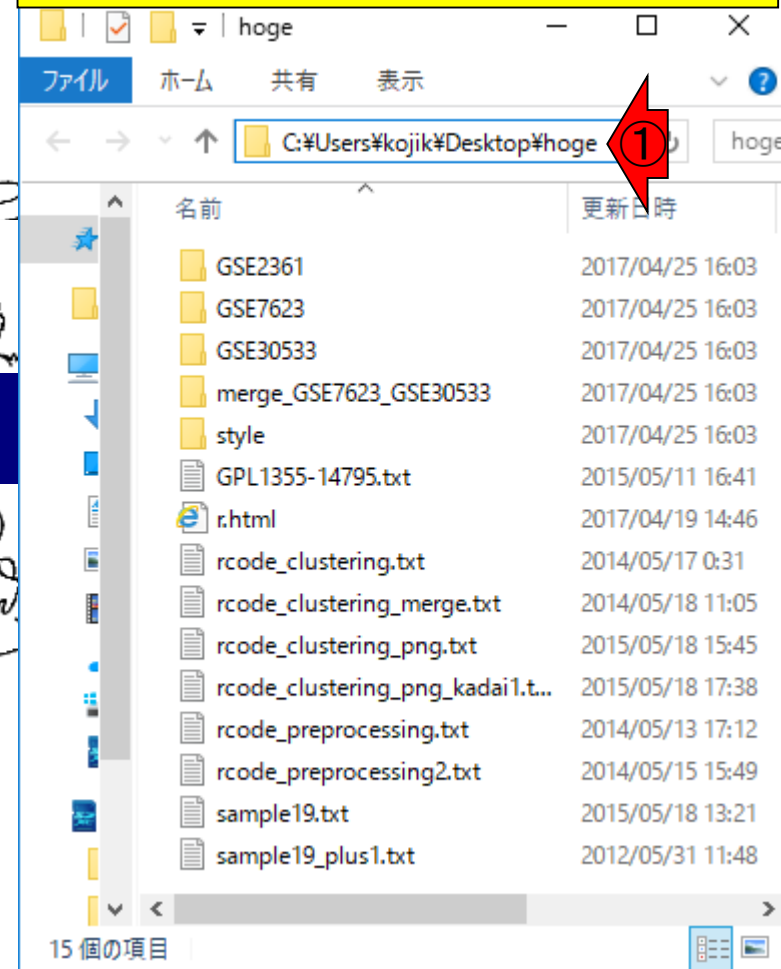


①持ち込みPCのヒトは後ろにあるUSBメモリ中のhogeフォルダをデスクトップにコピーしておいてください



機能ゲノム学第2回

大学院農学生命科学研究科
アグリバイオインフォマティクス教育研究プログラム

門田幸二(かどた こうじ)

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

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

Contents

- 実データ概観
 - GSE2361 (ヒト)、GSE7623 (ラット)、GSE30533 (ラット)
- クラスタリング (教科書の § 3.2.1)
 - 対数変換の有無 (Spearman相関係数を使う場合は気にしなくてよい)
 - 階層的 vs. 非階層的、様々な選択肢
 - 距離の定義: ベクトル間、クラスター間
 - 実データで実行: GSE2361 (ヒト)、課題1
 - 実データで実行: GSE7623 (ラット)、GSE30533 (ラット)
 - 同一プラットフォームデータ(GSE7623 + GSE30533)をマージして実行、課題2
- 実験デザイン (教科書の § 3.2.2)
- 2群間比較: 発現変動遺伝子 (DEG) 検出
 - パターンマッチング法 (相関係数の利用)
 - コードの中身をおさらい、apply関数の基本的な利用法など

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)

実データ概観

■ 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

実データ概観

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

□ 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: 3T3-L1
 - WAT 8サンプル: 3T3-L1
 - 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 (Iron

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

GSE2361(ヒト)

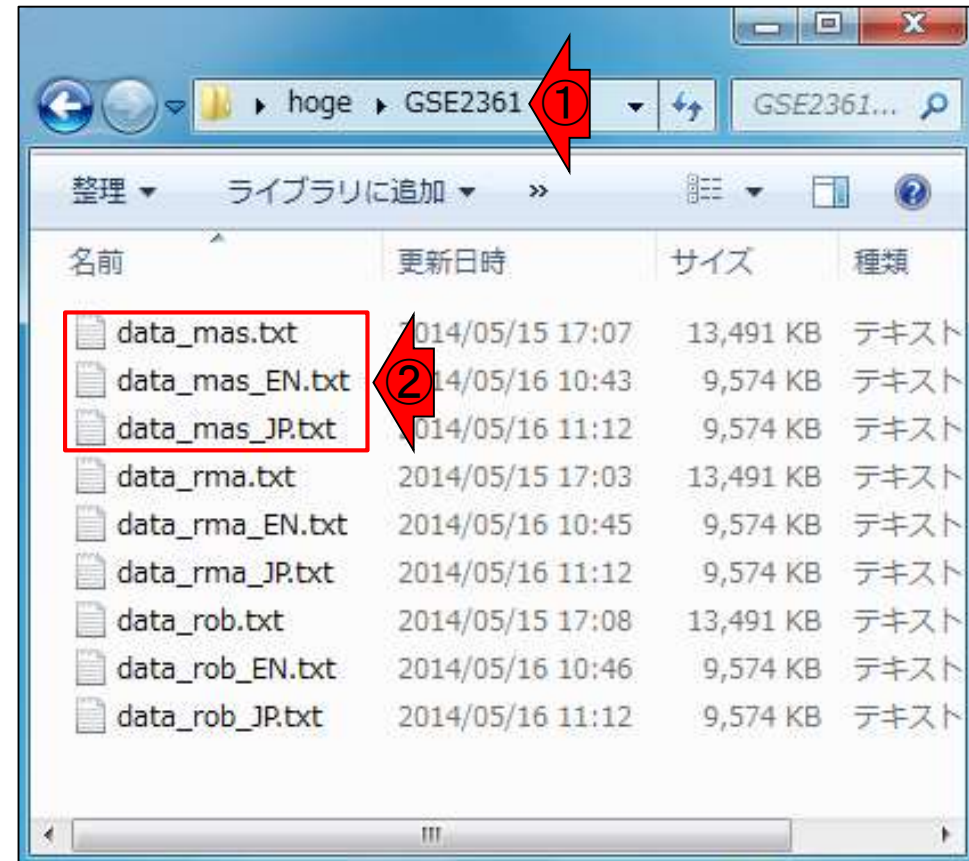
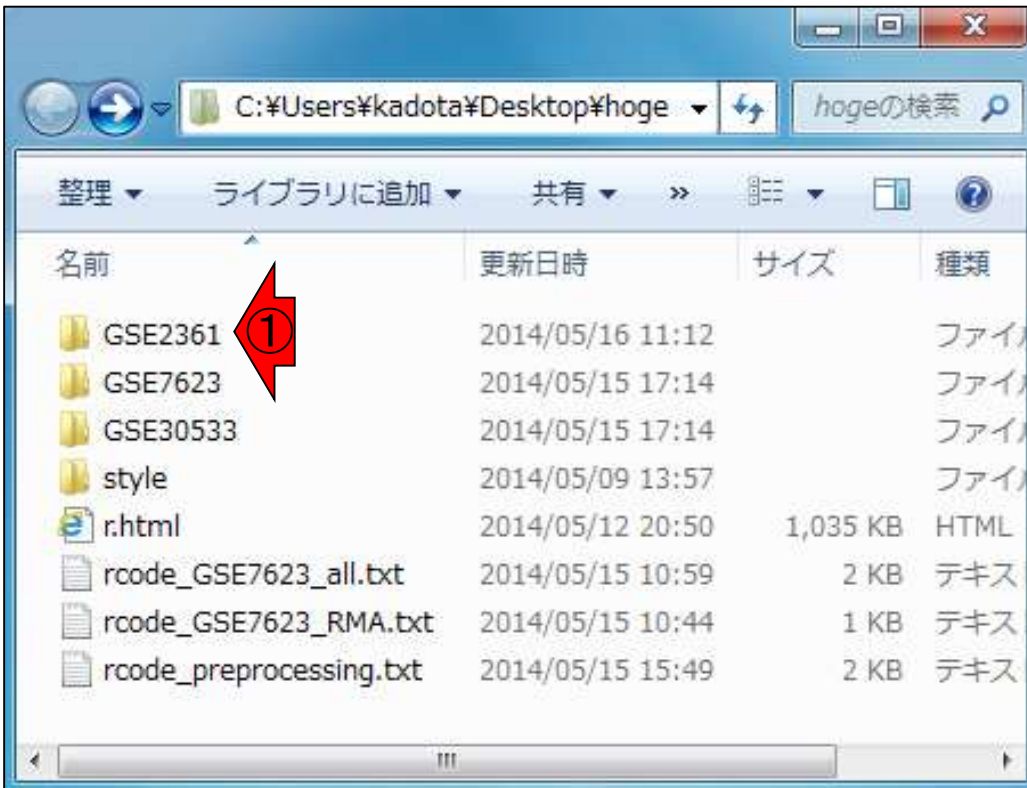
Affymetrix GeneChip

□ 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)、…

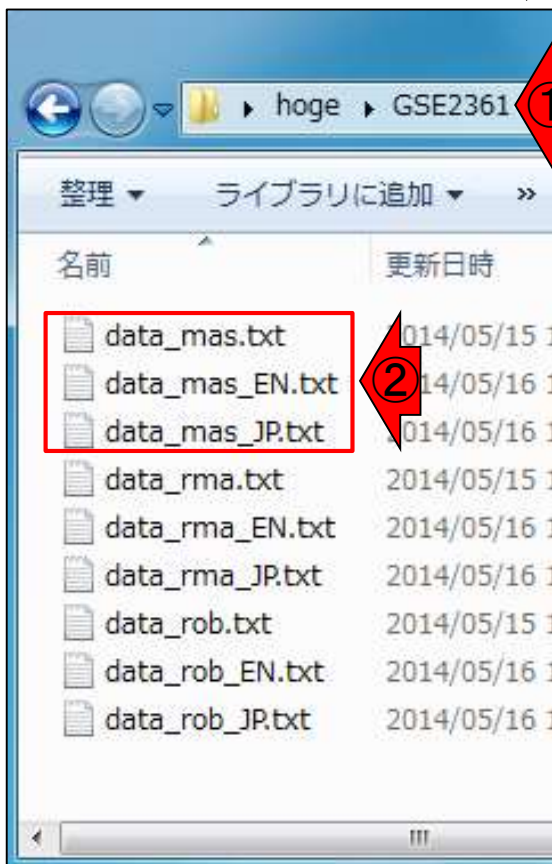
①GSE2361フォルダの中身。②

data_mas.txtは前処理法をそのまま適用した結果。*_EN.txtはサンプル名を英語で、*_JP.txtは日本語で書き換えたもの



GSE2361(ヒト)

*_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。好きなものをご利用ください。いずれも対数変換後のデータ(log-transformed data)



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	11.48001
3	1053_at	6.621657	6.47488	7.371465	4.168622	7.350244	7.209918	8.362777	7.176571	7.75001
4	117_at	8.259603	8.647364	8.689724	7.875649	5.083001	8.165044	7.96707	8.418082	5.71901
5	121_at	11.26699	11.04186	11.39999	11.02855	13.13267	11.39138	12.32899	11.0559	11.28001
6	1255_g_at	7.1757	6.477278	6.781766	7.048799	7.30767	6.600267	6.42359	6.694412	7.30001
7	1294_at	9.137586	9.718507	9.083742	9.014997	8.813377	8.402562	9.146946	8.421351	10.12001
8	1316_at	8.235789	7.418895	8.07469	7.280095	8.238176	7.600147	7.422269	7.288894	7.67001

data_mas_EN.txt

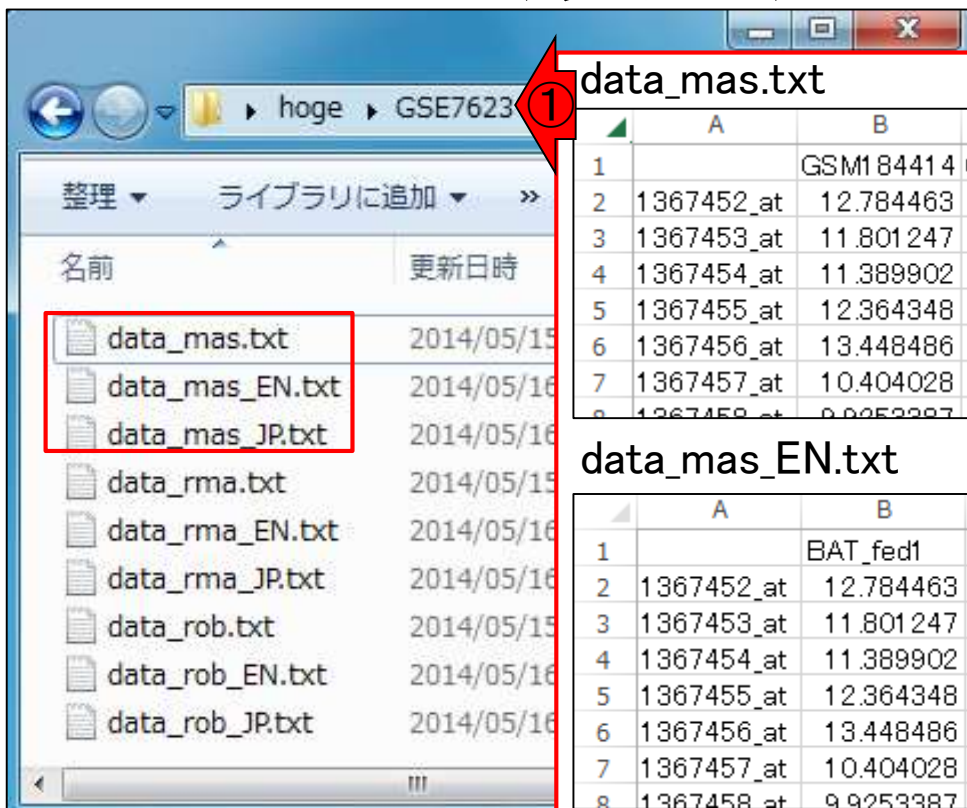
	A	B	C	D	E	F	G	H	I	J
1		Heart	Thymus	Spleen	Ovary	Kidney	Skeletal_Mu	Pancreas	Prostate	Small_I
2	1007_s_at	10.82435	11.21261	9.898072	10.8948	11.75531	10.92032	11.61243	11.22101	11.48001
3	1053_at	6.621657	6.47488	7.371465	4.168622	7.350244	7.209918	8.362777	7.176571	7.75001
4	117_at	8.259603	8.647364	8.689724	7.875649	5.083001	8.165044	7.96707	8.418082	5.71901
5	121_at	11.26699	11.04186	11.39999	11.02855	13.13267	11.39138	12.32899	11.0559	11.28001
6	1255_g_at	7.1757	6.477278	6.781766	7.048799	7.30767	6.600267	6.42359	6.694412	7.30001
7	1294_at	9.137586	9.718507	9.083742	9.014997	8.813377	8.402562	9.146946	8.421351	10.12001
8	1316_at	8.235789	7.418895	8.07469	7.280095	8.238176	7.600147	7.422269	7.288894	7.67001

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	11.48001
3	1053_at	6.621657	6.47488	7.371465	4.168622	7.350244	7.209918	8.362777	7.176571	7.75001
4	117_at	8.259603	8.647364	8.689724	7.875649	5.083001	8.165044	7.96707	8.418082	5.71901
5	121_at	11.26699	11.04186	11.39999	11.02855	13.13267	11.39138	12.32899	11.0559	11.28001
6	1255_g_at	7.1757	6.477278	6.781766	7.048799	7.30767	6.600267	6.42359	6.694412	7.30001
7	1294_at	9.137586	9.718507	9.083742	9.014997	8.813377	8.402562	9.146946	8.421351	10.12001
8	1316_at	8.235789	7.418895	8.07469	7.280095	8.238176	7.600147	7.422269	7.288894	7.67001

GSE7623(ラット)

①GSE7623 (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	12.439875
3	1367453_at	11.801247	12.152935	11.942227	11.968477	11.845375	11.681727	12.078672	12.048875
4	1367454_at	11.389902	11.160757	11.145987	11.212088	11.540652	11.308877	11.49885	11.402875
5	1367455_at	12.364348	12.529744	12.432574	12.604011	12.441991	12.249935	12.281827	12.190875
6	1367456_at	13.448486	13.543046	13.552794	13.629799	13.36913	13.244278	13.424371	13.329875
7	1367457_at	10.404028	10.69632	10.475078	10.45579	10.141921	10.290666	10.146529	10.260875
8	1367458_at	9.9253387	10.244544	9.9720008	9.9576072	8.702884	9.3578792	9.2134367	9.499875

data_mas_EN.txt

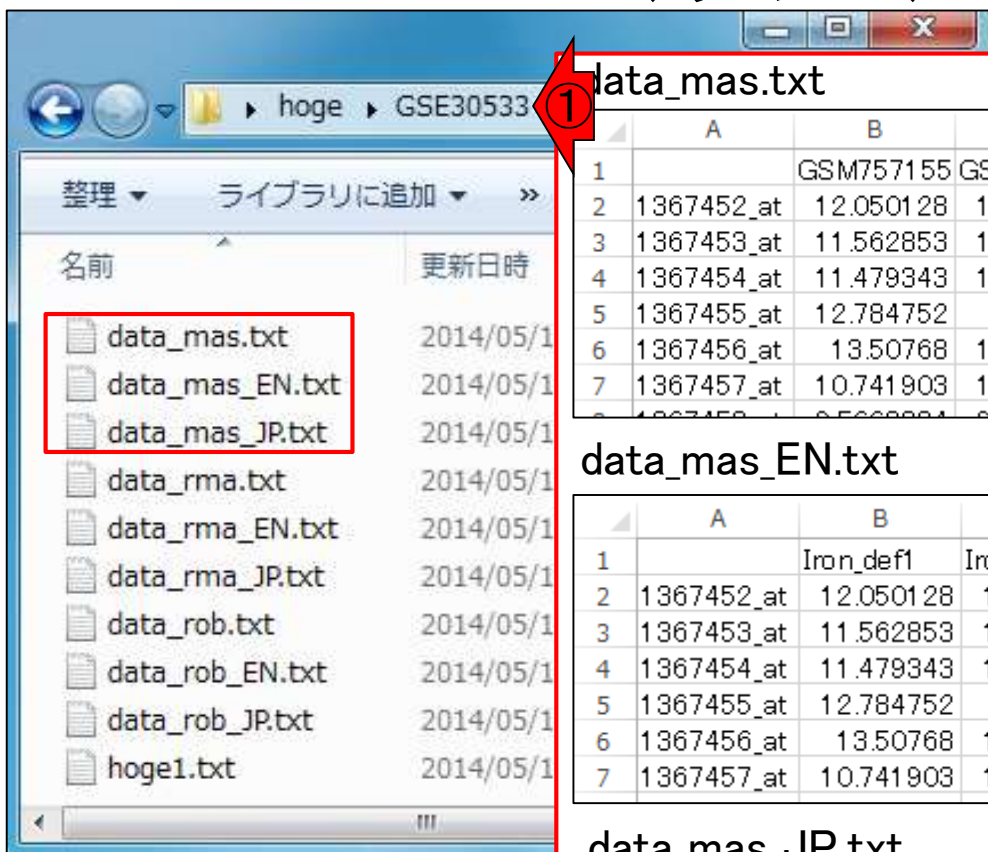
	A	B	C	D	E	F	G	H	I
1		BAT_fed1	BAT_fed2	BAT_fed3	BAT_fed4	BAT_fas1	BAT_fas2	BAT_fas3	BAT_fas4
2	1367452_at	12.784463	12.447082	12.805908	12.304718	12.589425	12.607532	11.815378	12.439875
3	1367453_at	11.801247	12.152935	11.942227	11.968477	11.845375	11.681727	12.078672	12.048875
4	1367454_at	11.389902	11.160757	11.145987	11.212088	11.540652	11.308877	11.49885	11.402875
5	1367455_at	12.364348	12.529744	12.432574	12.604011	12.441991	12.249935	12.281827	12.190875
6	1367456_at	13.448486	13.543046	13.552794	13.629799	13.36913	13.244278	13.424371	13.329875
7	1367457_at	10.404028	10.69632	10.475078	10.45579	10.141921	10.290666	10.146529	10.260875
8	1367458_at	9.9253387	10.244544	9.9720008	9.9576072	8.702884	9.3578792	9.2134367	9.499875

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	9.92533749	10.24454359	9.97200015	9.957607169	8.70288404	9.35787919

①GSE30533 (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

GSE30533 (ラット)



	A	B	C	D	E	F	G	H	I
1		GSM757155	GSM757156	GSM757157	GSM757158	GSM757159	GSM757160	GSM757161	GSM7571
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	11.62175	11.747779	11.776342	11.633797	11.50868	11.4794
4	1367454_at	11.479343	11.676545	11.608875	11.652704	11.86008	11.70589	11.975747	12.0095
5	1367455_at	12.784752	12.58787	12.696588	12.789029	13.002298	12.678645	12.780148	12.5522
6	1367456_at	13.50768	13.533852	13.483843	13.524296	13.452217	13.472369	13.594191	13.603
7	1367457_at	10.741903	10.136369	10.609696	10.256467	10.312411	10.49749	10.304808	10.4274

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	11.979
3	1367453_at	11.562853	11.587924	11.62175	11.747779	11.776342	11.633797	11.50868	11.4794
4	1367454_at	11.479343	11.676545	11.608875	11.652704	11.86008	11.70589	11.975747	12.0095
5	1367455_at	12.784752	12.58787	12.696588	12.789029	13.002298	12.678645	12.780148	12.5522
6	1367456_at	13.50768	13.533852	13.483843	13.524296	13.452217	13.472369	13.594191	13.603
7	1367457_at	10.741903	10.136369	10.609696	10.256467	10.312411	10.49749	10.304808	10.4274

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	11.979
3	1367453_at	11.562853	11.587924	11.62175	11.747779	11.776342	11.633797	11.50868	11.4794
4	1367454_at	11.479343	11.676545	11.608875	11.652704	11.86008	11.70589	11.975747	12.0095
5	1367455_at	12.784752	12.58787	12.696588	12.789029	13.002298	12.678645	12.780148	12.5522
6	1367456_at	13.50768	13.533852	13.483843	13.524296	13.452217	13.472369	13.594191	13.603
7	1367457_at	10.741903	10.136369	10.609696	10.256467	10.312411	10.49749	10.304808	10.4274

Contents

- 実データ概観
 - GSE2361 (ヒト)、GSE7623 (ラット)、GSE30533 (ラット)
- クラスタリング (教科書の § 3.2.1)
 - 対数変換の有無 (Spearman相関係数を使う場合は気にしなくてよい)
 - 階層的 vs. 非階層的、様々な選択肢
 - 距離の定義: ベクトル間、クラスター間
 - 実データで実行: GSE2361 (ヒト)、課題1
 - 実データで実行: GSE7623 (ラット)、GSE30533 (ラット)
 - 同一プラットフォームデータ(GSE7623 + GSE30533)をマージして実行、課題2
- 実験デザイン (教科書の § 3.2.2)
- 2群間比較: 発現変動遺伝子 (DEG) 検出
 - パターンマッチング法 (相関係数の利用)
 - コードの中身をおさらい、apply関数の基本的な利用法など

