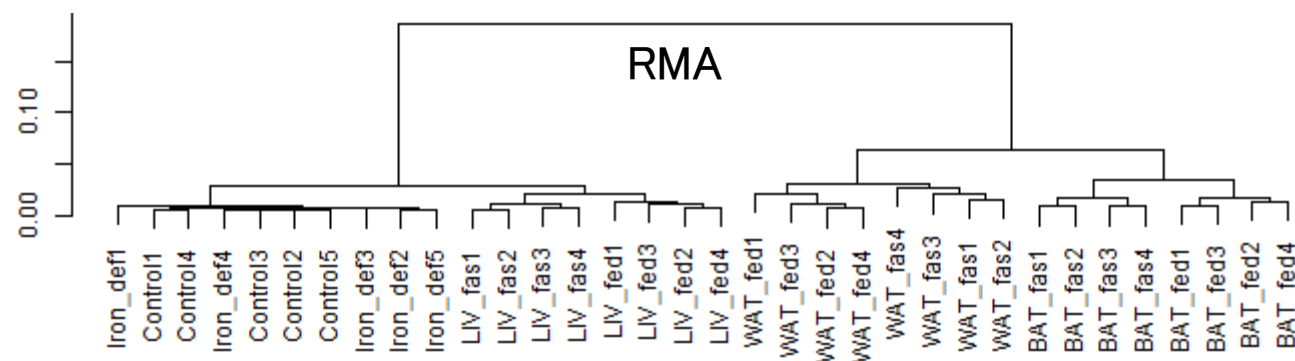
param1 <- "average"          #方法(method)を指定↓
param2 <- "pearson"         #相関係数の種類を指定↓
param_fig <- c(720, 310)    #ファイル出力時の横幅と縦幅を指定(単位はピクセル)↓
↓
#####↓
### MAS5データのクラスタリング, 距離: 1-Spearman相関係数, 方法: 平均連結法 ###↓
#####↓
in_f <- "data_mas_EN.txt"   #入力ファイル名を指定してin_fに格納↓
out_f <- "data_mas_JP.png" #出力ファイル名を指定してout_fに格納↓
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
data.dist <- as.dist(1 - cor(data, method=param2))#サンプル間の距離を計算した結果をdata.distに格納↓
out <- hclust(data.dist, method=param1) #階層的クラスタリングを実行した結果をoutに格納↓
png(out_f, pointsize=12, width=param_fig[1], height=param_fig[2])#出力ファイルの各種パラメータを指定↓
plot(out)                   #樹形図 (デンドログラム) の表示↓