対数変換の有無

①で用いている②hoge1.txtは、GSE30533 (Kamei et al., 2013)の対数変換(log₂変換)前のMAS5データ

- 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\(プロブレレベルデータ\)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p99-107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ペ。

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

対数変換の有無

- 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

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	758.1	575.5	568.6	494.5	691.0	610.6	442.0	565

対数変換の有無

①data_mas.txtは、GSE30533 (Kamei et al., 2013)の、②対数変換(\log_2 変換)後のMAS5データ

	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	11.979128
3	1367453_at	11.562853	11.587924	11.62175	11.747779	11.776342	11.633797	11.50868	11.479418
4	1367454_at	11.479343	11.676545	11.608875	11.652704	11.86008	11.70589	11.975747	12.009512
5	1367455_at	12.784752	12.58787	12.696588	12.789029	13.002298	12.678645	12.780148	12.552212
6	1367456_at	13.50768	13.533852	13.483843	13.524296	13.452217	13.472369	13.594191	13.603128
7	1367457_at	10.741903	10.136369	10.609696	10.256467	10.312411	10.49749	10.304808	10.427412

```
#####↓
### 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↓
↓
```


② hoge - GSE30533フォルダ中の③ hoge1.txtのサンプル間クラスタリングをやってみよう

対数変換の有無

- 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p99-100

「ファイル」-「ディレクトリの変更」でデスクトップ上の"E-GEODATA"フォルダに

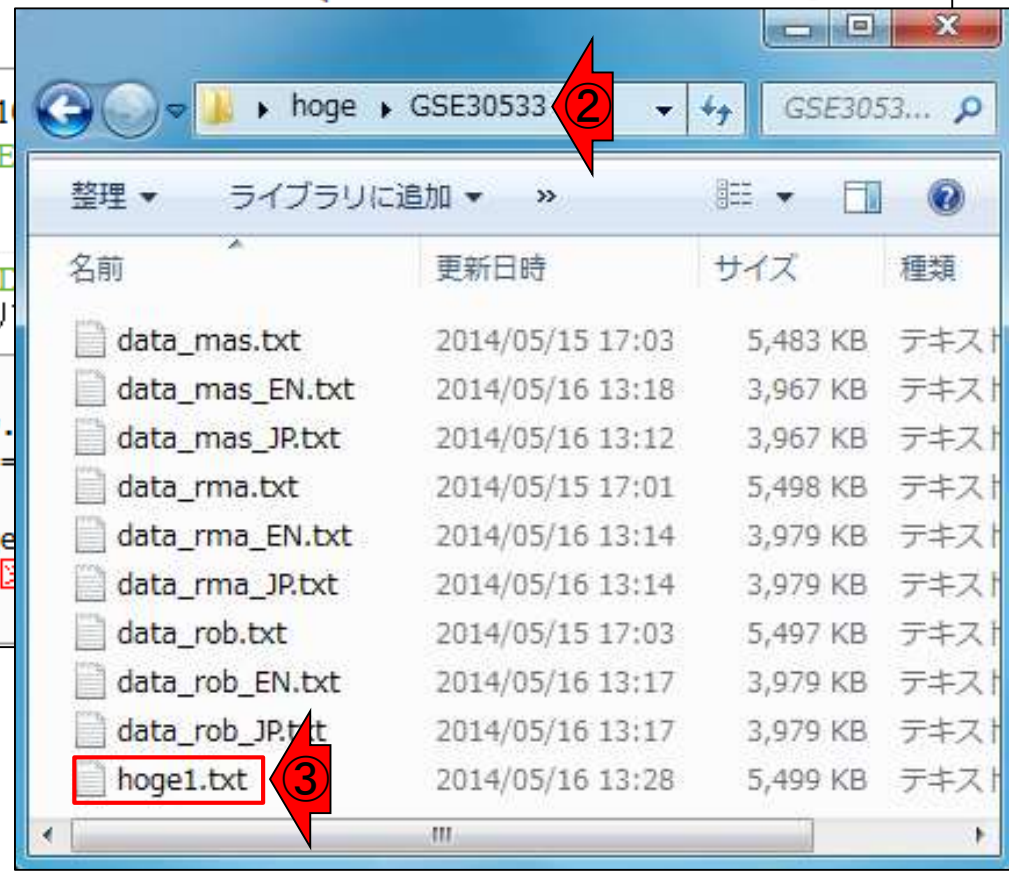
p99の網掛け部分:

書籍中では作業ディレクトリがデスクトップ上の"E-GEODATA"フォルダ内のMAS5データファイル(hoge1.txt)を置いてあるディレクトリ

```

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

```



参考

① list.files() のみだと全ファイルになるが、②や③で示すように任意の文字列を含むファイル名のみにもすることもできる

Tips: list.files

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...) **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(hogel.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 <- "hogel.txt"
data <- read.table(in_f, header=TRUE, as.is=TRUE)
colnames(data) <- c(paste("G1", 1:10))
data.dist <- as.dist(1 - cor(data))
out <- hclust(data.dist, method="ward.D2")
plot(out)

```

```

R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/GSE30533"
> list.files()
[1] "data_mas.txt"      "data_mas_EN.txt"  "data_mas_JP.txt"
[4] "data_rma.txt"      "data_rma_EN.txt"  "data_rma_JP.txt"
[7] "data_rob.txt"      "data_rob_EN.txt"  "data_rob_JP.txt"
[10] "hogel.txt"
① > list.files(pattern="hoge")
[1] "hogel.txt"
② > list.files(pattern="EN")
[1] "data_mas_EN.txt"  "data_rma_EN.txt"  "data_rob_EN.txt"
③ > |

```

サンプル名のところが変わっている理由を説明します

対数変換の有無

	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				7.7	4123
5	1367455_at	7056.6	6156.4	6638.3				4.1	6006
6	1367456_at	11647.1	11860.3	11456.2				6.9	12449
7	1367457_at	1712.5	1125.5	1562.6	1223.2	1271.6	1445.6	1264.9	1377
8	1367458_at	759.4	575.5	569.6	494.5	691.0	610.6	442.0	565

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

R Console

```

[10] "hogel.txt"
> list.files(pattern="hoge")
[1] "hogel.txt"
> list.files(pattern="EN")
[1] "data_mas_EN.txt" "data_
> in_f <- "hogel.txt"
> data <- read.table(in_f, he
> colnames(data) <- c(paste(
> data.dist <- as.dist(1 - co
> out <- hclust(data.dist, me
> plot(out)
> |
                    
```

R Graphics: Device 2 (ACTIVE)

Cluster Dendrogram

data.dist
hclust(*, "average")

Tips

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
```

```
R Console
> data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
① > 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"
② > 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"
> |
```


黒下線部分がG1群のサンプル名作成に相当する部分。①その部分のみ実行。②同じ結果だがやり方を微妙に変えている。③Iron_defに文字を変更

Tips

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
```

R Console

```
> 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"
① > paste("G1_", 1:5, sep="")
[1] "G1_1" "G1_2" "G1_3" "G1_4" "G1_5"
② > paste("G1", 1:5, sep=" ")
[1] "G1_1" "G1_2" "G1_3" "G1_4" "G1_5"
③ > paste("Iron_def", 1:5, sep="")
[1] "Iron_def1" "Iron_def2" "Iron_def3" "Iron_def4" "Iron_def5"
> paste("Iron_def", 1:3, sep="")
[1] "Iron_def1" "Iron_def2" "Iron_def3"
> paste("Iron_def", c(2,4,5), sep="")
[1] "Iron_def2" "Iron_def4" "Iron_def5"
> |
```

対数変換の有無

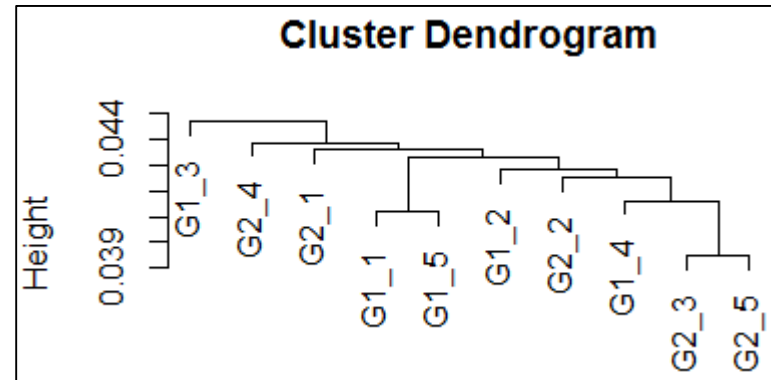
① GSE30533 (Kamei et al., 2013)の対数変換後のMAS5データ(data_mas_EN.txt)でもクラスタリングを行い、②対数変換前(hoge1.txt)の結果と比較する

rancode_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)
```

	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	11.62175	11.747779	11.776342	11.633797	11.50868	11.4794
4	1367454_at	11.479343	11.676545	1			0589	11.975747	12.0095
5	1367455_at	12.784752	12.58787	1			8645	12.780148	12.5522
6	1367456_at	13.50768	13.533852	1			2369	13.594191	13.603
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



対数変換の有無

①対数変換後の結果。対数変換の有無にかかわらずクラスタリング結果(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

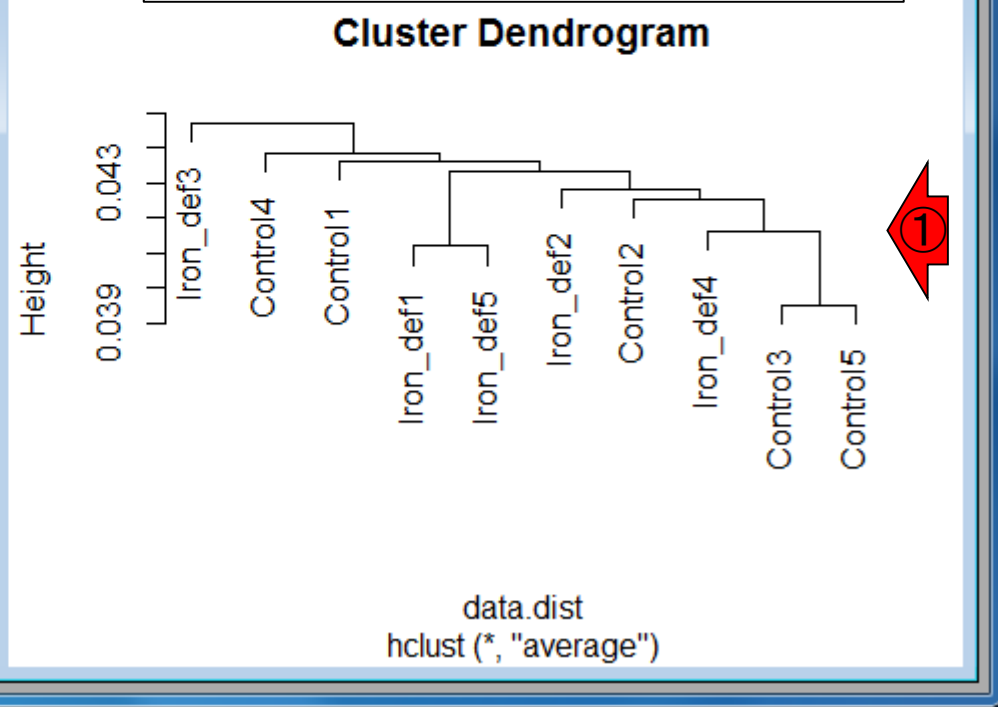
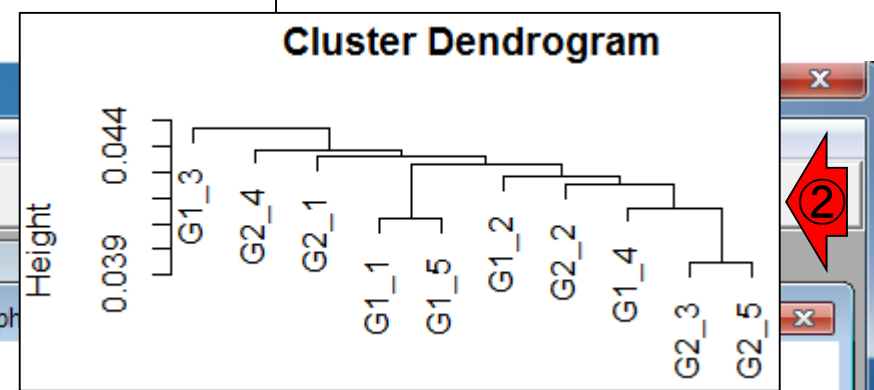
```
#####↓
### MAS5データのクラスタリング ###↓
#####↓
in_f <- "data_mas_EN.txt"
data <- read.table(
#colnames(data) <-
data.dist <- as.dist
out <- hclust(data,
plot(out)↓
←
```

rancode_clustering.txt

The screenshot shows the R GUI interface. The R Console window contains the following code and output:

```
[1] "Iron_def1" "Iron_def2"
> paste("Iron_def", 1:3, sep="")
[1] "Iron_def1" "Iron_def2"
> paste("Iron_def", c(2,4,5), sep="")
[1] "Iron_def2" "Iron_def4"
> getwd()
[1] "C:/Users/kadota/Desktop"
> in_f <- "data_mas_EN.txt"
> data <- read.table(in_f, header=TRUE)
> #colnames(data) <- c(paste("Iron_def", 1:5, sep=""),
> #                    paste("Control", 1:5, sep=""))
> data.dist <- as.dist(1 - cor(data))
> out <- hclust(data.dist, method="average")
> plot(out)
```

The R Graph window displays a dendrogram with the following labels: Iron_def3, Control4, Control1, Iron_def1, Iron_def5, Iron_def2, Control2, Iron_def4, Control3, Control5. The y-axis is labeled 'Height' with values 0.039 and 0.043. The plot title is 'data.dist hclust(*, "average")'.



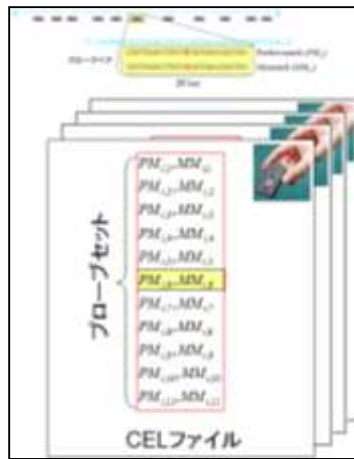
Contents

- 実データ概観
 - GSE2361 (ヒト)、GSE7623 (ラット)、GSE30533 (ラット)
- クラスタリング (教科書の § 3.2.1)
 - 対数変換の有無 (Spearman相関係数を使う場合は気にしなくてよい)
 - 階層的 vs. 非階層的、様々な選択肢
 - 距離の定義: ベクトル間、クラスター間
 - 実データで実行: GSE2361 (ヒト)、課題1
 - 実データで実行: GSE7623 (ラット)、GSE30533 (ラット)
 - 同一プラットフォームデータ(GSE7623 + GSE30533)をマージして実行、課題2
- 実験デザイン (教科書の § 3.2.2)
- 2群間比較: 発現変動遺伝子 (DEG) 検出
 - パターンマッチング法 (相関係数の利用)
 - コードの中身をおさらい、apply関数の基本的な利用法など

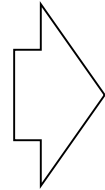
階層的 vs. 非階層的

- 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)

様々な選択肢



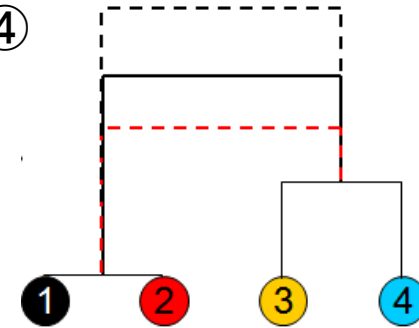
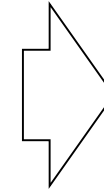
①前処理法



	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②

クラスタリング③と④



①前処理法

- ・MAS5
- ・RMA
- ・RobLoxBioC

...

...

×

②スケールング

- ・対数変換
- ・相対値(0~1)
- ・Z-score化

...

...

×

③距離

- ・1-相関係数
- ・ユークリッド
- ・マンハッタン
- ・キャンベラ

...

×

④群の併合

- ・単連結法
- ・完全連結法
- ・平均連結法
- ・ワード法

...

様々な選択肢

■ 決めておくべき2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

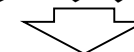
- クラスタ間の距離を定義する方法、とほぼ同じ
- 最短距離法、平均連結法、ワード法など

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。私は、「外れサンプルのチェック」や「発現変動遺伝子の有無や数」の見当をつける目的で行う

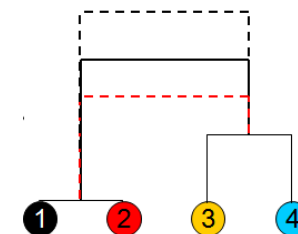
入力例

	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



Danielsson et al., *Brief Bioinform.*, **16**: 941-949, 2015

Tang et al., *BMC Bioinformatics*, **16**: 361, 2015

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 $D(x,y)$

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

- xとyの発現パターンが酷似 $\rightarrow r \approx 1$
- xとyの発現パターンがばらばら $\rightarrow r \approx 0$
- xとyの発現パターンがほぼ正反対 $\rightarrow r \approx -1$

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

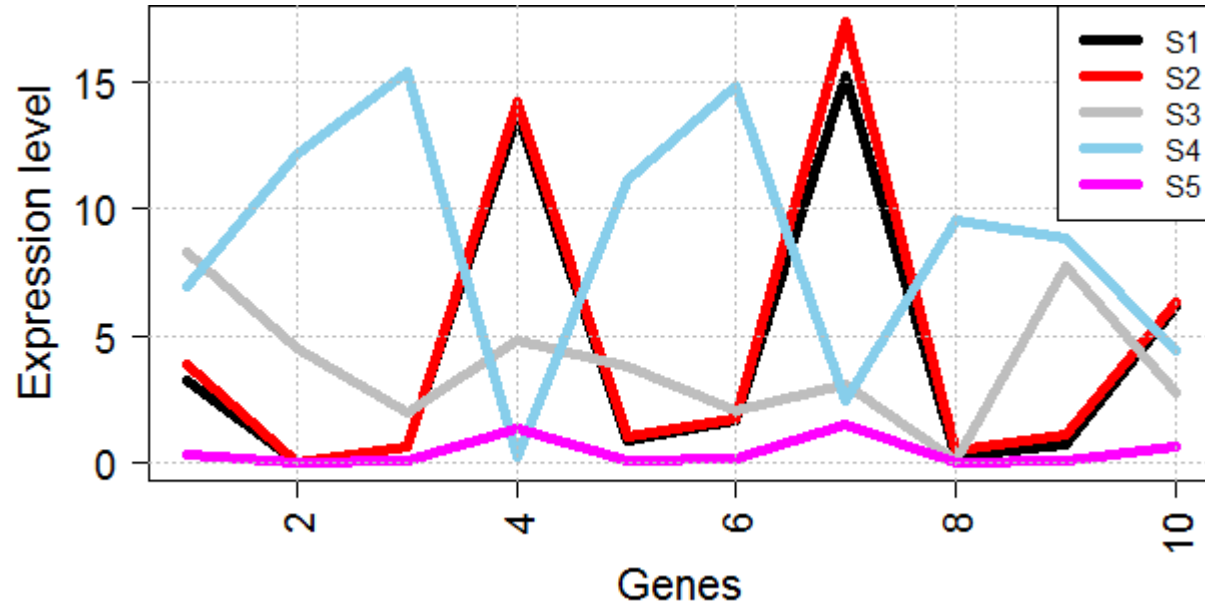
$$\text{距離 } D(x,y) = 1 - r \quad (0 \leq D \leq 2) \quad \begin{cases} r = 1 \rightarrow D = 1 - 1 = 0 \\ r = 0 \rightarrow D = 1 - 0 = 1 \\ r = -1 \rightarrow D = 1 - (-1) = 2 \end{cases}$$

パターンが似ていれば0に近い値、逆パターンに近ければ最大値の2に近い値になっていることが分かる

相関係数 → 距離 (計算例)

■ ベクトルxとyの発現パターンの距離 $D(x,y)$

	S ¹	S ²	S ³	S ⁴	S ⁵
g1	3.24	3.89	8.27	6.93	0.32
g2	0.01	0.03	4.55	12.17	0
g3	0.65	0.69	1.98	15.39	0.07
g4	13.73	14.21	4.83	0.28	1.38
g5	0.89	1.05	3.84	11.16	0.09
g6	1.65	1.74	2.11	14.82	0.17
g7	15.21	17.33	3.13	2.49	1.51
g8	0.26	0.52	0.08	9.53	0.03
g9	0.73	1.11	7.76	8.88	0.07
g10	6.18	6.36	2.81	4.47	0.62



相関係数 $r_{S^1S^2} = 0.998 \rightarrow$ 距離 $D_{S^1S^2} = 1 - 0.998 = 0.002$

相関係数 $r_{S^1S^3} = 0.035 \rightarrow$ 距離 $D_{S^1S^3} = 1 - (0.035) = 0.965$

相関係数 $r_{S^1S^4} = -0.851 \rightarrow$ 距離 $D_{S^1S^4} = 1 - (-0.851) = 1.851$

②の赤枠全体をコピーで実行。
③距離計算結果を格納した data.dist の中身を眺める

相関係数 → 距離

- 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p99-107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。

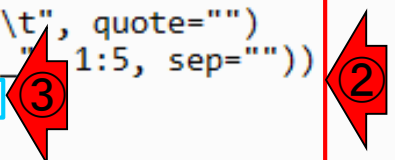
p99の網掛け部分:

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。

```

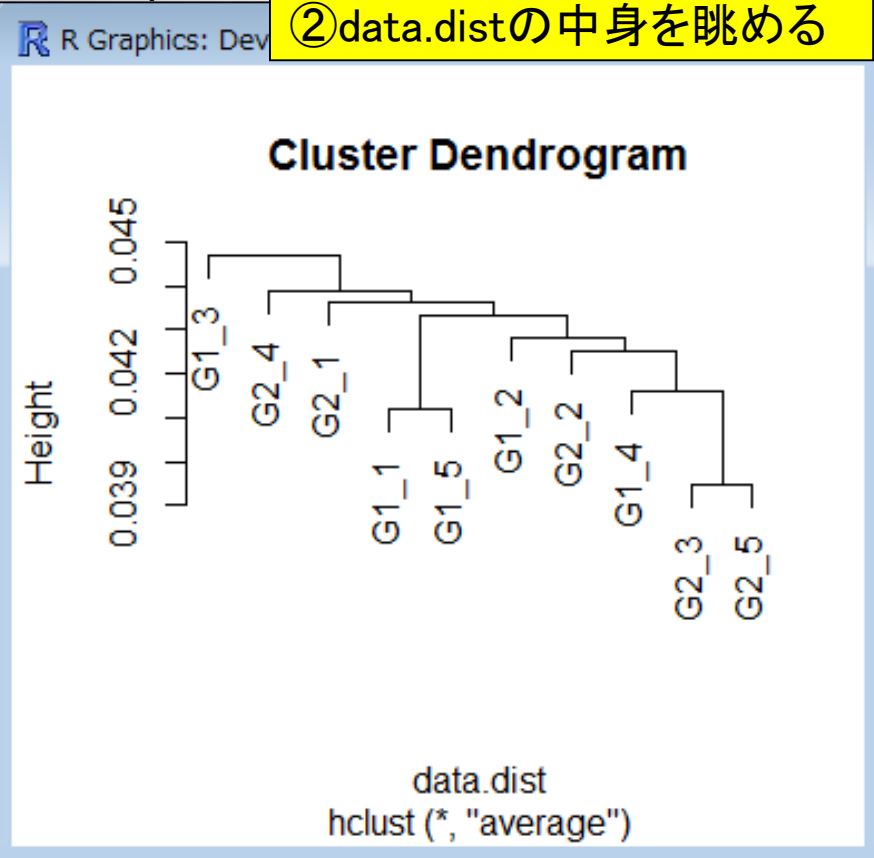
in_f <- "hoge1.txt"
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))
data.dist <- as.dist(1 - cor(data, method = "spearman"))
out <- hclust(data.dist, method = "average")
plot(out)
# 図3-1作成部分

```



- ① 距離計算結果を格納した
- ② data.distの中身を眺める

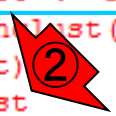
相関係数 → 距離



```

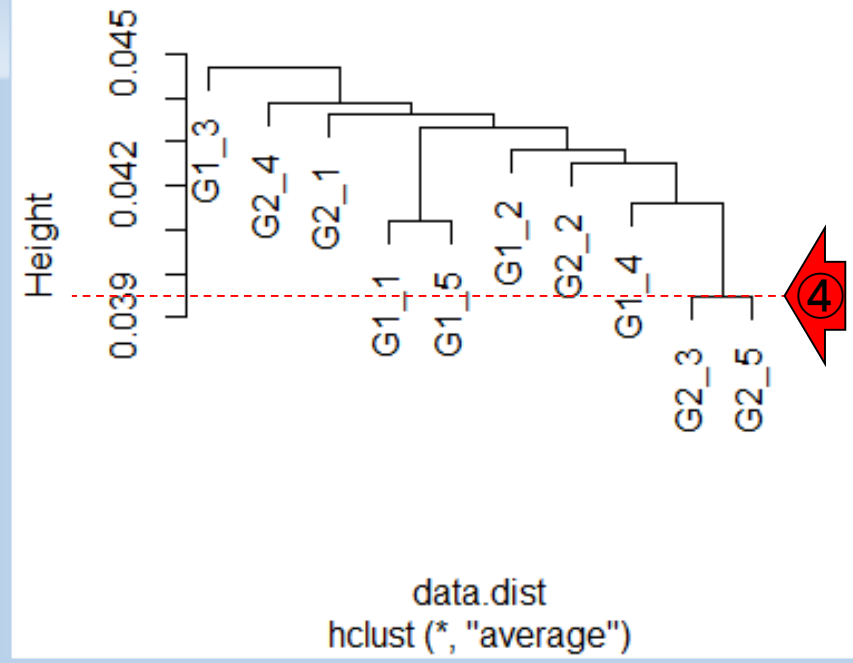
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/GSE30533"
> in_f <- "hogel.txt"
> data <- read.table(in_f, header=TRUE, row.names=1, sep="\t",
> colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1:5, sep=""))
> data.dist <- as.dist(1 - cor(data, method = "spearman"))
> out <- hclust(data.dist,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

①の0.03949559がクラスタリング結果の「②G2_3と③G2_5の発現ベクトル間の距離」に相当し、それはクラスタリング結果の④高さ部分に相当する



```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/GSE30533"
> in_f <- "hogel.txt"
> data <- read.table(in_f, header=TRUE, row.names=1, sep="\t",
> colnames(data) <- c(paste("G1_", 1:5, sep=""), paste("G2_", 1
> 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	

Tips: 相関係数

教科書p100 参考

Spearman相関係数とPearson相関係数の関係。①Spearman相関係数、② $1 - \text{Spearman相関係数}$ 、③Pearson相関係数、④rank関数を用いて順位変換後のPearson相関係数、⑤列名で計算することもできるというTips、⑥順位変換後のSpearman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  
> |
```


他の類似性尺度

■ 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 $D = \sum_{i=1}^n \frac{|x_i - y_i|}{|x_i + y_i|}$

□ ...

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。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

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

名前が仰々しいだけで計算自体は大したことありません

計算例 (サンプルxとy間の距離D)

	A	B	C	D	E	F
1		x	y		xi - yi	xi - yi / xi + yi
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

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

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

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

?関数名

(Rで)マイクロアレイデータ解析

(last modified 2015/05/16, since 2005)

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 | ID変換 | について (last modified 2014/06/03)

前処理 | ID変換 | probe ID -> gene symbol (last modified 2014/06/03)

前処理 | ID変換 | probe ID -> Entrez ID (last modified 2014/06/03)

前処理 | ID変換 | probe ID -> その他 (last modified 2014/06/03)

前処理 | ID変換 | 同じ遺伝子名を持つものを抽出 (last modified 2014/06/03)

解析 | 基礎 | 共通遺伝子の抽出 (last modified 2014/06/03)

解析 | 基礎 | ベクトル間の距離 (last modified 2014/06/03)

解析 | 基礎 | 遺伝子ごとの各種統計量 (last modified 2014/06/03)

解析 | 基礎 | 最大発現量を示す組織の特定 (last modified 2014/06/03)

解析 | 基礎 | 似た発現パターンを持つ遺伝子の抽出 (last modified 2014/06/03)

解析 | 基礎 | 平均-分散プロット (last modified 2014/06/03)

解析 | 基礎 | クラスターリング | 階層的 | について (last modified 2014/06/03)

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" #入力ファイル名を指定してin_fに格納

#データファイルの読み込み
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")#i

#本番
dist(t(data), method="euclidean") #ユークリッド(Euclidean)距離
dist(t(data), method="manhattan") #マンハッタン(Manhattan)距離
dist(t(data), method="maximum") #チェビシェフ(Chebyshev)距離
dist(t(data), method="canberra") #キャンベラ(Canberra)距離
1 - cor(data, method="pearson") #1 - Pearson相関係数

dist(t(data), method="binary") #ハミング(Hamming)距離
dist(t(data), method="minkowski") #ミンコフスキー(Minkowski)距離
1 - cor(data, method="spearman") #1 - Spearman相関係数
```

?関数名

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計算し、「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリにコピー

1. 10 genes × 2 samplesのデータファイル([sample19.txt](#))の場合:

```
in_f <- "sample19.txt" #入力ファイル名
#データファイルの読み込み
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t")
#本番
dist(t(data), method="euclidean") #ユークリッド(Euclid)距離
dist(t(data), method="manhattan") #マンハッタン(Manhattan)距離
dist(t(data), method="maximum") #チェビシェフ(Chebyshev)距離
dist(t(data), method="canberra") #キャンベラ(Canberra)距離
1 - cor(data, method="pearson") #1 - Pearson相関係数
dist(t(data), method="binary") #ハミング(Hamming)距離
dist(t(data), method="minkowski") #ミンコフスキー(Minkowski)距離
1 - cor(data, method="spearman") #1 - Spearman相関係数
```

```
> dist(t(data), method="euclidean") #ユークリッド (Euclid)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") #チェビシェフ (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
> |
```

```
R Console
> ?dist
starting httpd help server ... done
> dist(t(data), method="euclidean")
      sample1
sample2 2.792848
> dist(t(data))
      sample1
sample2 2.792848
> |
```

参考

①「?dist」。②ユークリッド距離でよければ、「method="xxx"」のところを記述しなくてもいいようだ。③“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

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)

## 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

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.



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?!(ワード法も?!)...。いろいろ試して総合的に判断することが重要

クラスター間の距離の定義

- 最短距離法(単連結法; single-linkage)
- 最長距離法(完全連結法; complete-linkage)
- 群平均法(平均連結法; average-linkage)
- 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...

Contents

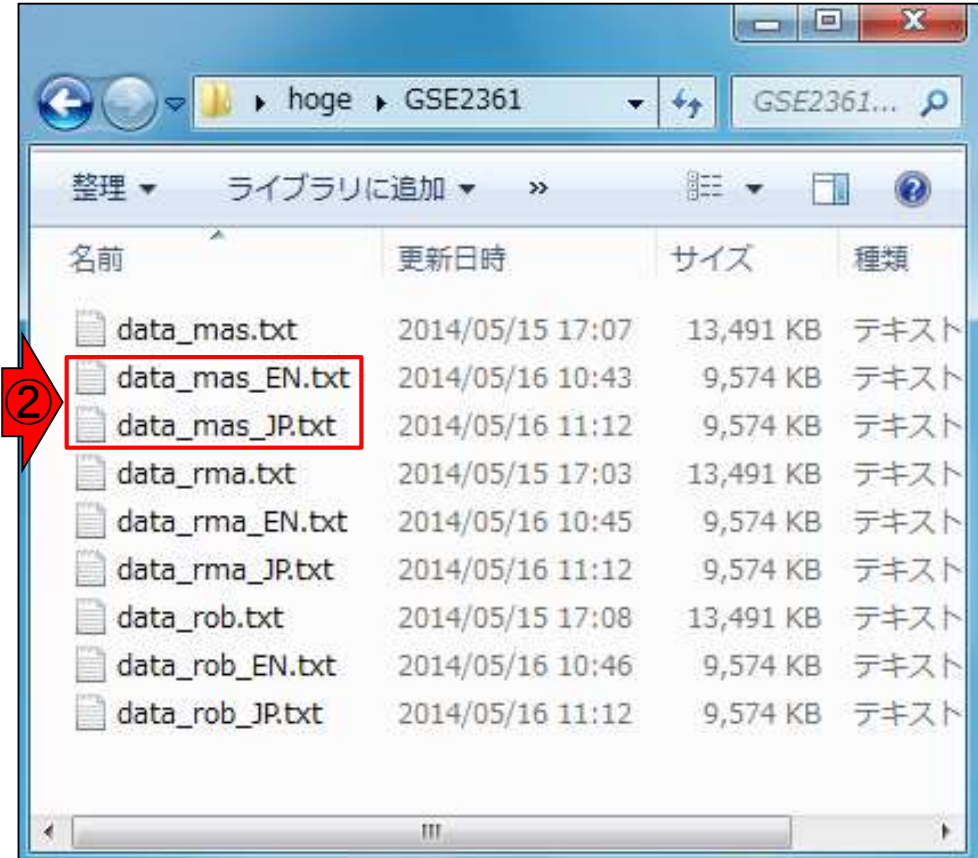
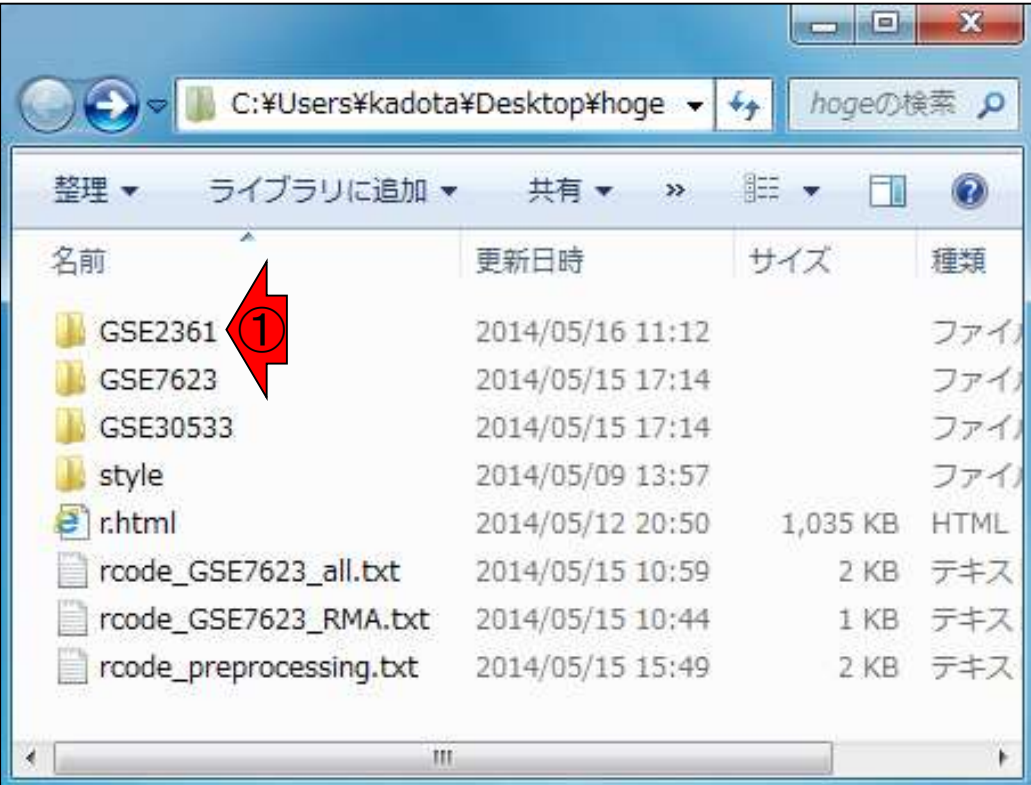
- 実データ概観
 - GSE2361 (ヒト)、GSE7623 (ラット)、GSE30533 (ラット)
- クラスタリング (教科書の § 3.2.1)
 - 対数変換の有無 (Spearman相関係数を使う場合は気にしなくてよい)
 - 階層的 vs. 非階層的、様々な選択肢
 - 距離の定義: ベクトル間、クラスター間
 - 実データで実行: GSE2361 (ヒト)、**課題1**
 - 実データで実行: GSE7623 (ラット)、GSE30533 (ラット)
 - 同一プラットフォームデータ(GSE7623 + GSE30533)をマージして実行、**課題2**
- 実験デザイン (教科書の § 3.2.2)
- 2群間比較: 発現変動遺伝子 (DEG) 検出
 - パターンマッチング法 (相関係数の利用)
 - コードの中身をおさらい、apply関数の基本的な利用法など

①hoge - GSE2361フォルダ中の、②MAS5データを用いてサンプル間クラスタリングをやってみよう

GSE2361(ヒト)

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)、...



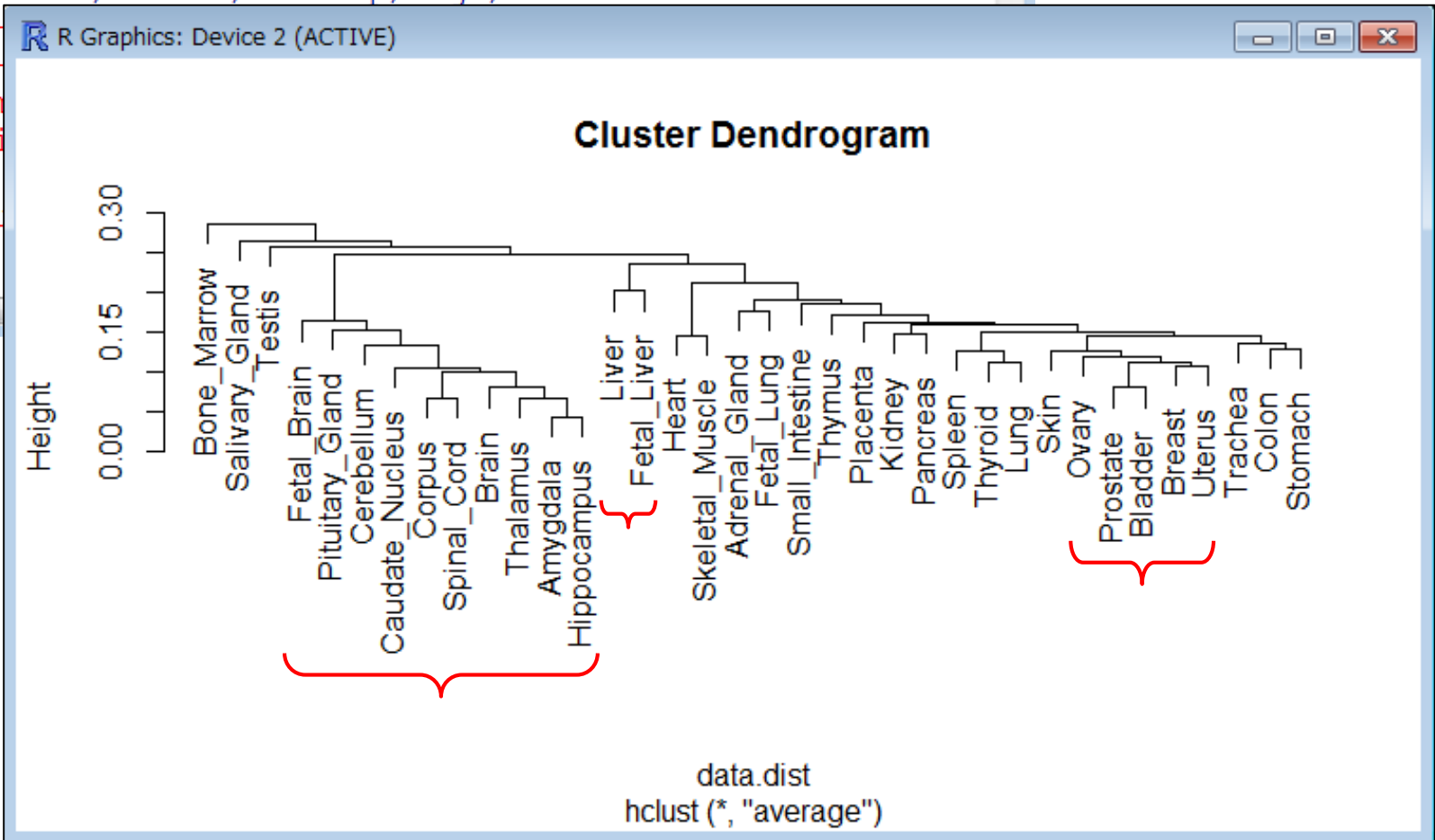
GSE2361 (ヒト)

rancode_clustering.txt

```
#####↓
### MAS5データのクラスタリング ###↓
#####↓
in_f <- "data.mas.EN.txt"↓
```

```
data <- readR
#colnames(da
data.dist <-
out <- hclus
plot(out)↓
<
```

```
R R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/GSE2361"
```



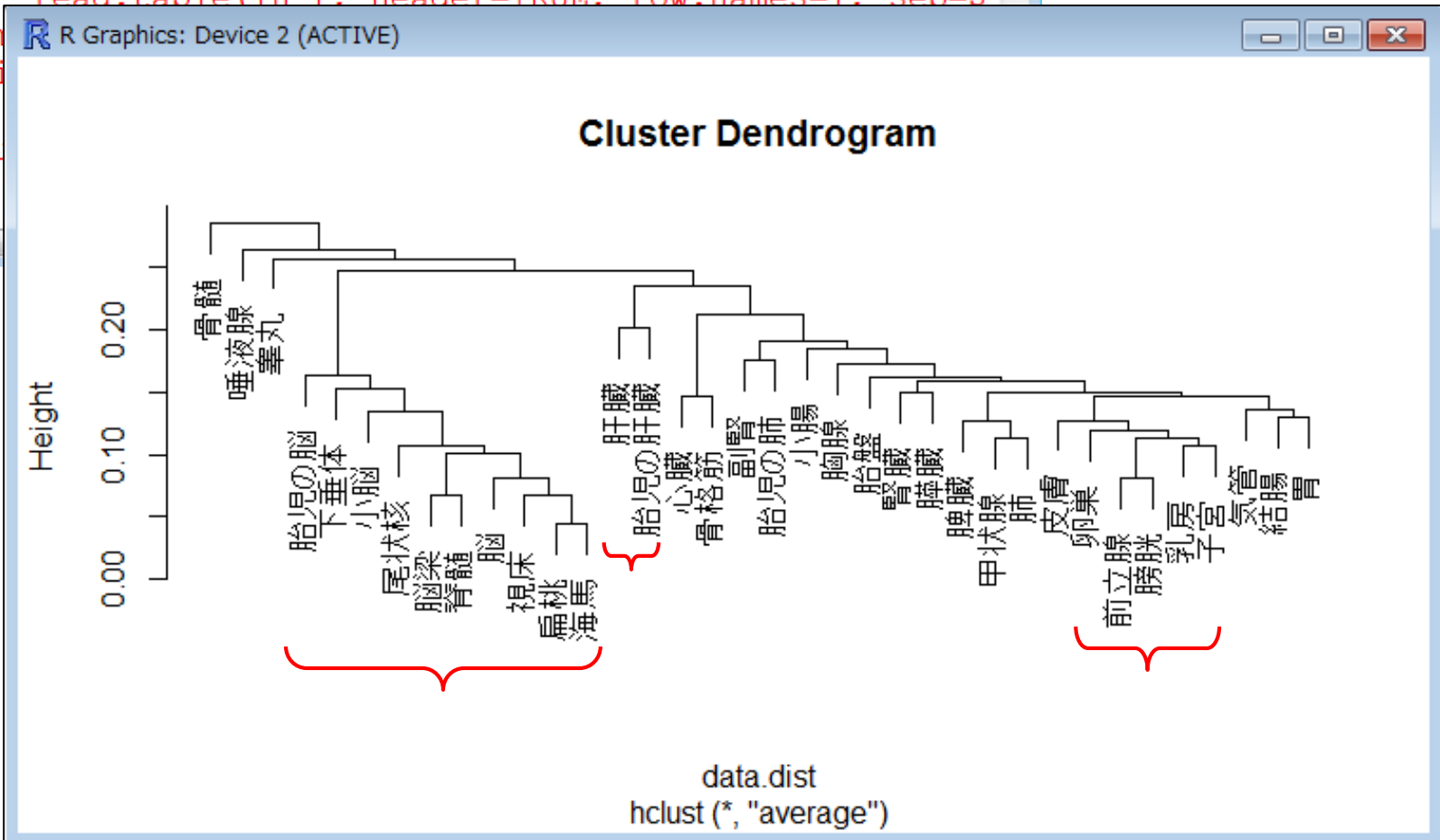
GSE2361 (ヒト)

日本語環境ではない場合?!
文字化けすることもあるよ
うですのでご注意ください

rcode_clustering.txt

```
#####↓
### MAS5データのクラスタリング ###↓
#####↓
in_f <- "data_mas_JP.txt"
data <- readTable(in_f, header=TRUE, row.names=1, sep=$)
#colnames(data)
data.dist <- dist(data)
out <- hclust(data.dist)
plot(out)
```

```
R Console
> in_f <- "data_mas_JP.txt"
> data <- readTable(in_f, header=TRUE, row.names=1, sep=$)
> #colnames(data)
> data.dist <- dist(data)
> out <- hclust(data.dist)
> plot(out)
```



Tips (ファイル保存)

①例題3。PNG形式ファイルとして縦横の大きさを指定して保存することもできる。テンプレートとの違いは赤矢印部分。②rcode_clustering_png.txt

解析 | クラスタリング | 階層的 | hclust

①

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。

1. サンプルデータに表に表 in par #入 dat #本 dat out plo

```
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
```