dev.off()                   #おまじない↓
↓
#####↓
### RMAデータのクラスタリング, 距離: 1-Spearman相関係数, 方法: 平均連結法 ###↓
#####↓
in_f <- "data_rma_EN.txt"   #入力ファイル名を指定してin_fに格納↓
out_f <- "data_rma_JP.png" #出力ファイル名を指定してout_fに格納↓
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
data.dist <- as.dist(1 - cor(data, method=param2))#サンプル間の距離を計算した結果をdata.distに格納↓
out <- hclust(data.dist, method=param1) #階層的クラスタリングを実行した結果をoutに格納↓

```

サンプル間クラスタリング (マージ後)

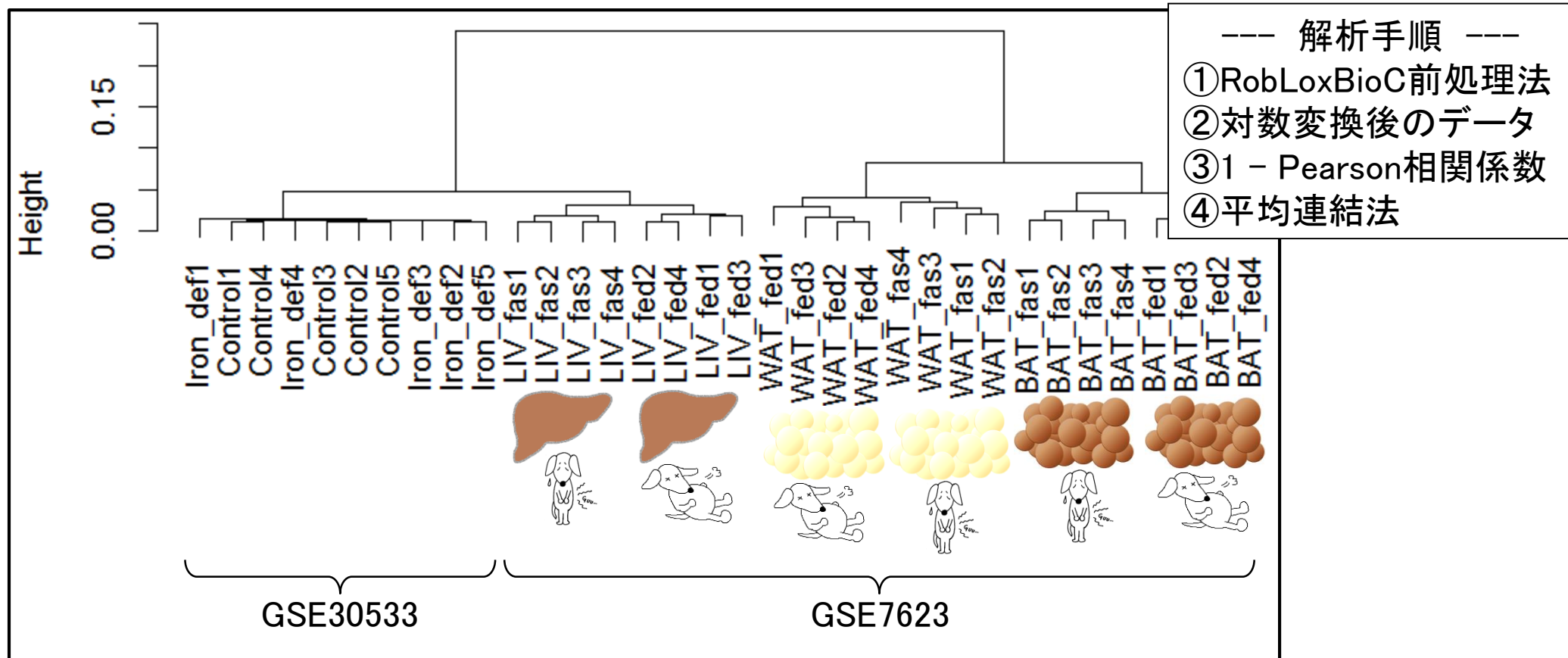


「1 - Pearson相関係数」の結果
English version



課題(クラスタリング結果の解釈)

■ ラット(24サンプル+10サンプル)のクラスタリング結果について考察せよ。

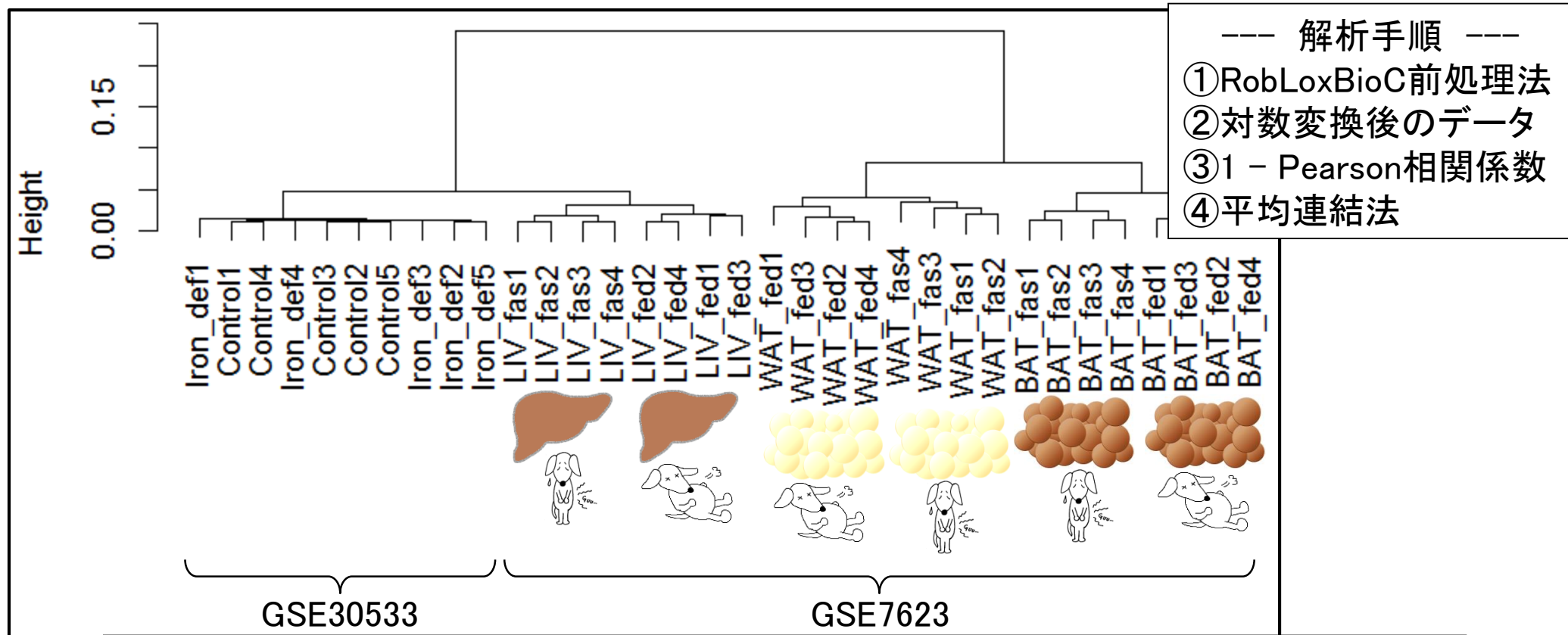


Q1:なぜGSE7623とGSE30533のデータはマージ可能か?