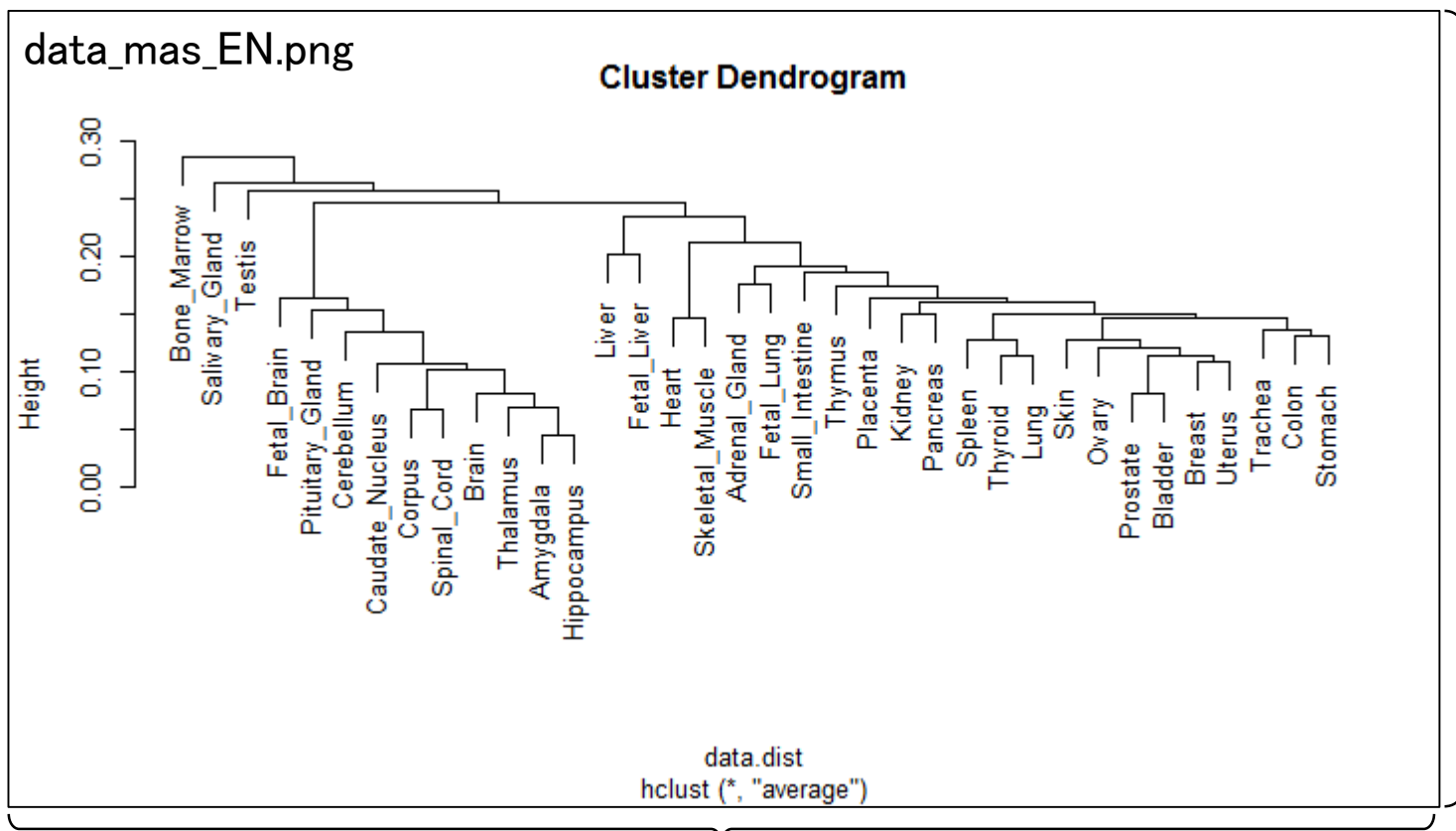
rcode_clustering_png.txt ②

```
##### ↓
### 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) #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 ↓
```

Tips (ファイル保存)

rancode_clustering_png.txt ②

```
##### ↓
### 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, head=1)
  ↓
#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)
→ dev.off()
```



720

400

課題1

GSE2361のサンプル間クラスタリングをRMA, およびRMX前処理法を適用したデータについても行い、結果を考察せよ。距離の定義はデフォルトのままでよい。①の部分をdata_rma...やdata_rob...に変更すればよいが...

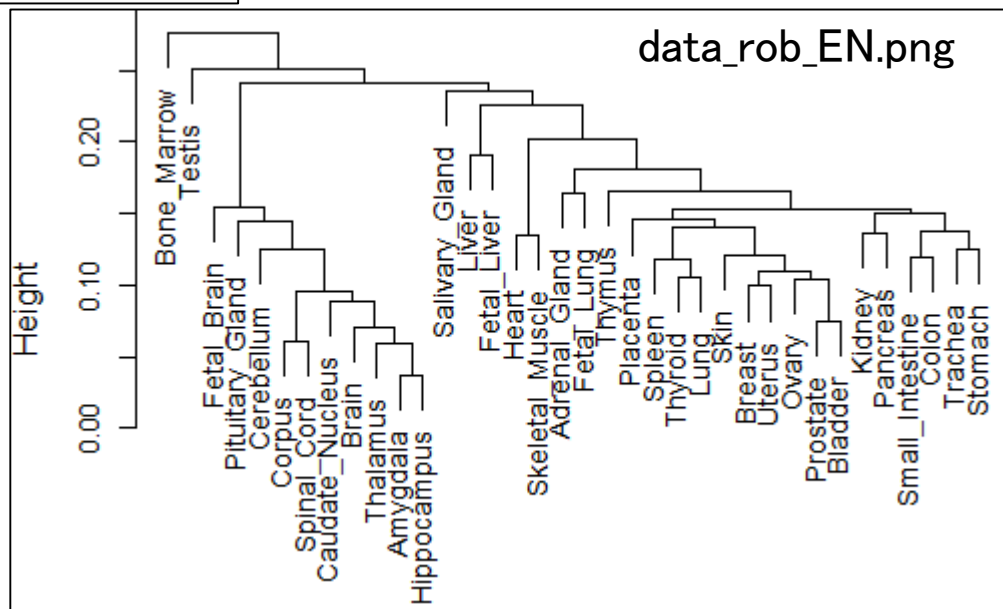
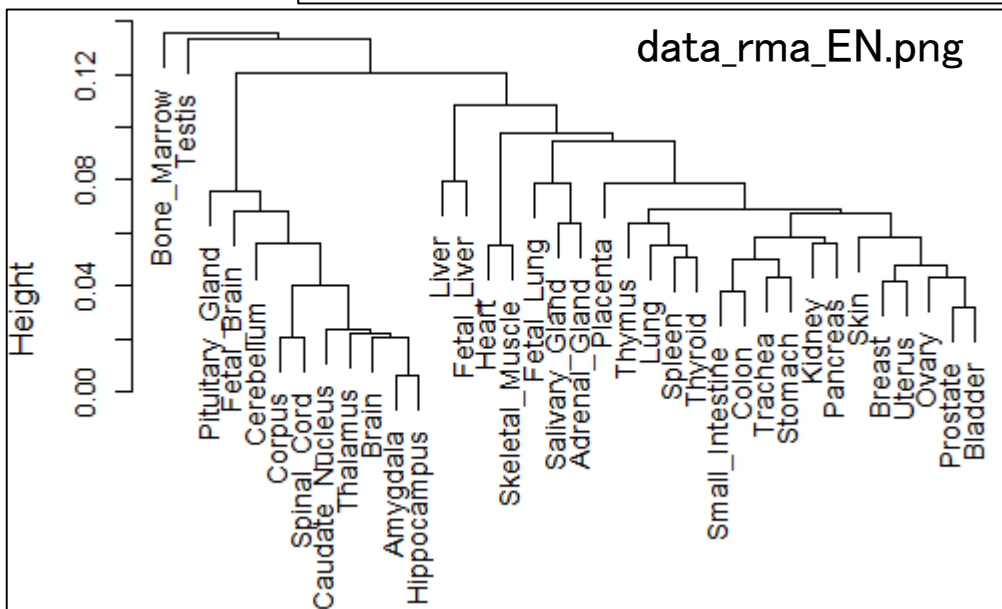
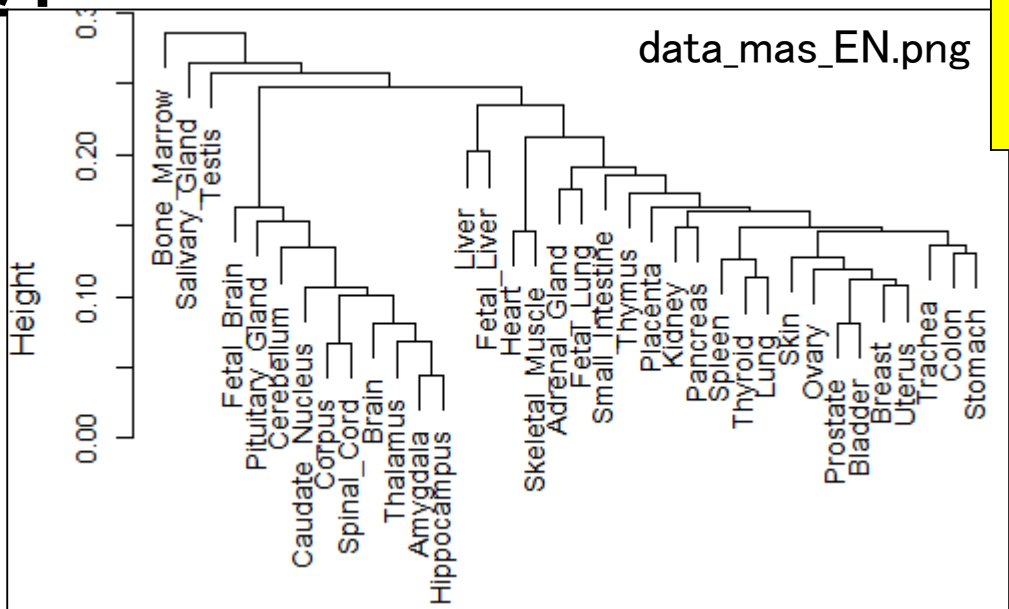
rcode_clustering_png.txt

```
#####↓
### 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) #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↓
```



課題1

①rcode_clustering_png_kadai2.txt
をコピペで実行すれば、このような結果
が得られる。実際の論文用の図は、
余白指定や文字の大きさなどをいろ
いろ変えて、見栄えをよくします



①例題5はユークリッド距離を用いる場合のテンプレートです。②param1(やparam2)などをいろいろいじって結果を眺め、可能な範囲で考察してください

発展課題 (optional)

解析 | クラスタリング | 階層的 | hclust

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、以下をコピー

1. サンプルデータ 5. 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の結果を表示するやり方

```
in_f <- "sample3.txt"
param <- "average"
#入力ファイル名を指定してin_fに格納
data <- read.table(in_f, header=TRUE, row.names=1, sep="\t", quote="")
#本番
data.dist <- dist(t(data), method=param)
out <- hclust(data.dist, method="average")
plot(out)
```

サンプル間クラスタリング(距離: ユークリッド距離(euclidean)、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"
out_f <- "hoge5.png"
param1 <- "euclidean"
param2 <- "average"
param_fig <- c(500, 400)

#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を読み込み

#本番
data.dist <- dist(t(data), method=param1)#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=param2)#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)
dev.off()