Q2: GSE30533の10サンプルからなるクラスターは、GSE7623の3種類の組織(LIV, WAT, BAT)のどの発現パターンに近いか?

課題(クラスタリング結果の解釈)

■ ラット(24サンプル+10サンプル)のクラスタリング結果について考察せよ。



Q3: GSE30533のみのクラスタリング結果は「鉄欠乏 (Iron_def) 状態と通常 (Control) 状態」が入り混じっている。その一方で、「満腹 (fed) 状態と空腹 (fas) 状態」の違いは3種類の組織 (LIV, WAT, BAT) いずれにおいても明瞭に分かれている。鉄欠乏 (Iron_def) 状態と空腹 (fas) 状態の発現プロファイル変化への影響度はどちらか大きいと思われるか？

実験デザイン (§ 3.2.2)

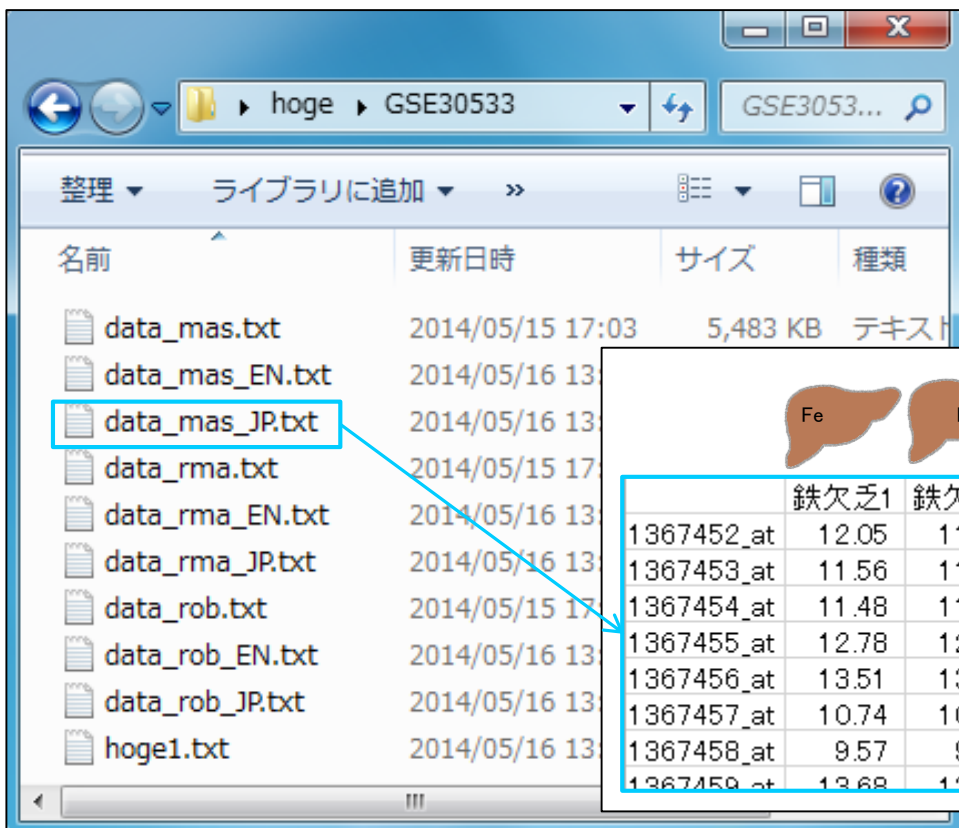
Affymetrix GeneChip

- Ge et al., *Genomics*, **86**: 127–141, 2005
 - GSE2361、GPL96 (Affymetrix Human Genome U133A Array)、22,283 probesets
 - ヒト36サンプル: Heart (心臓)、Thymus (胸腺)、Spleen (脾臓)、Ovary (卵巣)、Kidney (腎臓)、Skeletal Muscle (骨格筋)、Pancreas (膵臓)、Prostate (前立腺)、…