#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
#距離(dist)を指定
#方法(method)を指定
#ファイル出力時の横幅と縦幅を指定(単位はピクセル)
#樹形図(デンドログラム)の表示
#おまじない
```

Contents

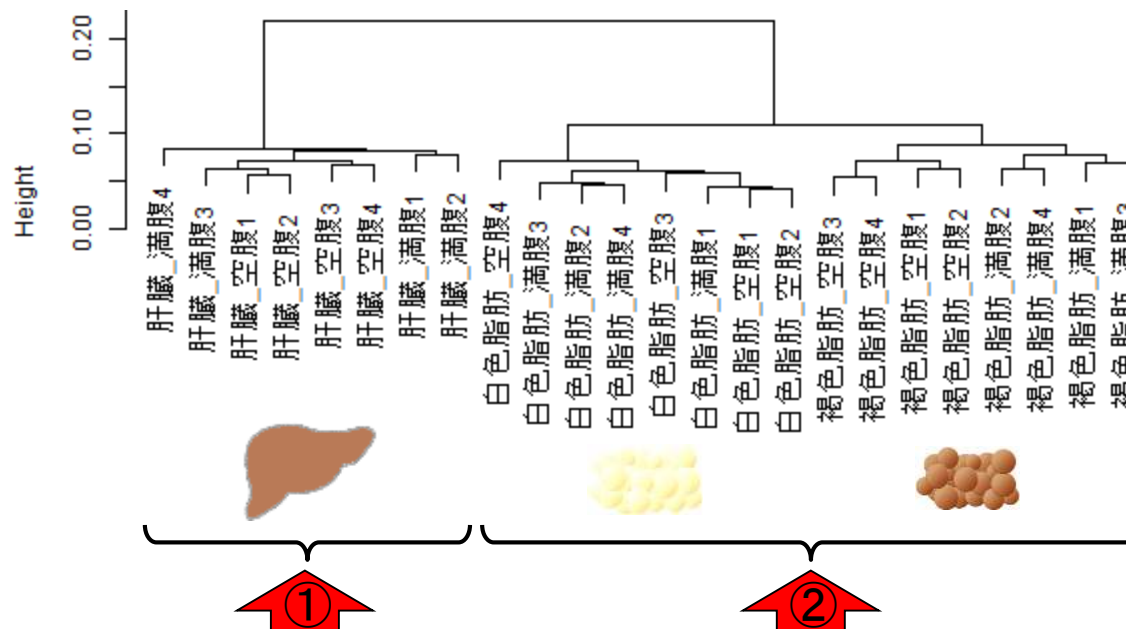
- 実データ概観
 - GSE2361 (ヒト)、GSE7623 (ラット)、GSE30533 (ラット)
- クラスタリング (教科書の § 3.2.1)
 - 対数変換の有無 (Spearman相関係数を使う場合は気にしなくてよい)
 - 階層的 vs. 非階層的、様々な選択肢
 - 距離の定義: ベクトル間、クラスター間
 - 実データで実行: GSE2361 (ヒト)、課題1
 - 実データで実行: GSE7623 (ラット)、GSE30533 (ラット)
 - 同一プラットフォームデータ(GSE7623 + GSE30533)をマージして実行、課題2
- 実験デザイン (教科書の § 3.2.2)
- 2群間比較: 発現変動遺伝子 (DEG) 検出
 - パターンマッチング法 (相関係数の利用)
 - コードの中身をおさらい、apply関数の基本的な利用法など

GSE7623 (ラット)

MAS5データを(1 - Spearman)相関係数でクラスタリングした結果。①肝臓と②脂肪間で大きく2つのクラスターに分かれていることがわかる

Nakai et al., *BBB*, 72: 139–148, 2008

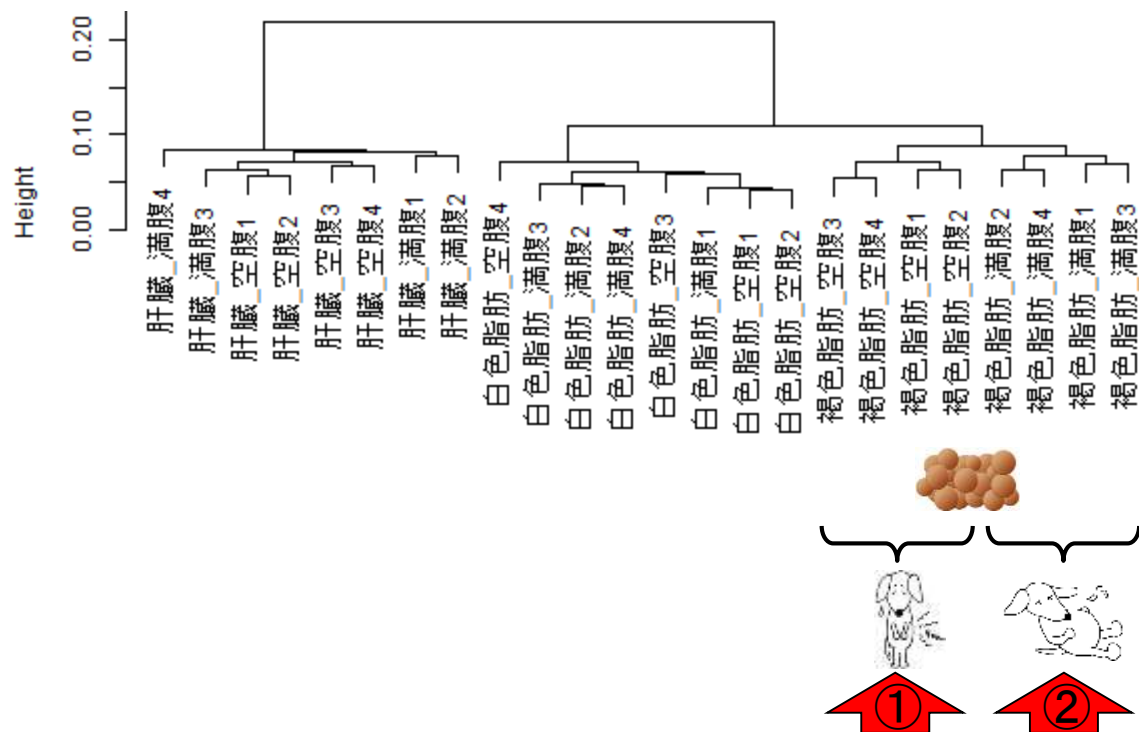
- GSE7623、GPL1355 (Affymetrix Rat Genome 230 2.0 Array)、31,099 probesets
- ラット24サンプル: Brown adipose tissue (褐色脂肪組織; BAT) 8サンプル、White adipose tissue (白色脂肪組織; WAT) 8サンプル、Liver (肝臓; LIV) 8サンプル
 - BAT 8サンプル: 通常(BAT_fed) 4サンプル 対 24時間絶食(BAT_fas) 4サンプル
 - WAT 8サンプル: 通常(WAT_fed) 4サンプル 対 24時間絶食(WAT_fas) 4サンプル
 - 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



GSE7623 (ラット)

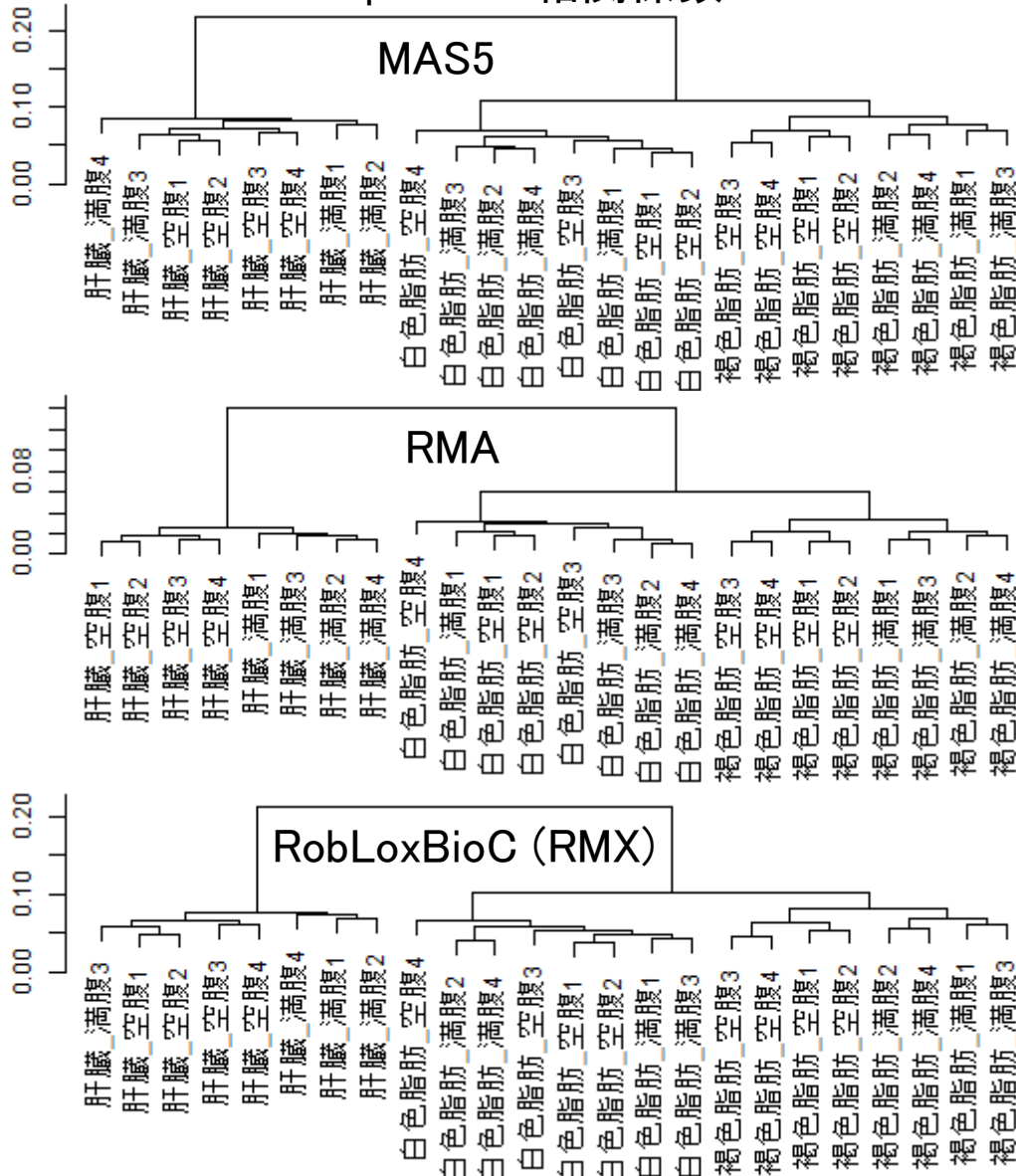
Nakai et al., *BBB*, 72: 139–148, 2008

- 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

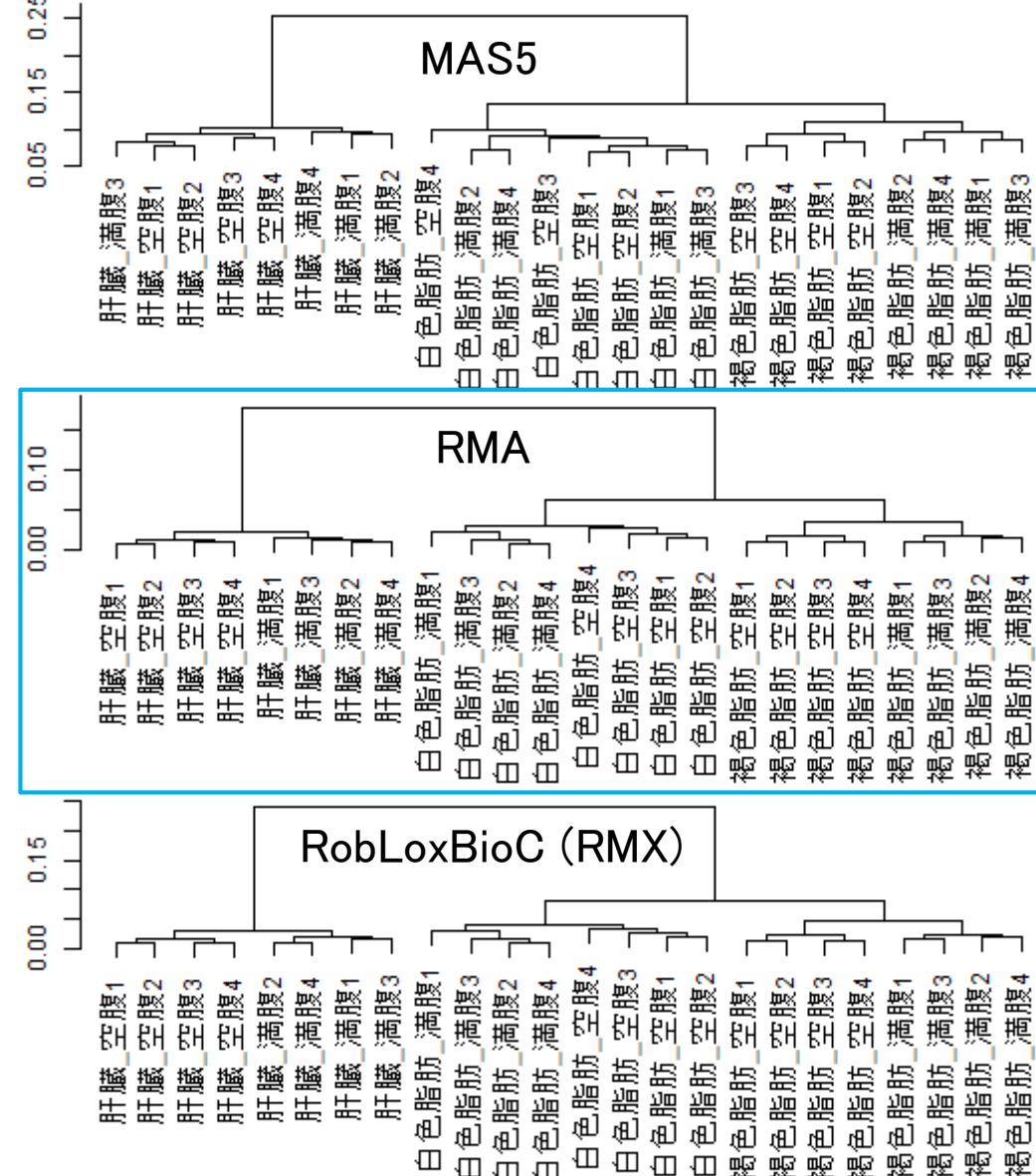


GSE7623 (ラット)

1 - Spearman相関係数



1 - Pearson相関係数



GSE7623 (ラット)

①原著論文のFig. 1(の一部)を再現できてます

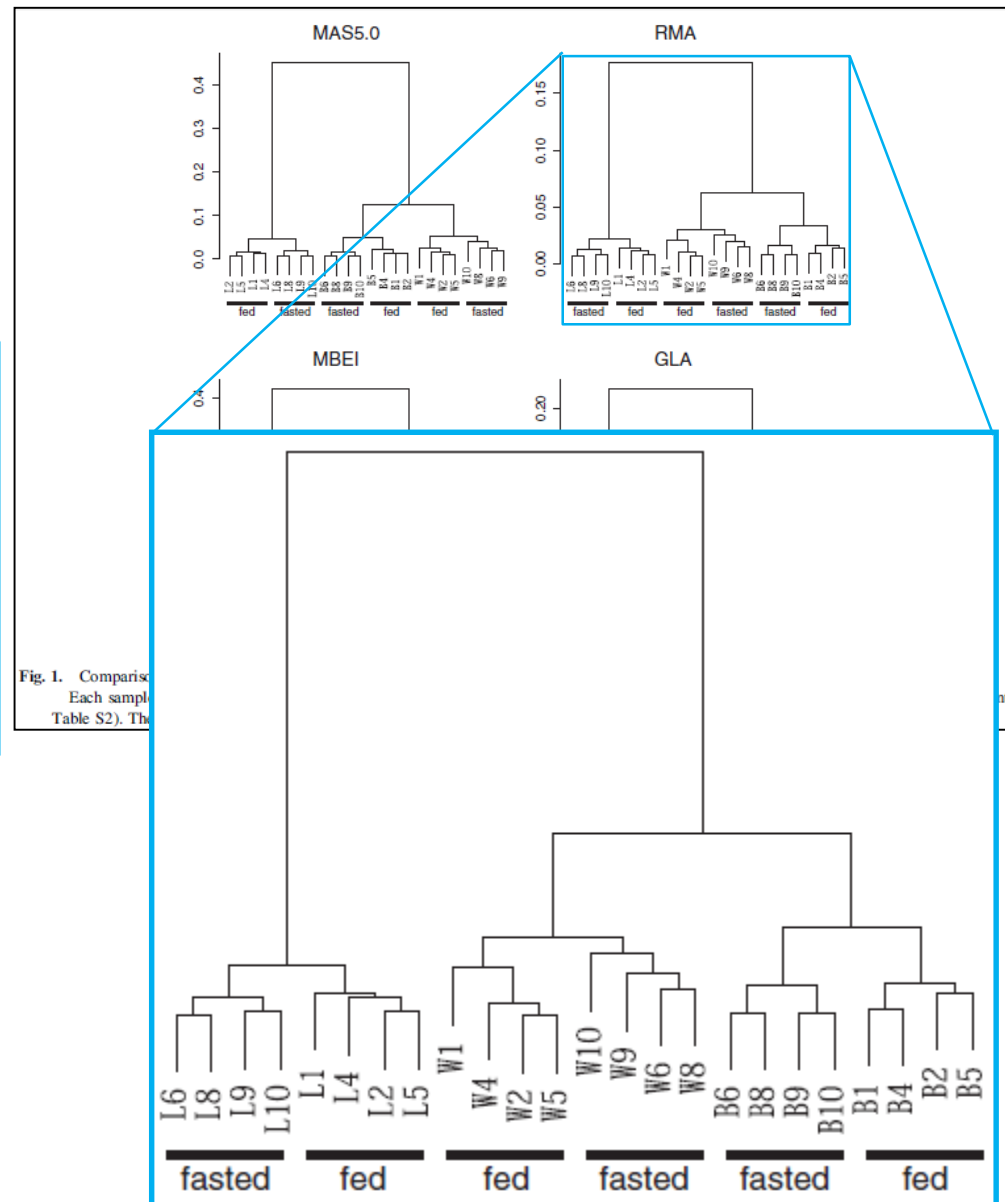
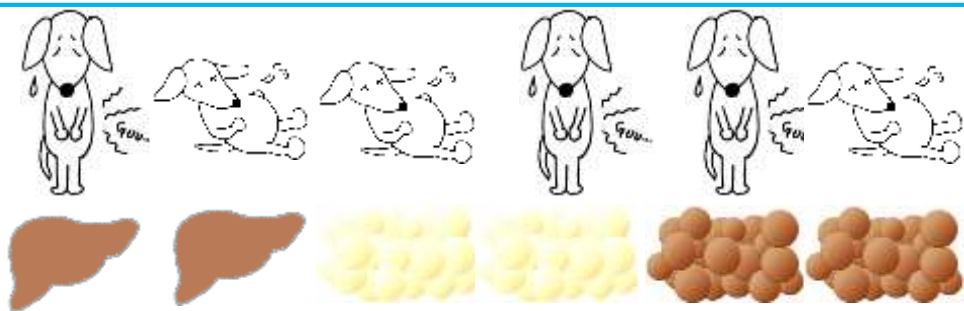
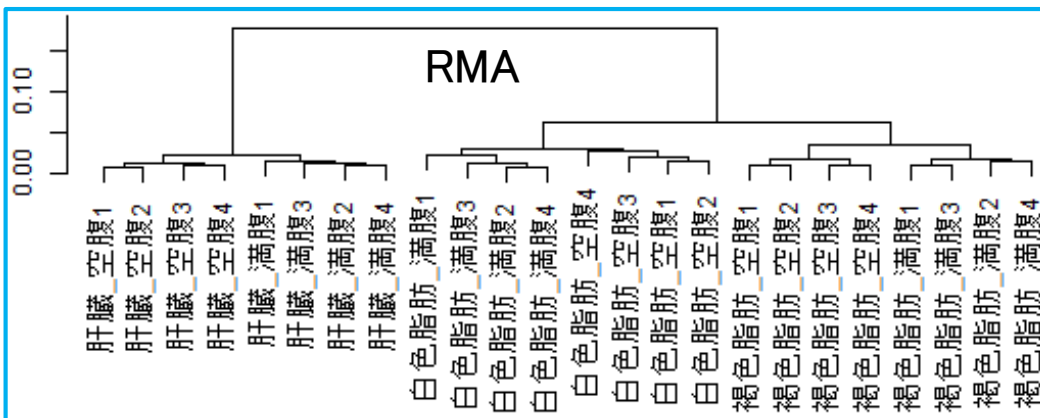


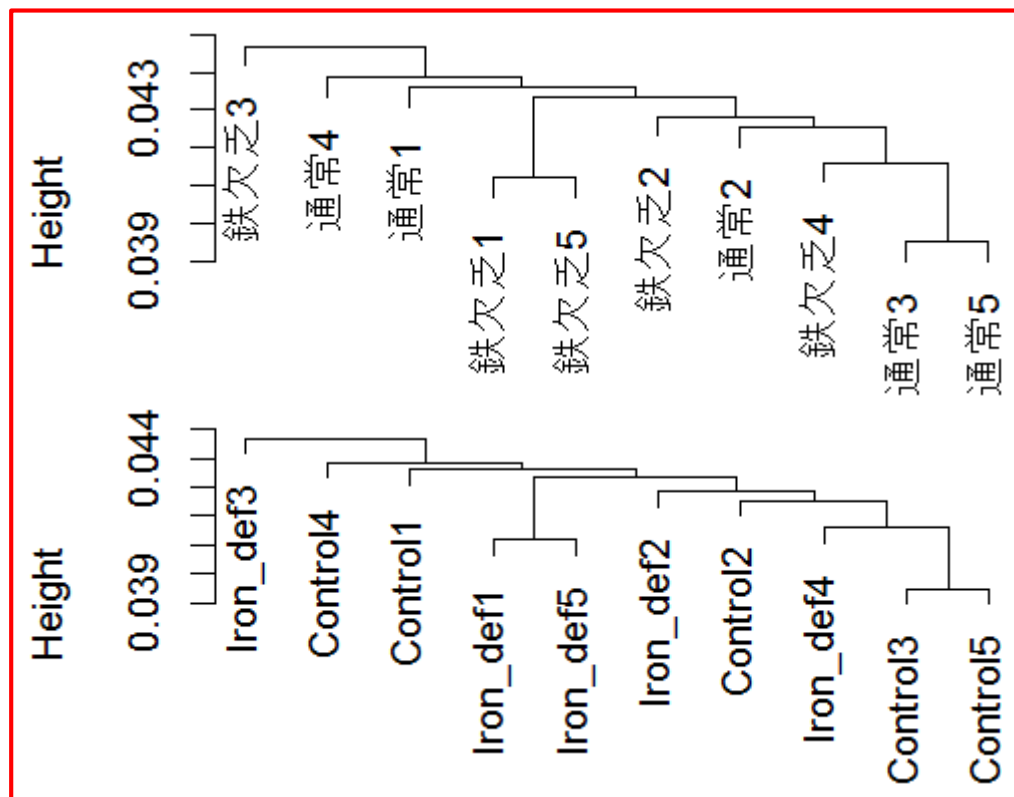
Fig. 1. Comparison of clustering methods. Each sample is labeled as in Table S2. The

GSE30533 (ラット)

①教科書p99の図3-1と基本同じ。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…

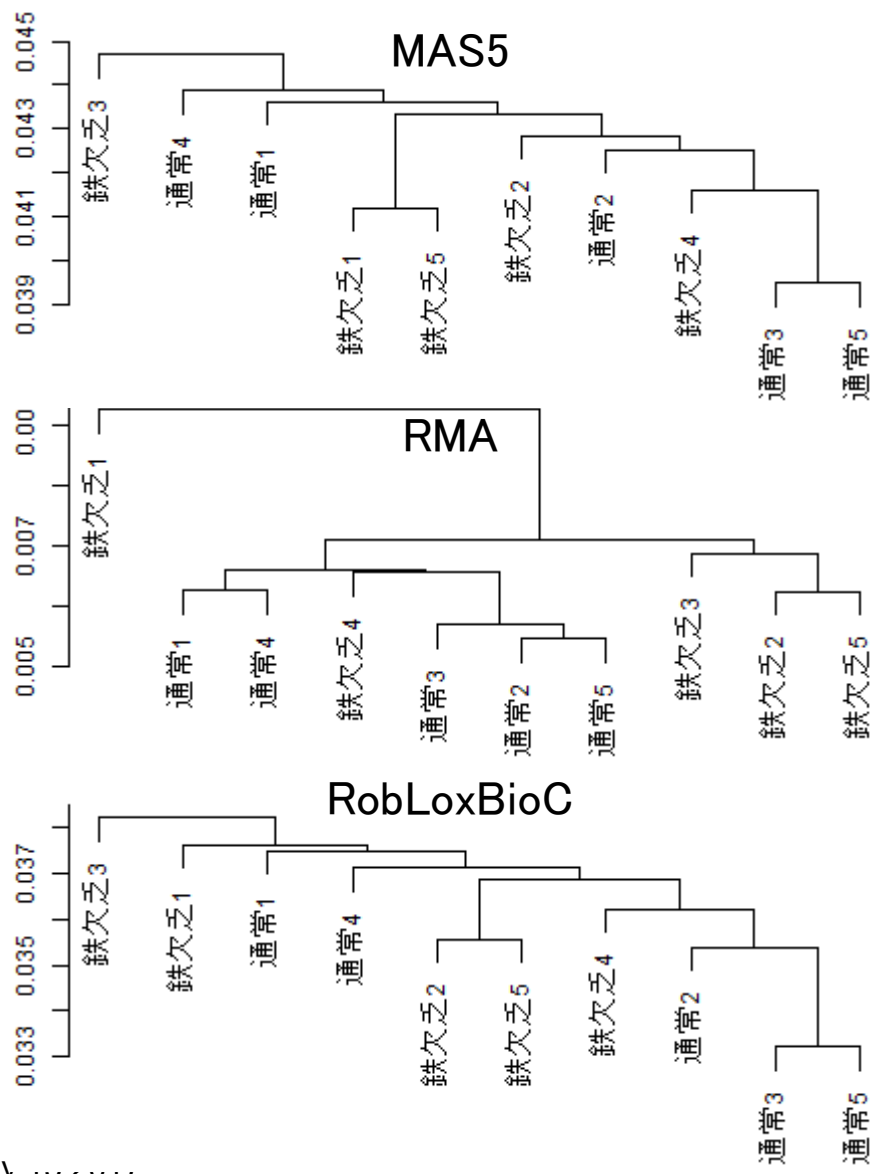
■ Kamei et al., PLoS One, 8: e65732, 2013

- 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

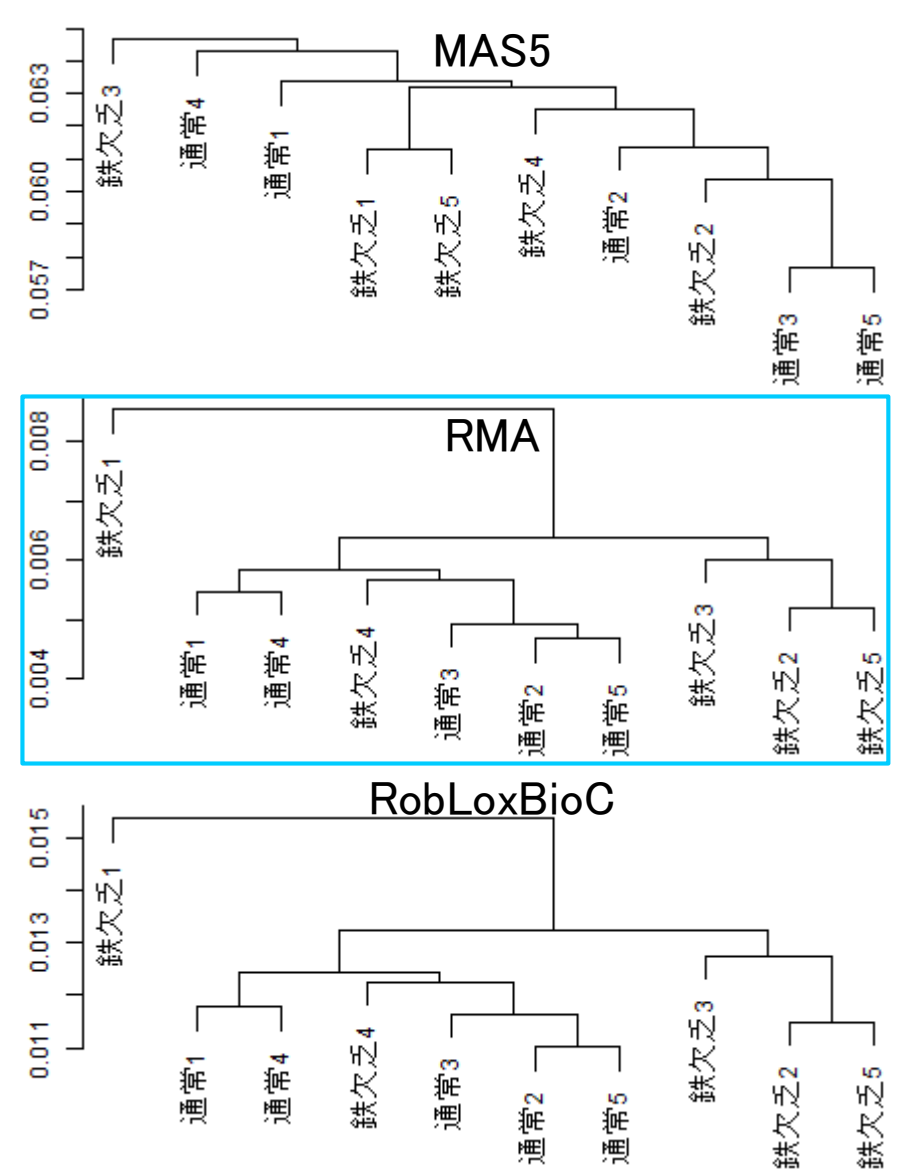


GSE30533 (ラット)

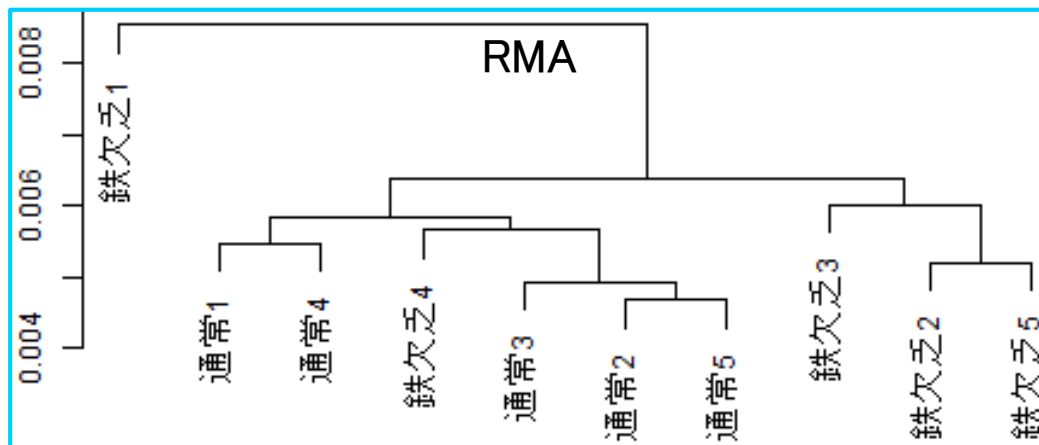
1 - Spearman相関係数



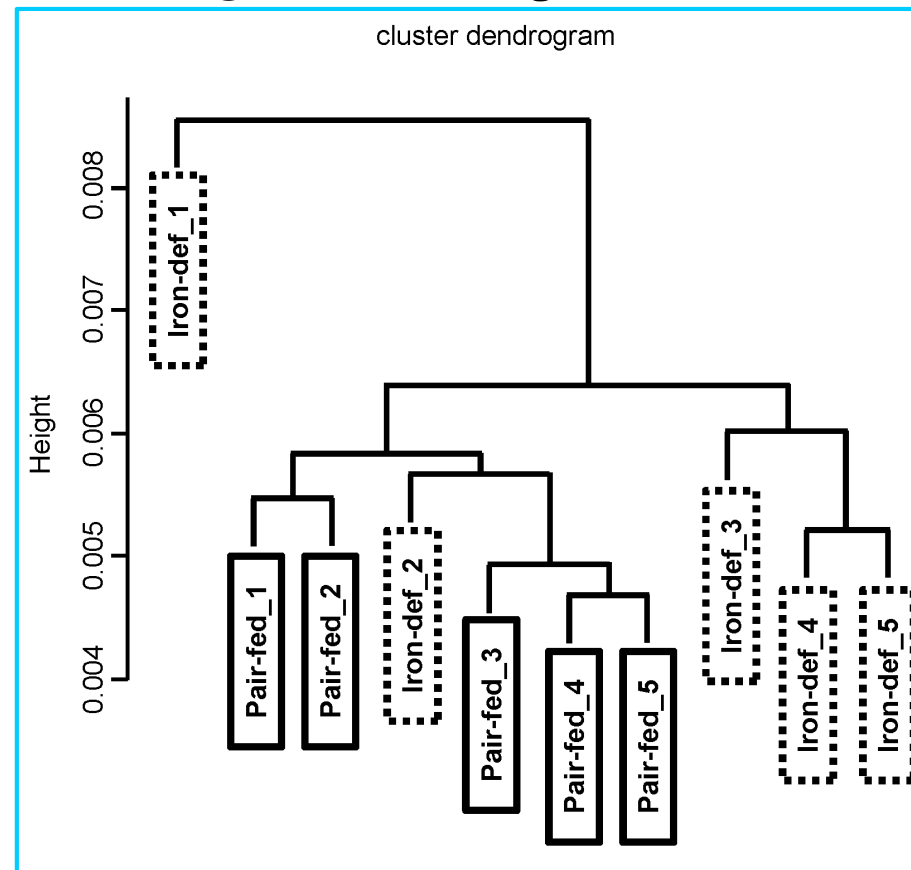
1 - Pearson相関係数



GSE30533 (ラット)



①原著論文のFigure S1



Contents

- 実データ概観
 - GSE2361 (ヒト)、GSE7623 (ラット)、GSE30533 (ラット)
- クラスタリング (教科書の § 3.2.1)
 - 対数変換の有無 (Spearman相関係数を使う場合は気にしなくてよい)
 - 階層的 vs. 非階層的、様々な選択肢
 - 距離の定義: ベクトル間、クラスター間
 - 実データで実行: GSE2361 (ヒト)、課題1
 - 実データで実行: GSE7623 (ラット)、GSE30533 (ラット)
 - 同一プラットフォームデータ(GSE7623 + GSE30533)をマージして実行、課題2
- 実験デザイン (教科書の § 3.2.2)
- 2群間比較: 発現変動遺伝子 (DEG) 検出
 - パターンマッチング法 (相関係数の利用)
 - コードの中身をおさらい、apply関数の基本的な利用法など

同一アレイデータはマージ可能

①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。3'発現アレイを用いることで、他の多くのデータセットとの比較が可能

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



①merge_GSE7623_GSE30533上にGSE7623とGSE30533の計34 CELファイルを入力として前処理法を適用した結果のファイルがあります(見るだけ)。②rcode_preprocessing2.txtが左記のコード

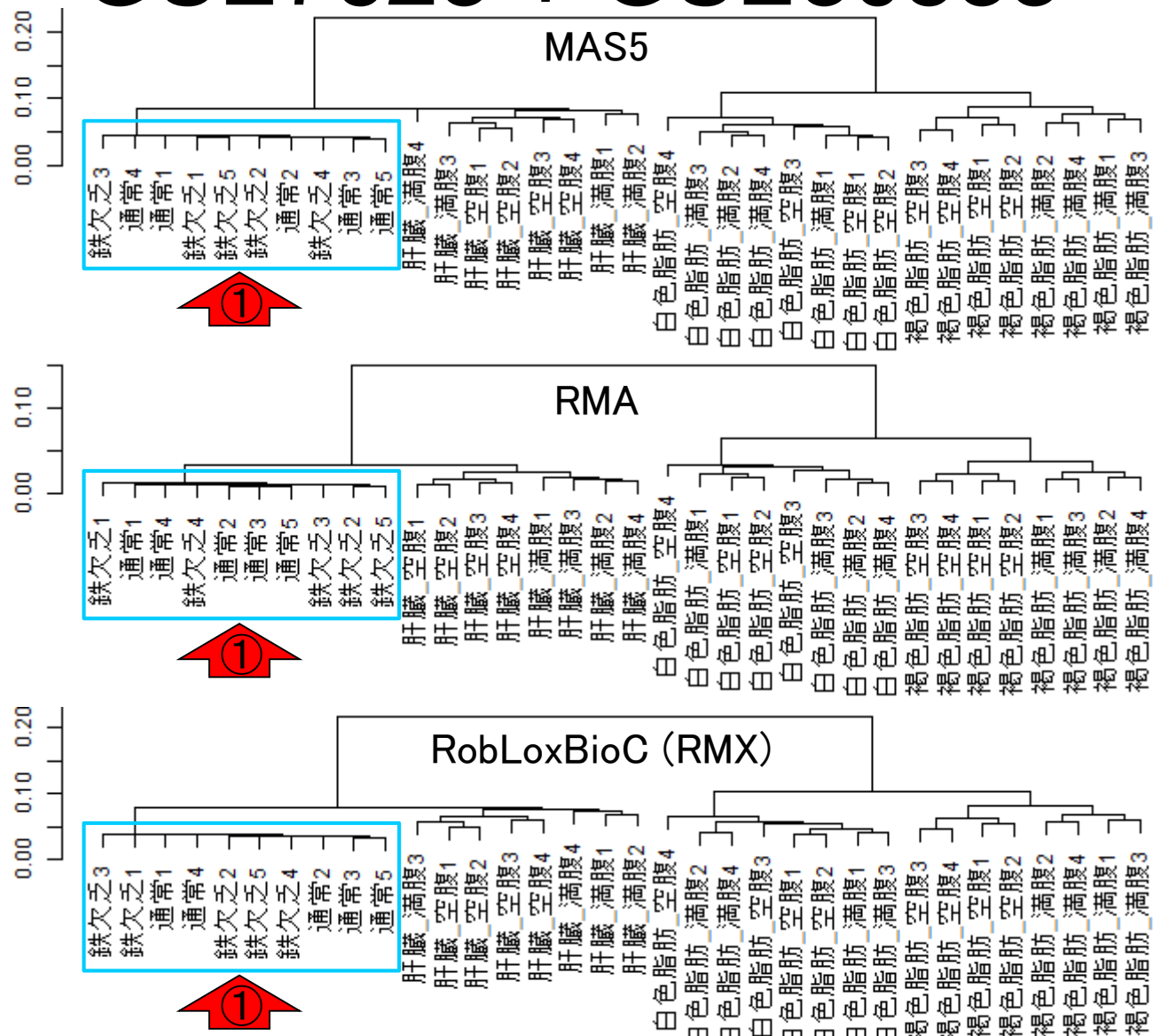