- Nakai et al., *Biosci Biotechnol Biochem.*, **72**: 139–148, 2008 8匹のラットを使用
 - GSE7623、GPL1355 (Affymetrix Rat Genome 230 2.0 Array)、31,099 probesets
 - ラット24サンプル: Brown adipose tissue (褐色脂肪組織; BAT) 8サンプル、White adipose tissue (白色脂肪組織; WAT) 8サンプル、Liver (肝臓; LIV) 8サンプル
 - BAT 8サンプル: 通常 (BAT_fed) 4サンプル 対 24時間絶食 (BAT_fas) 4サンプル
 - WAT 8サンプル: 通常 (WAT_fed) 4サンプル 対 24時間絶食 (WAT_fas) 4サンプル
 - LIV 8サンプル: 通常 (LIV_fed) 4サンプル 対 24時間絶食 (LIV_fas) 4サンプル
- Kamei et al., *PLoS One*, **8**: e65732, 2013 10匹のラットを使用
 - GSE30533、GPL1355 (Affymetrix Rat Genome 230 2.0 Array)、31,099 probesets
 - ラット10サンプル: 全てLiver (肝臓) サンプル
 - iron-deficient diet (Iron_def) 5サンプル 対 control diet (Control) 5サンプル

実験デザイン (§ 3.2.2)

■ Kamei et al., *PLoS One*, 8: e65732, 2013

- GSE30533、GPL1355 (Affymetrix Rat Genome 230 2.0 Array)、31,099 probesets
- ラット10サンプル: 全てLiver (肝臓) サンプル
- iron-deficient diet (Iron_def) 5サンプル 対 control diet (Control) 5サンプル



2群間比較が主な目的であり、各群につき5反復 (five replicates) としている。生物学的なばらつき (biological variation) を考慮すべく、反復データは別々の個体からとっている (biological replicates)

	鉄欠乏1	鉄欠乏2	鉄欠乏3	鉄欠乏4	鉄欠乏5	通常1	通常2	通常3	通常4	通常5
1367452_at	12.05	11.92	11.99	11.92	11.73	12.08	12.06	11.98	12.03	12.03
1367453_at	11.56	11.59	11.62	11.75	11.78	11.63	11.51	11.48	11.57	11.68
1367454_at	11.48	11.68	11.61	11.65	11.86	11.71	11.98	12.01	11.59	11.95
1367455_at	12.78	12.59	12.70	12.79	13.00	12.68	12.78	12.55	12.68	12.87