```
#####
### 作業ディレクトリ中のCELファイルの読み込み ###
#####
library(affy) #パッケージの読み込み
hoge <- ReadAffy() #*.CELファイルの読み込み
↓
↓
#####
### RMA前処理法実行 ###
#####
out_f <- "data_rma.txt" #出力
library(affy) #パッ
eset <- rma(hoge) #RMAを
write.exprs(eset, file=out_f) #結果
↓
#####
### MAS5前処理法実行 ###
#####
out_f <- "data_mas.txt"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

名前	更新日時	種類	サイズ
data_mas.txt	2014/05/17 22:34	テキストドキュ...	17,773 KB
data_mas_EN.txt	2014/05/17 22:46	テキストドキュ...	12,612 KB
data_mas_JP.txt	2014/05/17 22:19	テキストドキュ...	12,612 KB
data_rma.txt	2014/05/17 22:28	テキストドキュ...	17,807 KB
data_rma_EN.txt	2014/05/17 23:25	テキストドキュ...	12,643 KB
data_rma_JP.txt	2014/05/17 22:42	テキストドキュ...	12,643 KB
data_rob.txt	2014/05/17 22:35	テキストドキュ...	17,808 KB
data_rob_EN.txt	2014/05/17 23:26	テキストドキュ...	12,644 KB
data_rob_JP.txt	2014/05/17 22:22	テキストドキュ...	12,645 KB
GSM184414.CEL	2010/07/22 23:24	CEL ファイル	6,795 KB
GSM184415.CEL	2010/07/22 23:22	CEL ファイル	6,795 KB
GSM184416.CEL	2010/07/22 23:24	CEL ファイル	6,795 KB
GSM184417.CEL	2010/07/22 23:25	CEL ファイル	6,794 KB
GSM184418.CEL	2010/07/22 23:25	CEL ファイル	6,795 KB
GSM184419.CEL	2010/07/22 23:22	CEL ファイル	6,795 KB
GSM184420.CEL	2010/07/22 23:26	CEL ファイル	6,795 KB
GSM184421.CEL	2010/07/22 23:22	CEL ファイル	6,795 KB

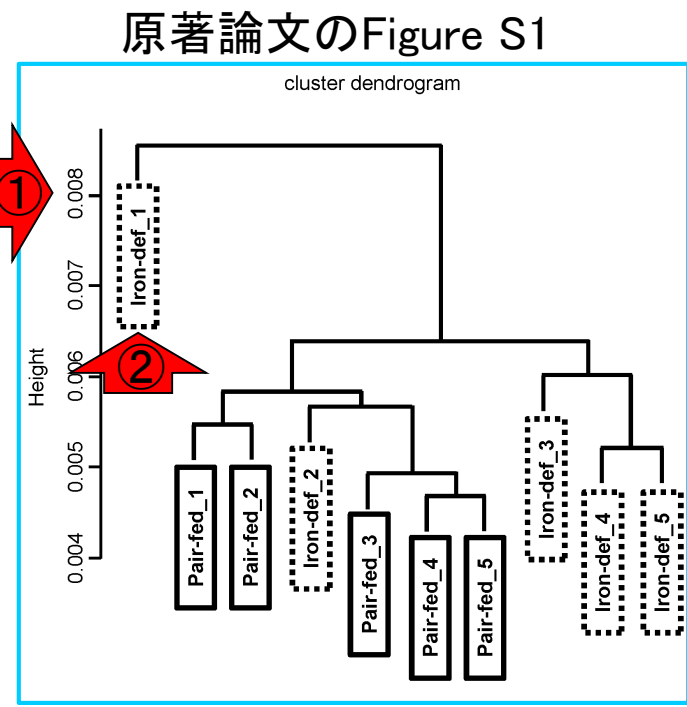
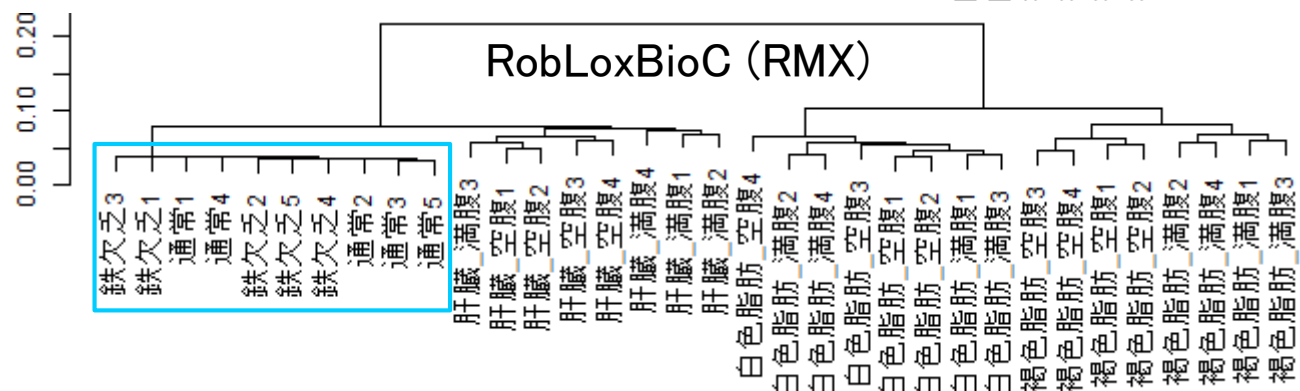
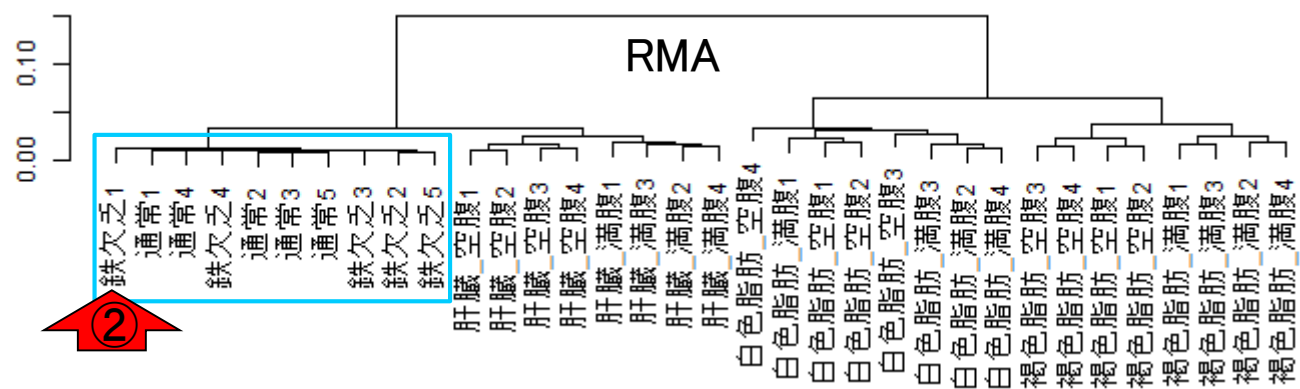
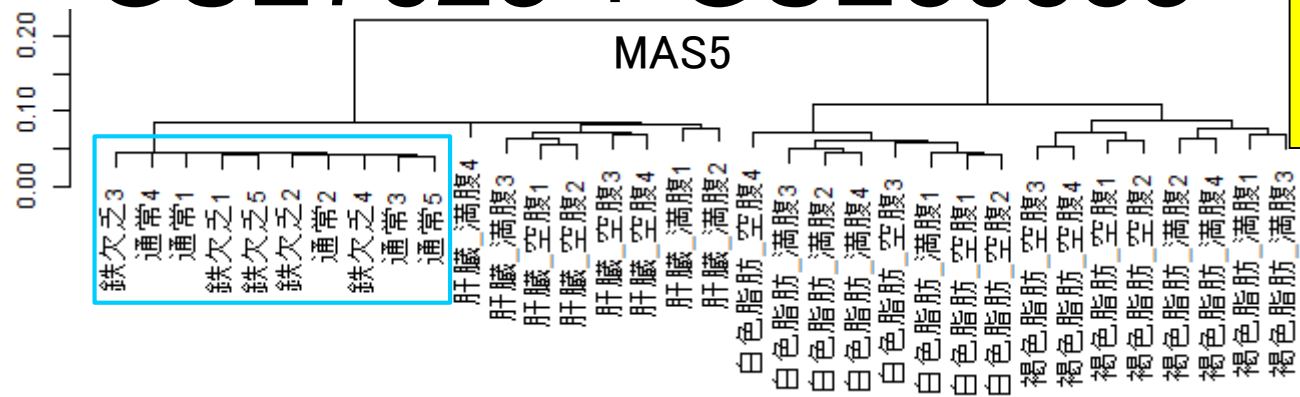
GSE7623 + GSE30533

「1 - Spearman相関係数」の結果。どの前処理法でも似たような結果となっているのが分かる。①ラット10サンプル(通常 対 鉄欠乏)の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)



GSE7623 + GSE30533

マージすることによって初めて、①距離が非常に近い(サンプル間の類似度が極めて高い)がために、「②Iron-def_1が外れサンプルっぽく見える」といった議論をしたことに気づく

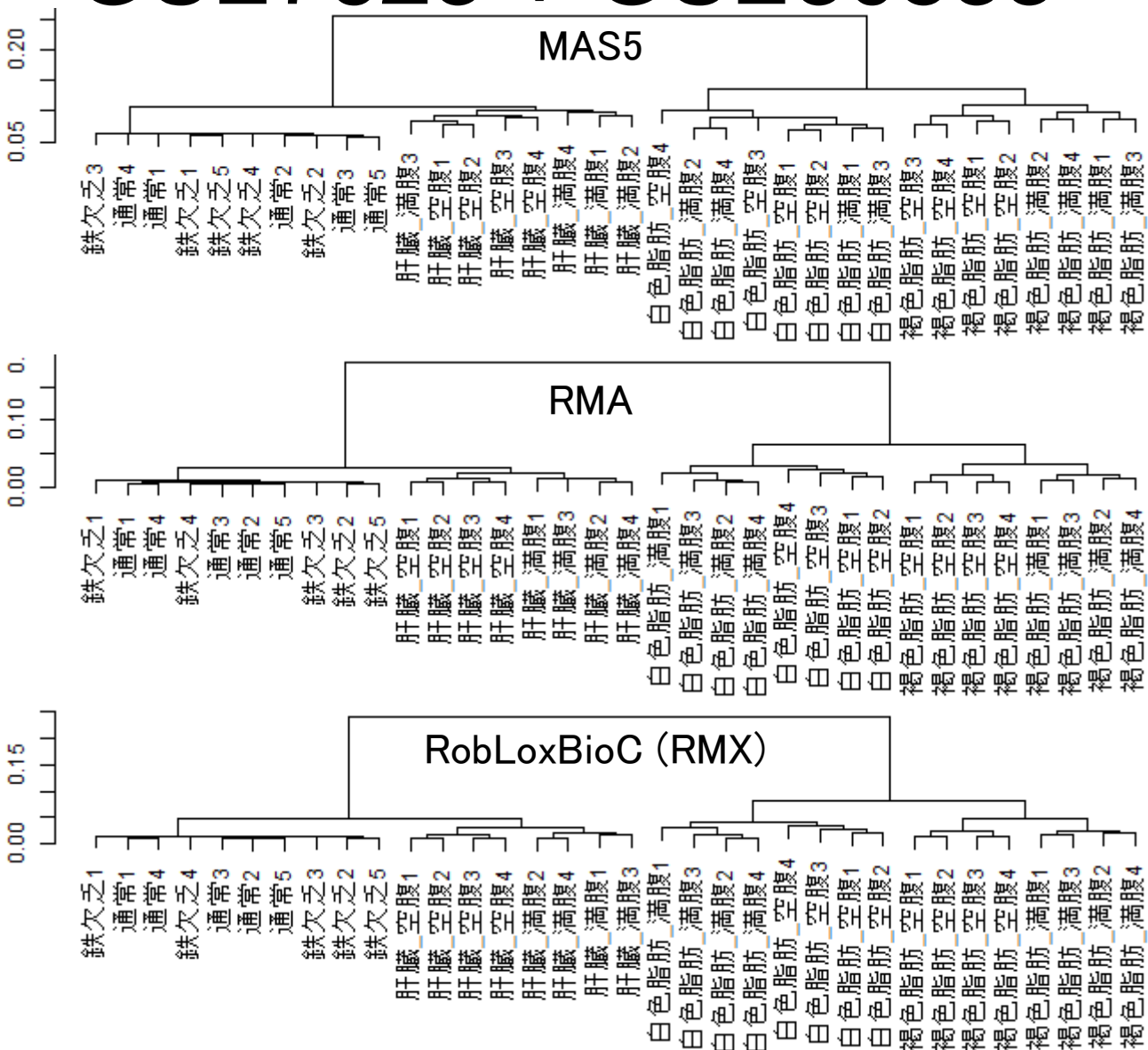


GSE7623 + GSE30533 参考

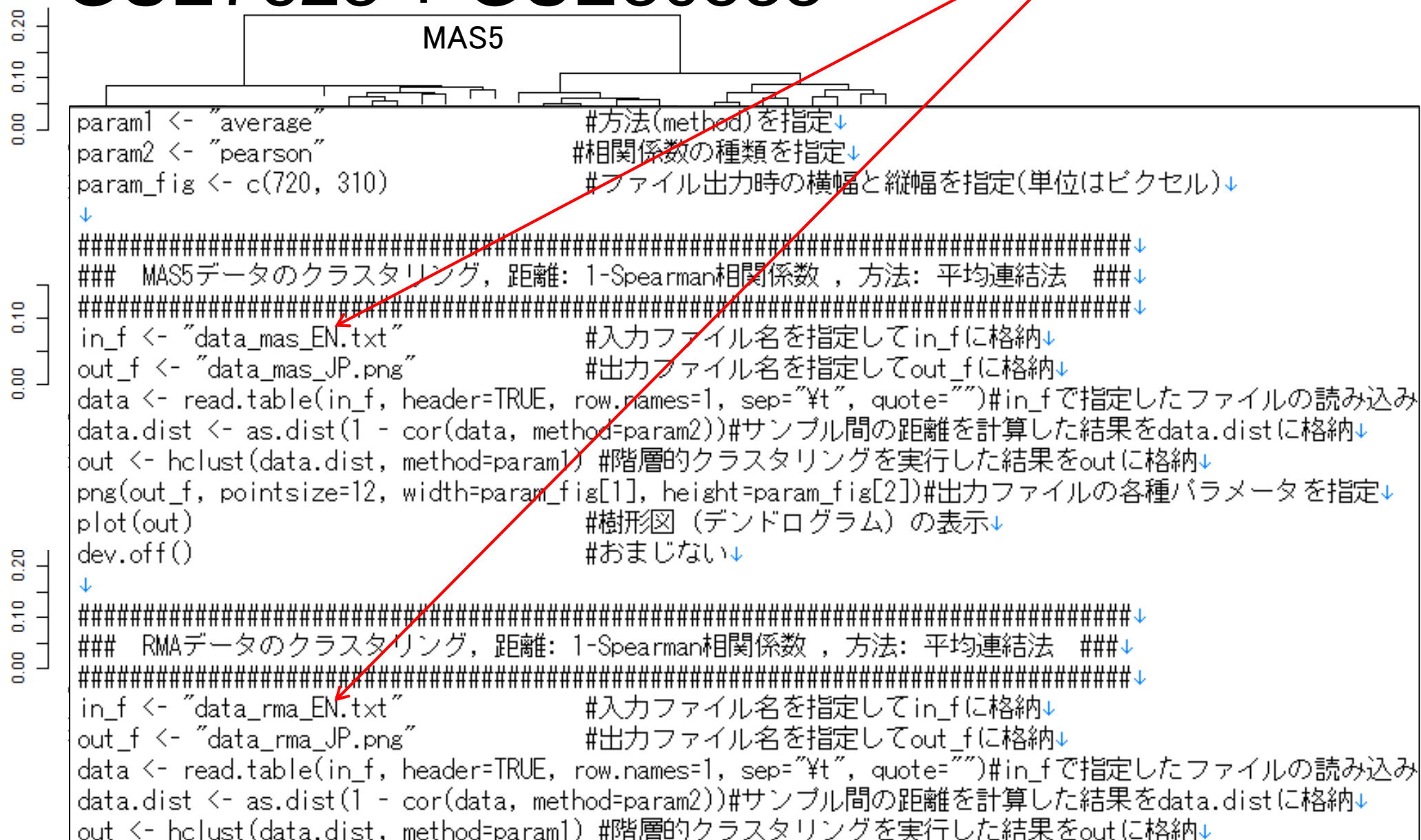


GSE7623 + GSE30533

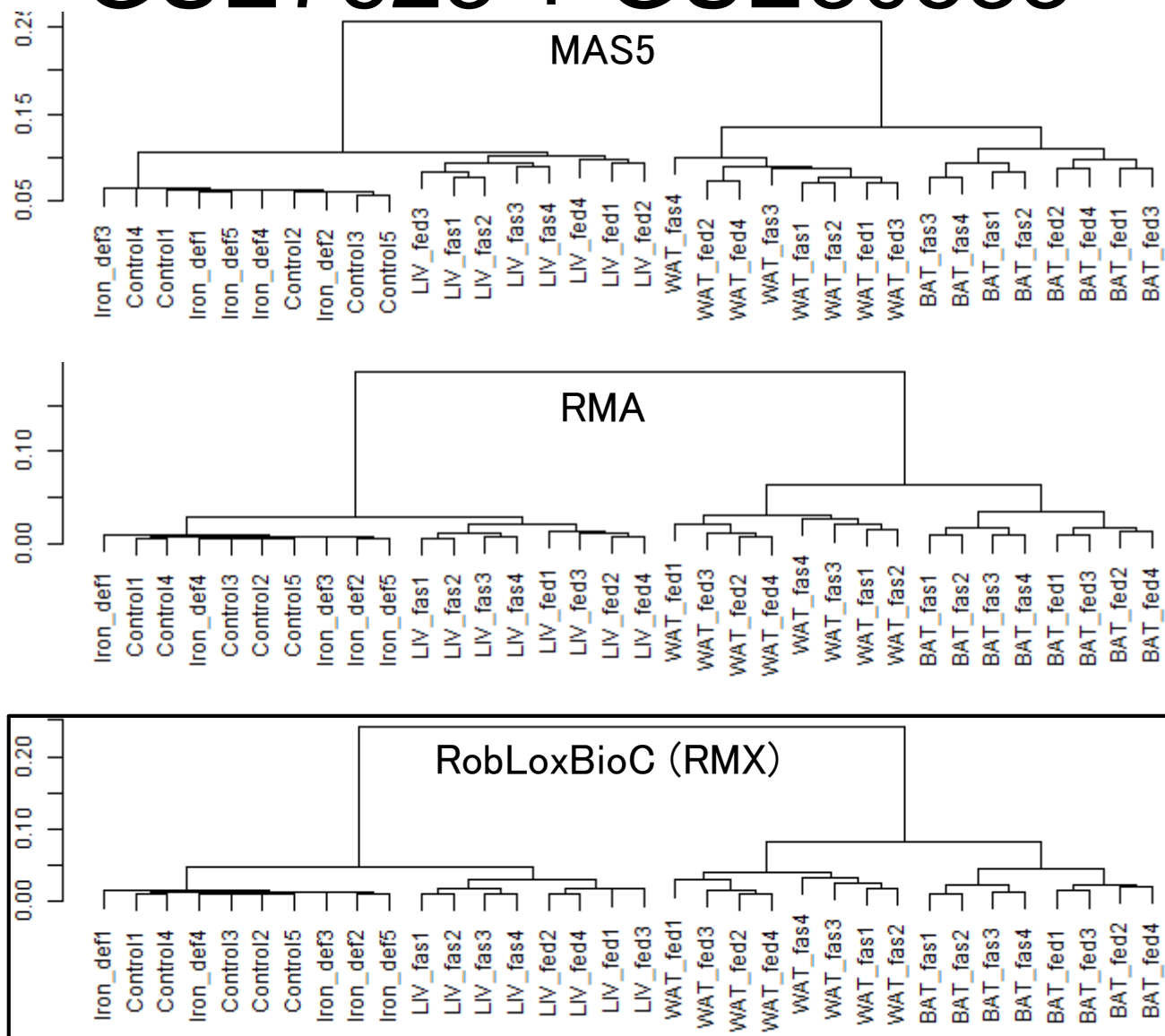
「1 - Pearson相関係数」の結果。
どの前処理法でも似たような結果となっているのが分かる