1367456_at	13.51	13.53	13.48	13.52	13.45	13.47	13.59	13.60	13.52	13.57
1367457_at	10.74	10.14	10.61	10.26	10.31	10.50	10.30	10.43	10.39	10.52
1367458_at	9.57	9.17	9.15	8.95	9.41	9.25	8.79	9.14	9.37	9.22
1367459_at	13.68	13.56	13.63	13.57	13.77	13.69	13.61	13.55	13.59	13.69

実験デザイン (§ 3.2.2)

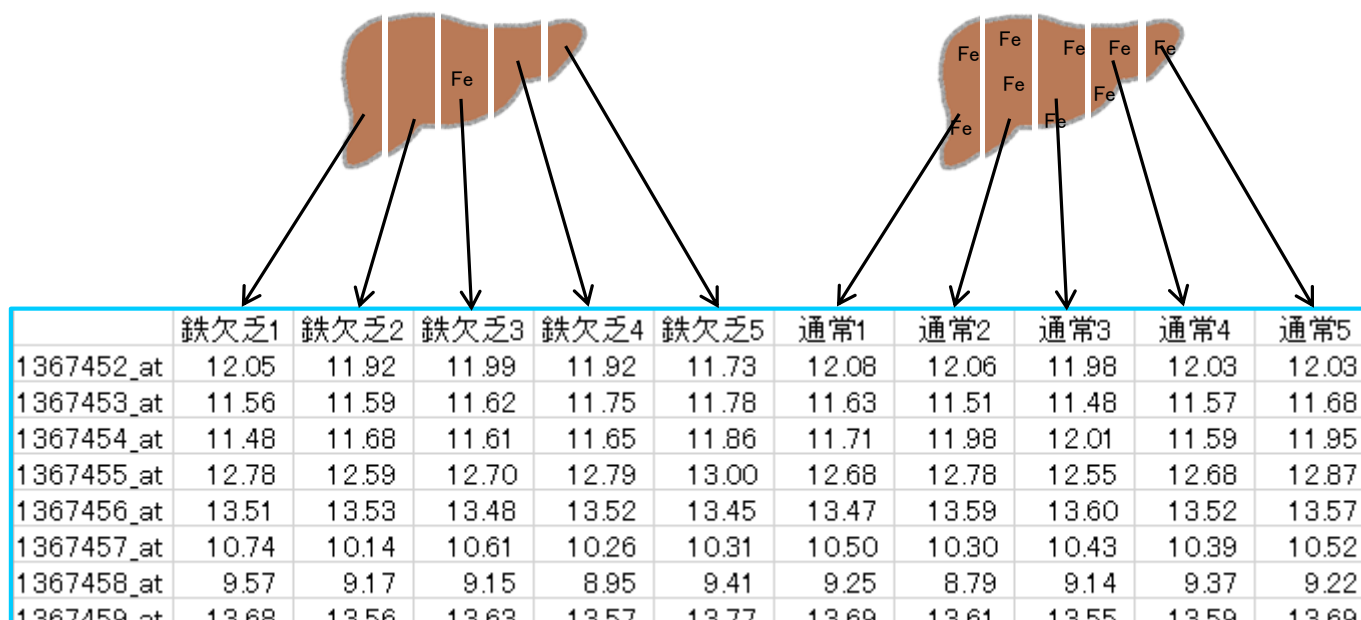
Kamei et al., *PLoS One* 8: e65722, 2012

- GSE30533、
- ラット10サン
- iron-deficient diet (iron_def)

対比的な用語は技術的なばらつき (technical variation) であり、同一個体由来サンプルを分割して得られた反復データ (technical replicates)

31,099 probesets

control diet (control) 5サンプル

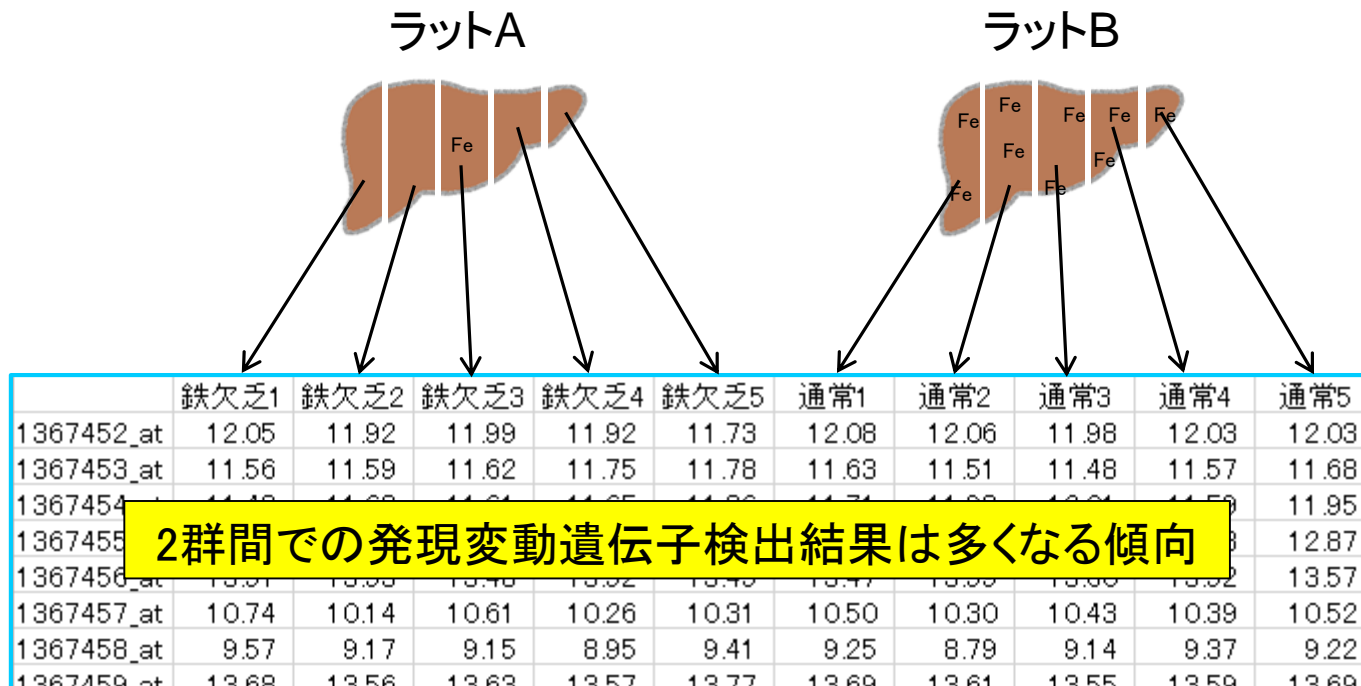


このやり方で得られる結論は限定的!できるだけ多様な別個体サンプルを沢山用いるべし!