GSE7623 + GSE30533



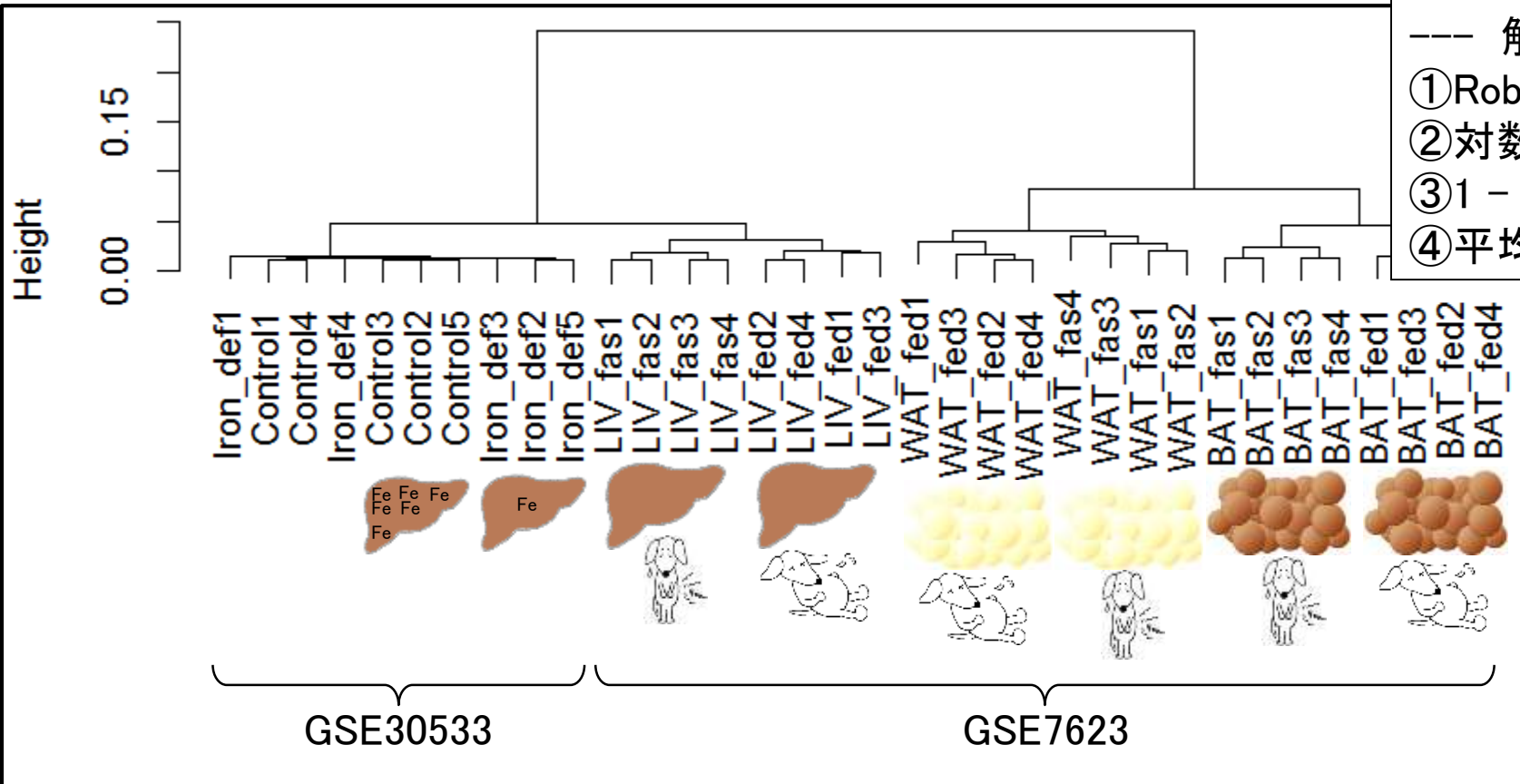
GSE7623 + GSE30533



課題2 (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簡単に考察せよ

- 解析手順 ---
- ① RobLoxBioC前処理法
 - ② 対数変換後のデータ
 - ③ 1 - Pearson相関係数
 - ④ 平均連結法



課題2

- GSE7623とGSE30533は独立した別々の論文
- GSE30533の由来サンプルは?
- 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?
- GSE30533のみでクラスタリングを行った結果のトポロジーは前処理法や距離の定義次第で変わりやすいが…。
- 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)で明瞭に分かれている(MAS5のWATサンプルを除く)。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?

Contents

- 実データ概観
 - GSE2361 (ヒト)、GSE7623 (ラット)、GSE30533 (ラット)
- クラスタリング (教科書の § 3.2.1)
 - 対数変換の有無 (Spearman相関係数を使う場合は気にしなくてよい)
 - 階層的 vs. 非階層的、様々な選択肢
 - 距離の定義: ベクトル間、クラスター間
 - 実データで実行: GSE2361 (ヒト)、課題1
 - 実データで実行: GSE7623 (ラット)、GSE30533 (ラット)
 - 同一プラットフォームデータ(GSE7623 + GSE30533)をマージして実行、課題2
- 実験デザイン (教科書の § 3.2.2)
- 2群間比較: 発現変動遺伝子 (DEG) 検出
 - パターンマッチング法 (相関係数の利用)
 - コードの中身をおさらい、apply関数の基本的な利用法など

実験デザイン (§ 3.2.2)

■ Affymetrix GeneChip

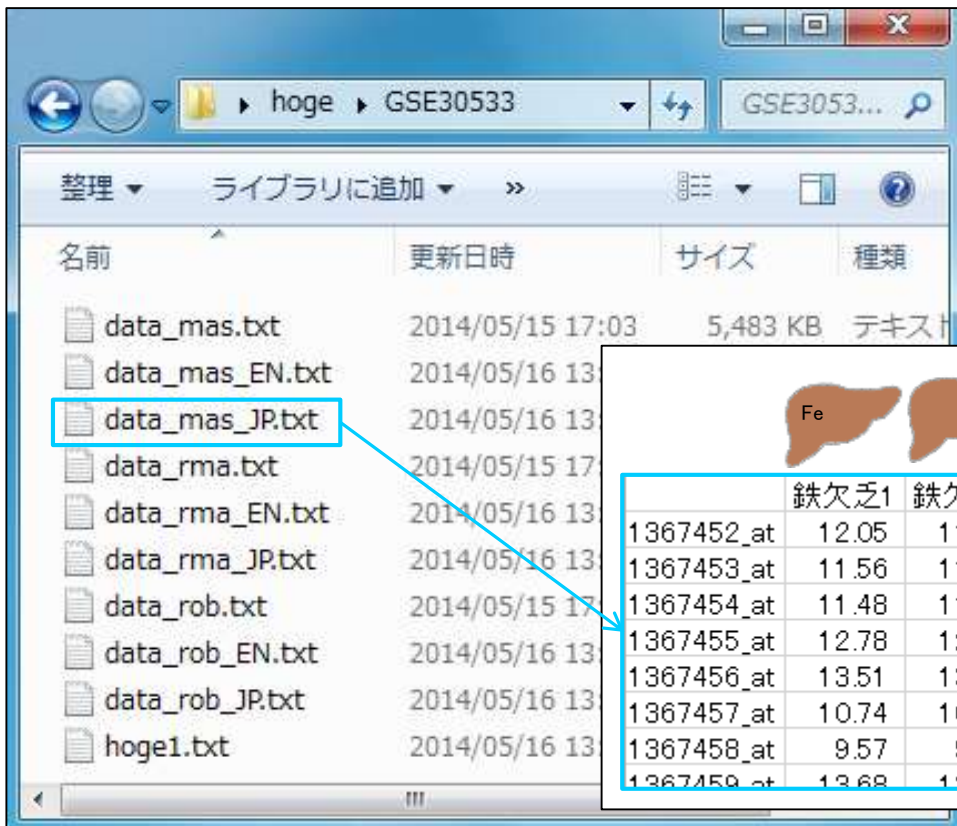
- Ge et al., *Genomics*, **86**: 127–141, 2005
 - GSE2361、GPL96 (Affymetrix Human Genome U133A Array)、22,283 probesets
 - ヒト36サンプル: Heart (心臓)、Thymus (胸腺)、Spleen (脾臓)、Ovary (卵巣)、Kidney (腎臓)、Skeletal Muscle (骨格筋)、Pancreas (膵臓)、Prostate (前立腺)、…
- Nakai et al., *Biosci Biotechnol Biochem.*, **72**: 139–148, 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)

2群間比較が主な目的であり、各群につき5反復(five replicates)とっている。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)

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

- GSE30533、GPL1355 (Affymetrix Rat Genome 230 2.0 Array)、81,000 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



	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

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)

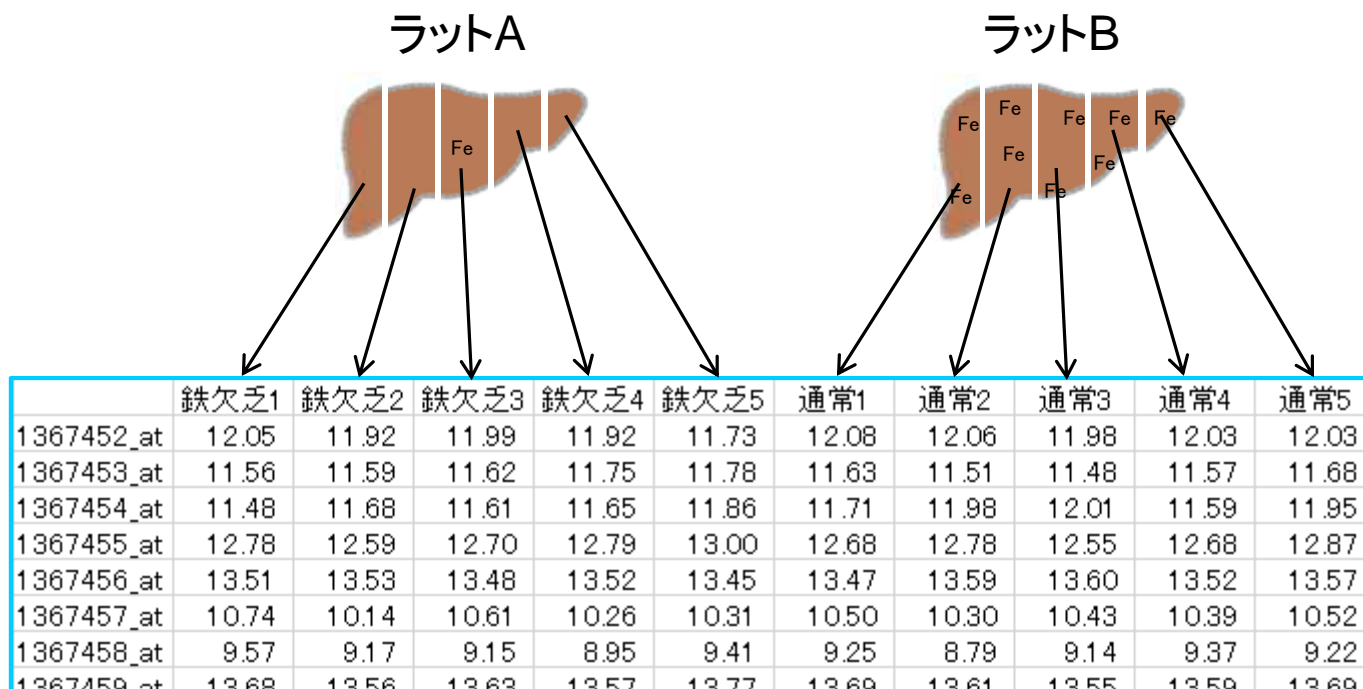
Kamei et al., PLoS One 8: e65732 2013

- 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

対 control diet (control) 5サンプル

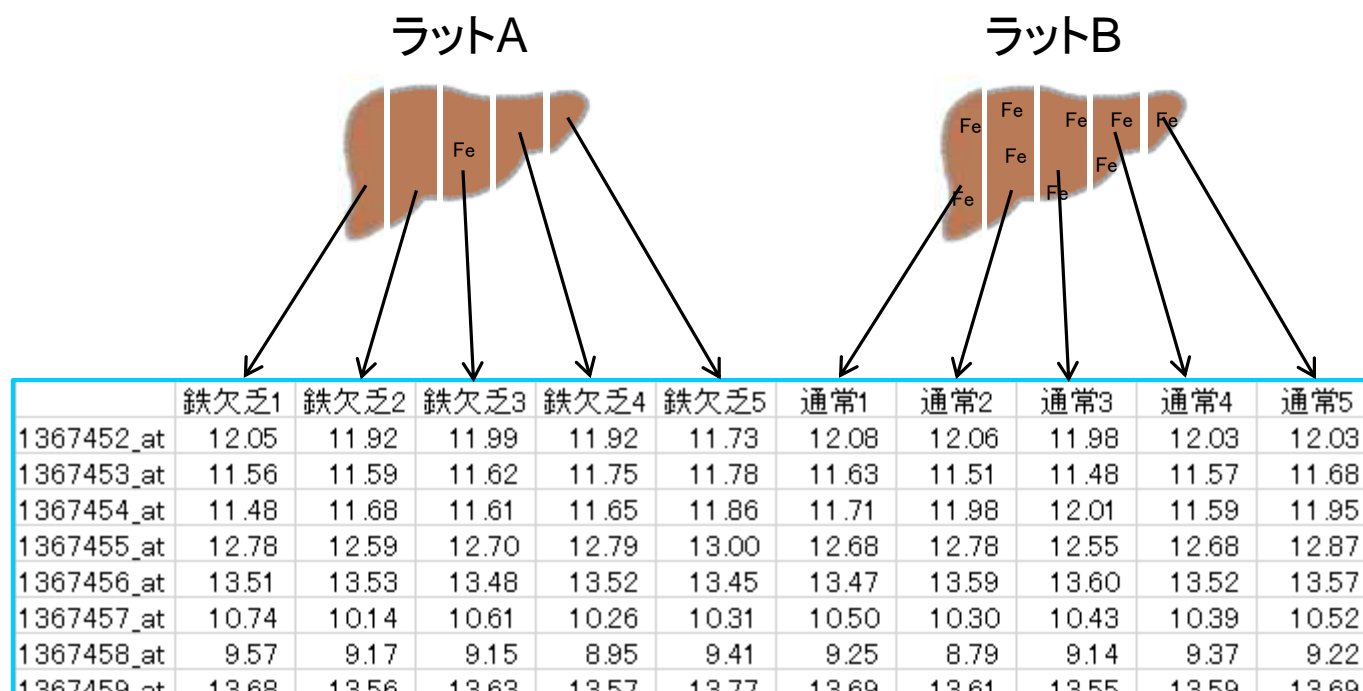


2群間での発現変動遺伝子(DEG) 検出結果は多くなる傾向。多ければいいというものではない!

実験デザイン (§ 3.2.2)

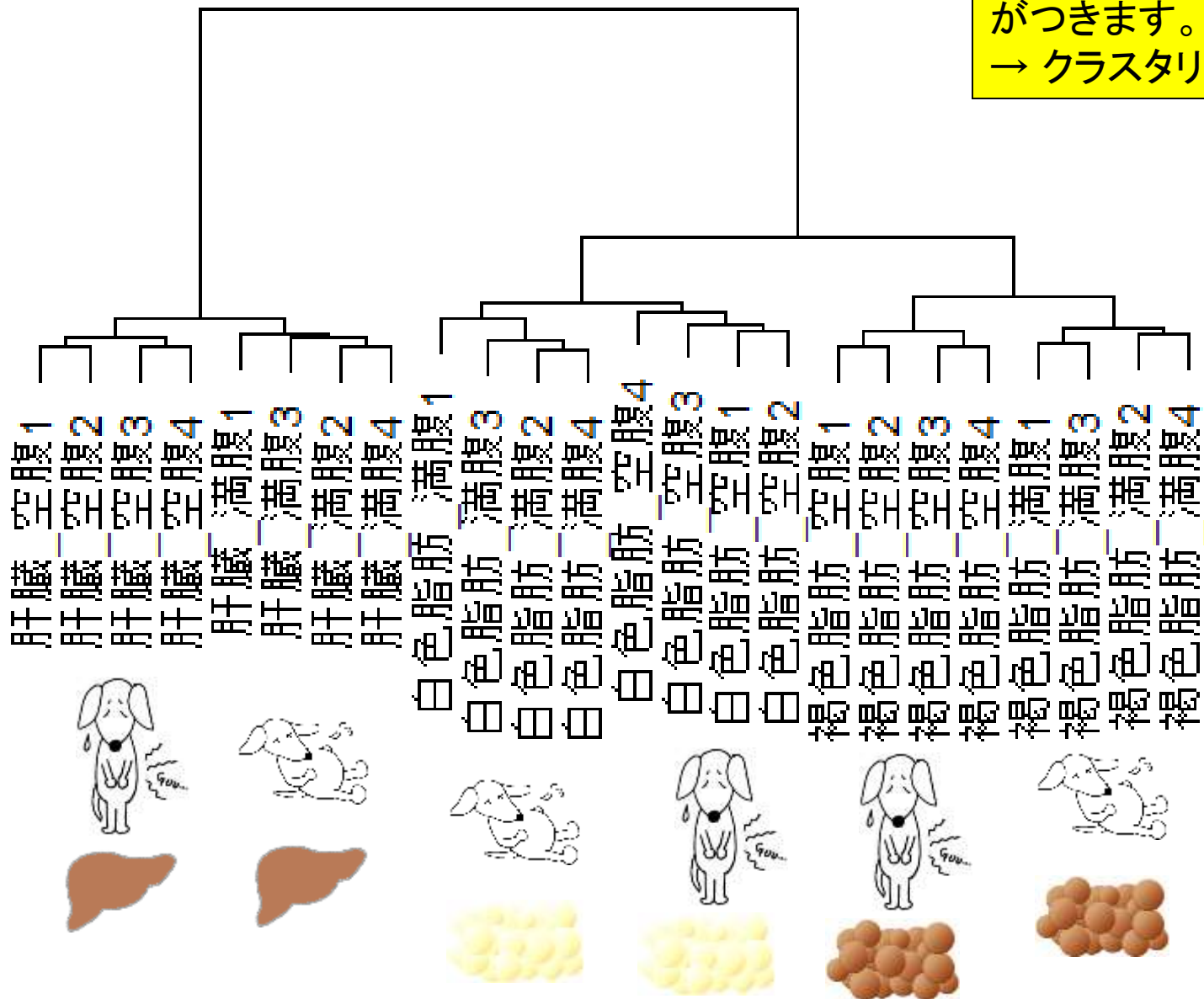
Technical replicatesだと...

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



クラスタリングと発現変動解析

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

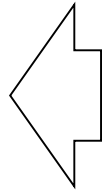
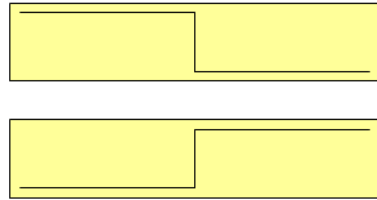


Contents

- 実データ概観
 - GSE2361 (ヒト)、GSE7623 (ラット)、GSE30533 (ラット)
- クラスタリング (教科書の § 3.2.1)
 - 対数変換の有無 (Spearman相関係数を使う場合は気にしなくてよい)
 - 階層的 vs. 非階層的、様々な選択肢
 - 距離の定義: ベクトル間、クラスター間
 - 実データで実行: GSE2361 (ヒト)、課題1
 - 実データで実行: GSE7623 (ラット)、GSE30533 (ラット)
 - 同一プラットフォームデータ(GSE7623 + GSE30533)をマージして実行、課題2
- 実験デザイン (教科書の § 3.2.2)
- 2群間比較: 発現変動遺伝子 (DEG) 検出
 - パターンマッチング法 (相関係数の利用)
 - コードの中身をおさらい、apply関数の基本的な利用法など

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	$x_{n,1}^A$	$x_{n,2}^A$	$x_{n,1}^B$	$x_{n,2}^B$

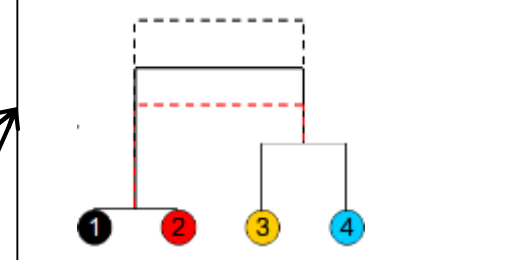
様々な組織(条件)

	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	$x_{n,1}$	$x_{n,2}$	$x_{n,3}$	$x_{n,4}$...