実験デザイン (§ 3.2.2)

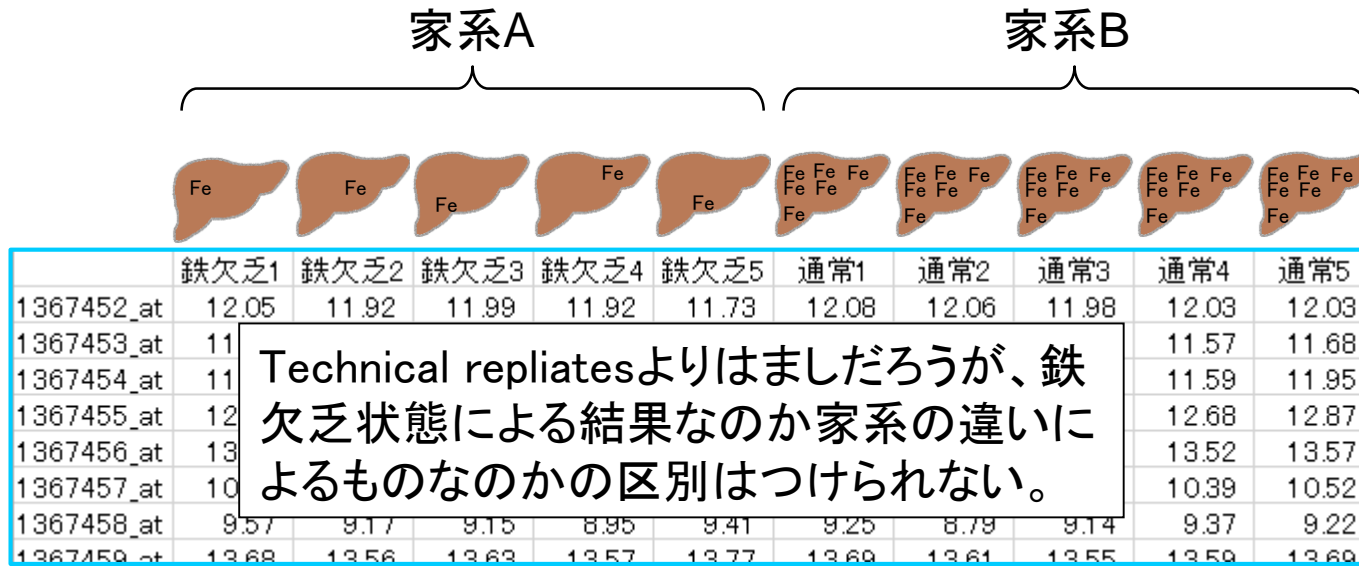
■ Technical replicatesだと...

1. 自分は「鉄欠乏 対 通常」の違いを見ているつもりでも、個体間の他の違い(身長、体重など)由来要因との区別がつかない
高身長 対 低身長、低体重 対 高体重、他の病気の有無、家系の違いなど
2. 得られる結果から導き出される結論は、そのラット間のみで成立する事象であり、ラットという生物種全体に適用可能なわけではない



実験デザイン(§ 3.2.2)

- Biological replicatesでも多様性が不十分な場合はイマイチ…



**普遍的な結果を得たいのなら、できるだけ多様な別個体サンプルを沢山用いるべし!
Expression Atlasも3 biological replicates以上を基本としているようだ。**

まとめ

- プロブレベルデータ取得および前処理法適用(発展形)
 - Rパッケージ経由でダウンロードしたデータをそのまま前処理法にかける
 - サンプル名の順番が変わることもあるので注意
 - 何度も同じ作業をしなくていいように有効活用
 - 3つの論文の遺伝子発現行列をコピーで作成
 - ヒトの36サンプルからなるデータ (Ge et al., 2005)
 - ラットの24サンプルからなるデータ (Nakai et al., 2008)
 - ラットの10サンプルからなるデータ (Kamei et al., 2013)
 - サンプル間クラスタリング
 - 階層的クラスタリング、距離の定義、クラスターをまとめる方法など
 - 前処理法 (MAS5, RMA, RobLoxBioC) と相関係数 (Spearman and Pearson) の違い
 - 同一アレイ由来データセットのマージ (ラット24サンプル+ラット10サンプル)
 - 3' 発現アレイの長所 (教科書p7)
 - ラット10サンプル (通常 対 鉄欠乏) クラスタリング結果の印象は、外群 (ラット24サンプル) の有無でずいぶん異なる (教科書p106-107)
 - 実験デザインは重要
- 6/4の組織特異的遺伝子検出で利用予定
- 5/28の発現変動遺伝子検出で利用予定

クラスタリング結果とDEG数の関係

クラスタリング結果を眺めれば、発現変動遺伝子 (DEG) 数に関するおおよその見当がつかます
→ クラスタリングって重要