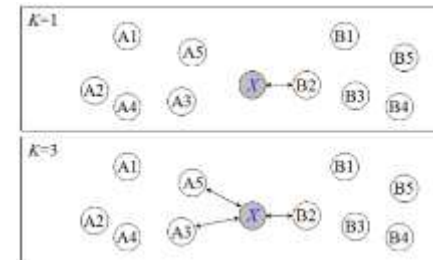
クラスタリング



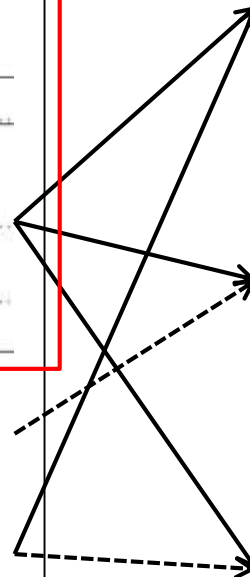
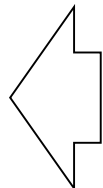
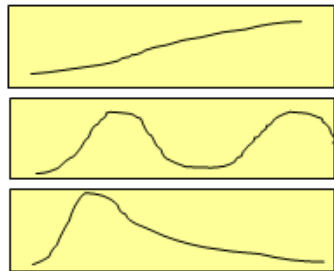
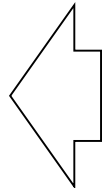
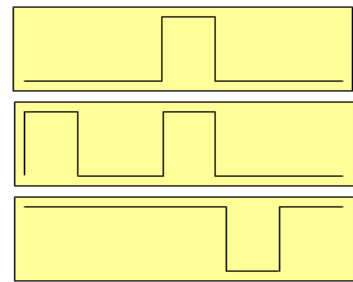
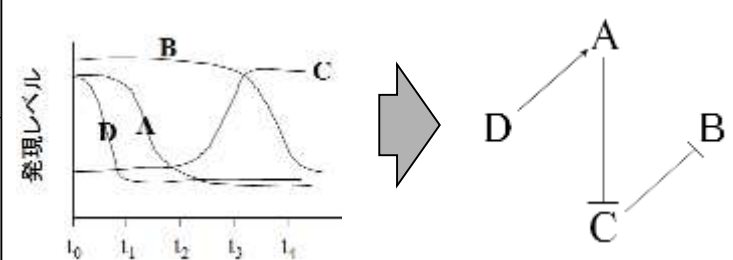
機能解析

- 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

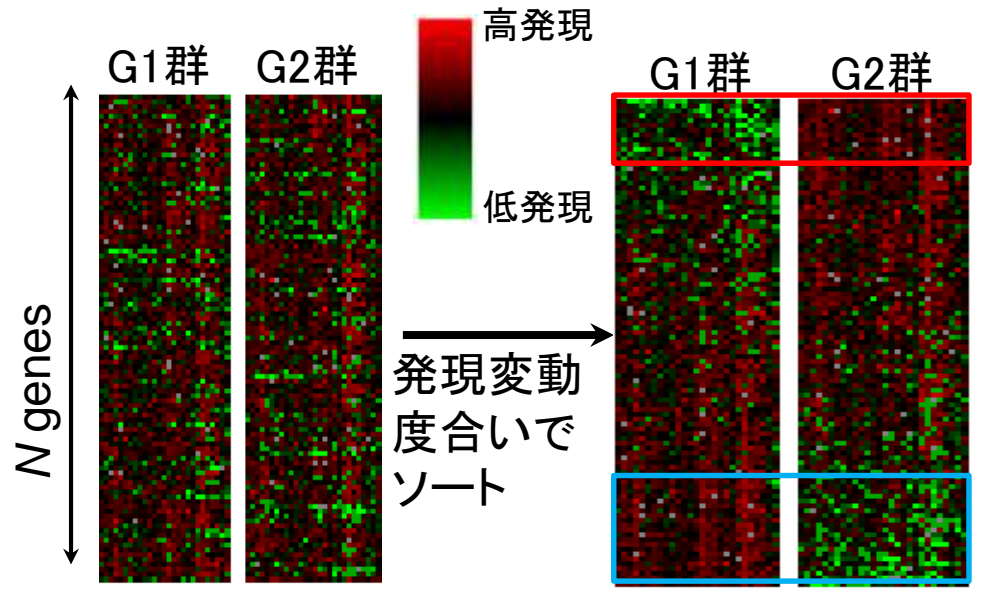


比較したいグループ(群)間で発現変動している遺伝子または転写物を同定することはデータ解析の基本

2群間比較

- 農産物の栽培条件の違い(通常 vs. 低温、通常 vs. 乾燥)
- 味の違い(おいしい vs. まずい)
- サンプルの状態の違い(癌 vs. 正常)

	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	$x_{n,1}^A$	$x_{n,2}^A$	$x_{n,1}^B$	$x_{n,2}^B$



G1群 G2群

2群間で発現の異なる遺伝子 (Differentially Expressed Genes; DEGs)を抽出

2群間比較

①理想的なパターンyと②遺伝子ベクトル間の
③相関係数rを遺伝子ごとに計算。rの絶対値が
大きいほど発現変動の度合いが大きいと解釈

パターンマッチング法

- 理想的なパターンyとの類似度が高い順に遺伝子をランキング

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

	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	$x_{n,1}^A$	$x_{n,2}^A$	$x_{n,1}^B$	$x_{n,2}^B$



y	1	1	1	1	1	1	0	0	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---



rownan	G1_1	G1_2	G1_3	G1_4	G1_5	G1_6	G2_1	G2_2	G2_3	G2_4	G2_5	r
--------	------	------	------	------	------	------	------	------	------	------	------	---



gene1	6.44	6.30	6.51	6.36	6.49	6.39	3.58	4.39	4.25	3.70	4.09	0.98
gene2	5.81	6.93	6.73	5.55	6.39	6.61	2.81	5.46	1.00	3.46	4.17	0.81
gene3	3.91	4.81	5.04	3.17	4.75	5.36	5.58	5.52	5.70	5.64	5.61	-0.71

パターンマッチング法

(Rで)マイクロアレイデータ解析

(last modified 2015/05/16, since 2005)

①パターンマッチング法の、②例題1。
③テンプレートパターン情報を含むファイルを読み込んでパターンマッチングを行ってみよう。まだコピペしなくてよい

What's new?

- ・ 門田幸二
- ・ イ解析に
- ・ どについ
- ・ このペー
- ・ お知らせ
- ・ イ関連の
- ・ ら迎れます
- ・ はじめに

- ・ 解析 | 発現変動 | 2群間 | 対応なし | [Student's t-test](#)(last modified 2014/05/25)
- ・ 解析 | 発現変動 | 2群間 | 対応なし | [Welch t-test](#)(last modified 2014/05/23)
- ・ 解析 | 発現変動 | 2群間 | 対応なし | [Mann-Whitney U-test](#)(last modified 2013/10/15)
- ・ 解析 | 発現変動 | 2群間 | 対応なし | [パターンマッチング法](#) (last modified 2014/05/23)
- ・ 解析 | 発現変動 | 2群間 | 対応あり | [について](#)(last modified 2009/11/11)
- ・ 解析 | 発現変動 | 2群間 | 対応あり | [SAM \(Tusher, 2001\)](#)(last modified 2014/06/02)

解析 | 発現変動 | 2群間 | 対応なし | パターンマッチング法

パターンマッチング法を用いて、2群間での発現変動遺伝子の同定を行うやり方を紹介します。
「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、以下をコピペ

1. サンプルデータ16のsample16_log.txt(対数変換後のデータ)の場合:

クラスラベル情報ファイル([sample16_cl.txt](#))を利用するやり方です。

```

in_f1 <- "sample16_log.txt" #入力ファイル名を指定してin_f1に格納(発現データ)
in_f2 <- "sample16_cl.txt" #入力ファイル名を指定してin_f2に格納(テンプレート)
out_f <- "hoge1.txt" #出力ファイル名を指定してout_fに格納
param <- "pearson" #相関係数の種類を指定("pearson"または"spearman")

```

#入力ファイルの読み込みとラベル情報の作成

```

data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t", quote="")#in_f1で指定したファイルの読み込み
hoge <- read.table(in_f2, sep="\t", quote="")#in_f2で指定したファイルの読み込み
data.cl <- hoge[,2] #テンプレートパターンベクトルdata.clを作成

```

#本番

```

r <- apply(data, 1, cor, y=data.cl, method=param)#各(行)遺伝子についてテンプレートパターンとの相関係数rを計算

```

#ファイルに保存

```

tmp <- cbind(rownames(data), data, r) #入力データの右側に相関係数rのベクトルを結合した結果を作成
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に保存

```

入出力の関係

①テンプレートパターン情報ファイル。2列目がテンプレートパターン。②(log変換後の)発現データファイル。赤枠で示すように、列の並びは同じ

解析 | 発現変動 | 2群間 | 対応なし | パターンマッチング法

パターンマッチング法を用いて、2群間での発現変動遺伝子の同定を行うやり方を紹介します。
「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16のsample16_log.txt(対数変換後のデータ)の場合:

クラスラベル情報ファイル([sample16_cl.txt](#))を利用するやり方です。

```
in_f1 ② "sample16_log.txt" ① #入力ファイル名を指定してin_f1に格納(発現データ)
in_f2  "sample16_cl.txt"    #入力ファイル名を指定してin_f2に格納(テンプレート)
out_f  <- "hoge1.txt"       #出力ファイル名を指定してout_fに格納
param <- "pearson"         #相関係数の種類を指定("pearson"または"spearman")

#入力ファイルの読み込みとラベル情報の作成
data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t", quote="")#in_f1で指定し
hoge <- read.table(in_f2, sep="\t", quote="")#in_f2で指定したファイルの読み込み
data.cl <- hoge[,2]          #テンプレートパターンベクトルdata.clを作成

#本番
r <- apply(data, 1, cor, y=data.cl, method=param)#各(行)遺伝子についてテンプレートパタ

#ファイルに保存
tmp <- cbind(rownames(data), data, r) #入力データの右側に相関係数rのベクトルを結合した結
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F)#tmpの中身を指定した
```

G1_1	1
G1_2	1
G1_3	1
G1_4	1
G1_5	1
G1_6	1
G2_1	0
G2_2	0
G2_3	0
G2_4	0
G2_5	0



rowname	G1_1	G1_2	G1_3	G1_4	G1_5	G1_6	G2_1	G2_2	G2_3	G2_4	G2_5
gene1	6.44	6.30	6.51	6.36	6.49	6.39	3.58	4.39	4.25	3.70	4.09
gene2	5.81	6.93	6.73	5.55	6.39	6.61	2.81	5.46	1.00	3.46	4.17
gene3	3.91	4.81	5.04	3.17	4.75	5.36	5.58	5.52	5.70	5.64	5.61



入出力の関係

解析 | 発現変動 | 2群間 | 対応なし | パターンマッチング法

パターンマッチング法を用いて、2群間での発現変動遺伝子の同定を行うやり方を紹介します。
「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16のsample16_log.txt(対数変換後のデータ)の場合:

クラスラベル情報ファイル([sample16_cl.txt](#))を利用するやり方です。

```
in_f1 ② "sample16_log.txt"      #入力ファイル名を指定してin_f1に格納(発現データ)
in_f2  "sample16_cl.txt"      #入力ファイル名を指定してin_f2に格納(テンプレート)
out_f  <- "hoge1.txt"         ① #出力ファイル名を指定してout_fに格納
param <- "pearson"           #相関係数の種類を指定("pearson"または"spearman")

#入力ファイルの読み込みとラベル情報の作成
data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t", quote="")#in_f1で指定し
hoge <- read.table(in_f2, sep="\t", quote="")#in_f2で指定したファイルの読み込み
data.cl <- hoge[,2]          #テンプレートパターンベクトルdata.clを作成

#本番
r <- apply(data, 1, cor, y=data.cl, method=param)#各(行)遺伝子についてテンプレートパタ

#ファイルに保存
tmp <- cbind(rownames(data), data, r) #入力データの右側に相関係数rのベクトルを結合した結
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F)#tmpの中身を指定した
```

G1_1	1
G1_2	1
G1_3	1
G1_4	1
G1_5	1
G1_6	1
G2_1	0
G2_2	0
G2_3	0
G2_4	0
G2_5	0



rowname	G1_1	G1_2	G1_3	G1_4	G1_5	G1_6	G2_1	G2_2	G2_3	G2_4	G2_5	r
gene1	6.44	6.30	6.51	6.36	6.49	6.39	3.58	4.39	4.25	3.70	4.09	0.984
gene2	5.81	6.93	6.73	5.55	6.39	6.61	2.81	5.46	1.00	3.46	4.17	0.811
gene3	3.91	4.81	5.04	3.17	4.75	5.36	5.58	5.52	5.70	5.64	5.61	-0.708

黒枠内をコピー。ここは、①入力ファイル情報読み込み部分。②dataオブジェクトを確認

コードの解説

1. サンプルデータ16の sample16_log.txt(対数変換後のデータ)の場合:

クラスラベル情報ファイル(sample16_cl.txt)を利用するやり方です。

```

in_f1 <- "sample16_log.txt"
in_f2 <- "sample16_cl.txt"
out_f <- "hoge1.txt"
param <- "pearson"

```

#入力ファイル名を指定してin_f1に格納(発現データ)
 #入力ファイル名を指定してin_f2に格納(テンプレート情報)
 #出力ファイル名を指定してout_fに格納
 #相関係数の種類を指定("pearson"または"spearman")

#入力ファイルの読み込みとラベル情報の作成

```

data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t")
hoge <- read.table(in_f2, sep="\t")
data.cl <- hoge[,2]

```

#本番

```
r <- apply(data, 1, cor, y=data.cl)
```

#ファイルに保存

```

tmp <- cbind(rownames(data), data)
write.table(tmp, out_f, sep="\t")

```

```

> getwd()
[1] "C:/Users/kadota/Desktop/hoge"
> list.files(pattern="16")
[1] "sample16 cl.txt" "sample16_log.txt"
> in_f1 <- "sample16_log.txt"
> in_f2 <- "sample16_cl.txt"
> out_f <- "hoge1.txt"
> param <- "pearson"
>
> #入力ファイルの読み込みとラベル情報の作成
> data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t")
> data
      G1_1 G1_2 G1_3 G1_4 G1_5 G1_6 G2_1 G2_2 G2_3 G2_4 G2_5
gene1 6.44 6.30 6.51 6.36 6.49 6.39 3.58 4.39 4.25 3.70 4.09
gene2 5.81 6.93 6.73 5.55 6.39 6.61 2.81 5.46 1.00 3.46 4.17
gene3 3.91 4.81 5.04 3.17 4.75 5.36 5.58 5.52 5.70 5.64 5.61
> |

```


コードの解説

1. サンプルデータ16の sample16_log.txt(対数変換後のデータ)の場合:

クラスラベル情報ファイル(sample16_cl.txt)を利用するやり方です。

```

in_f1 <- "sample16_log.txt"      #入力ファイル名を指定してin_f1に格納(発現データ)
in_f2 <- "sample16_cl.txt"      #入力ファイル名を指定してin_f2に格納(テンプレート情報)
out_f <- "hoge1.txt"            #出力ファイル名を指定してout_fに格納
param <- "pearson"              #相関係数の種類を指定("pearson"または"spearman")

#入力ファイルの読み込みとラベル情報の作成
data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t", quote="")#in_f1で指定したファイルの読み込み
hoge <- read.table(in_f2, sep="\t", quote="")#in_f2で指定したファイルの読み込み
data.cl <- hoge[,2]              #テンプレートパターンベクトルdata.clを作成

#本番
r <- apply(data, 1, cor, y=data.cl, method=param)#各(行)遺伝子についてテンプレートパターンdata.clとの相関

#ファイルに保存
tmp <- cbind(row.names(data), data)
write.table(tmp, out_f, sep="\t", as.is=TRUE)
    
```

G1_1	1
G1_2	1
G1_3	1
G1_4	1
G1_5	1
G1_6	1
G2_1	0
G2_2	0
G2_3	0
G2_4	0
G2_5	0

R Console

```

> data
      G1_1 G1_2 G1_3 G1_4 G1_5 G1_6 G2_1 G2_2 G2_3 G2_4 G2_5
gene1 6.44 6.30 6.51 6.36 6.49 6.39 3.58 4.39 4.25 3.70 4.09
gene2 5.81 6.93 6.73 5.55 6.39 6.61 2.81 5.46 1.00 3.46 4.17
gene3 3.91 4.81 5.04 3.17 4.75 5.36 5.58 5.52 5.70 5.64 5.61

> hoge <- read.table(in_f2, sep="\t", quote="")#in_f2で指定$
> data.cl <- hoge[,2]              #テンプレートパター$
> data.cl
 [1] 1 1 1 1 1 1 0 0 0 0 0
> |
    
```

①テンプレートパターン情報を、②の部分で読み込んで得られた、③hogeの中身は入力ファイルと同じだが…

コードの解説

1. サンプルデータ16の sample16_log.txt(対数変換後のデータ)の場合:

クラスラベル情報ファイル(sample16_cl.txt)を利用するやり方です。

```

in_f1 <- "sample16_log.txt"
in_f2 <- "sample16_cl.txt"
out_f <- "hoge1.txt"
param <- "pearson"

#入力ファイルの読み込みとラベル情報の作成
data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t")
hoge <- read.table(in_f2, sep="\t", quote="")
data.cl <- hoge[,2]

#本番
r <- apply(data, 1, cor, y=data.cl, method=param)#各(行)遺伝子間の相関係数

#ファイルに保存
tmp <- cbind(rownames(data), data, r) #入力データの右側に相関係数を追加
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=TRUE)
    
```

#入力ファイル名を指定してin_f1に格納(発現データ)
 #入力ファイル名を指定してin_f2に格納(テンプレート情報)
 #出力ファイル名を指定してout_fに格納
 #相関係数の種類を指定("pearson")

```

R Console
> hoge
      V1 V2
1  G1_1  1
2  G1_2  1
3  G1_3  1
4  G1_4  1
5  G1_5  1
6  G1_6  1
7  G2_1  0
8  G2_2  0
9  G2_3  0
10 G2_4  0
11 G2_5  0
> dim(hoge)
[1] 11  2
> hoge[3, ]
      V1 V2
3  G1_3  1
> hoge[, 2]
[1] 1 1 1 1 1 1 0 0 0 0 0
> |
    
```

G1_1	1
G1_2	1
G1_3	1
G1_4	1
G1_5	1
G1_6	1
G2_1	0
G2_2	0
G2_3	0
G2_4	0
G2_5	0



①欲しいのはhoge行列の2列目部分のみなので、②その部分のみ抽出

コードの解説

1. サンプルデータ16の sample16_log.txt(対数変換後のデータ)の場合:

クラスラベル情報ファイル(sample16_cl.txt)を利用するやり方です。

```
in_f1 <- "sample16_log.txt"
in_f2 <- "sample16_cl.txt"
out_f <- "hoge1.txt"
param <- "pearson"
```

#入力ファイル名を指定してin_f1に格納(発現データ)
 #入力ファイル名を指定してin_f2に格納(テンプレート情報)
 #出力ファイル名を指定してout_fに格納
 #相関係数の種類を指定("p

```
#入力ファイルの読み込みとラベル情報の作成
data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t",
hoge <- read.table(in_f2, sep="\t", quote="")#in_f2で指定したフ
data.cl <- hoge[,2]
```

#テンプレートパターンベク

```
#本番
r <- apply(data, 1, cor, y=data.cl, method=param)#各(行)遺伝
```

```
#ファイルに保存
tmp <- cbind(rownames(data), data, r) #入力データの右側に相関係係
write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=
```

```
R Console
> hoge
      V1 V2
1  G1_1  1
2  G1_2  1
3  G1_3  1
4  G1_4  1
5  G1_5  1
6  G1_6  1
7  G2_1  0
8  G2_2  0
9  G2_3  0
10 G2_4  0
11 G2_5  0
> dim(hoge)
[1] 11  2
> hoge[3, ]
      V1 V2
3  G1_3  1
> hoge[, 2]
[1] 1 1 1 1 1 1 0 0 0 0 0
> |
```

G1_1	1
G1_2	1
G1_3	1
G1_4	1
G1_5	1
G1_6	1
G2_1	0
G2_2	0
G2_3	0
G2_4	0
G2_5	0

み
関

もちろん、読み込み時に①row.names=1をつけてもよい。その場合は黒下線部分を②2から1に変更

Tips

```
in_f1 <- "sample16_log.txt"
in_f2 <- "sample16_cl.txt"
out_f <- "hoge1.txt"
param <- "pearson"
```

```
#入力ファイル名を指定してin_f1に格納(発現データ)
#入力ファイル名を指定してin_f2に格納(テンプレート情報)
#出力ファイル名を指定してout_fに格納
#相関係数の種類を指定("pearson"または"spearman")
```

#入力ファイルの読み込みとラベル情報の作成

```
data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t", quote="")#in_f1で指定したファイルの読み込み
hoge <- read.table(in_f2, sep="\t", quote="")#in_f2で指定したファイルの読み込み
data.cl <- hoge[,2]
```

```
#本番
r <- apply(data, 1, cor,
```

#ファイルに保存

```
tmp <- cbind(rownames(data), r)
write.table(tmp, out_f, sep="\t", quote="")
```

```
R Console
> hoge <- read.table(in_f2, row.names=1, sep="\t", quote="")
> hoge
      V2
G1_1  1
G1_2  1
G1_3  1
G1_4  1
G1_5  1
G1_6  1
G2_1  0
G2_2  0
G2_3  0
G2_4  0
G2_5  0
> data.cl <- hoge[, 1]
> data.cl
[1] 1 1 1 1 1 1 0 0 0 0 0
> |
```

G1_1	1
G1_2	1
G1_3	1
G1_4	1
G1_5	1
G1_6	1
G2_1	0
G2_2	0
G2_3	0
G2_4	0
G2_5	0



コードの解説: apply

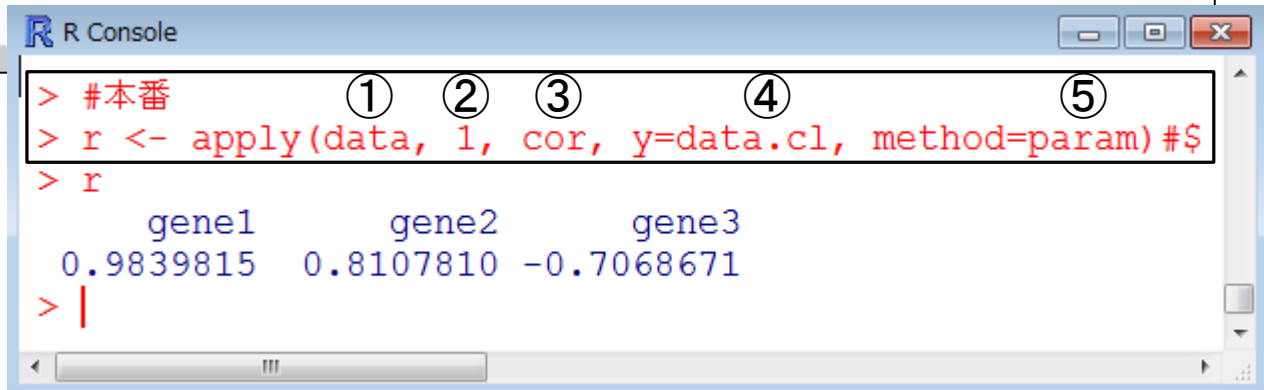
apply関数は行ごとや列ごとに同じ関数を繰り返し実行させたい場合に便利。①dataオブジェクトの、②各行に対して、③cor関数を適用せよ。その際、④テンプレートyはdata.clとし、⑤相関係数の種類はparamで指定したものとする

```
in_f1 <- "sample16_log.txt" #入力ファイル名を指定し
in_f2 <- "sample16_cl.txt" #入力ファイル名を指定し
out_f <- "hoge1.txt" #出力ファイル名を指定し
param <- "pearson" #相関係数の種類を指定("pearson" または "spearman")
```

```
#入力ファイルの読み込みとラベル情報の作成
data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t", quote="") #in_f1で指定したファイルの読み込み
hoge <- read.table(in_f2, sep="\t", quote="") #in_f2で指定したファイルの読み込み
data.cl <- hoge[,2] #テンプレートパターンベクトルdata.clを作成
```

```
#本番
r <- apply(data, 1, cor, y=data.cl, method=param) #各(行)遺伝子についてテンプレートパターンdata.clとの相関
```

```
#ファイル保存
tmp <- cbind(rownames(data), data, r) #入力データの右側に相関係数rのベクトルを結合した結果をtmpに格納。
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
```



```
R Console
> #本番 ① ② ③ ④ ⑤
> r <- apply(data, 1, cor, y=data.cl, method=param) # $
> r
      gene1      gene2      gene3
0.9839815  0.8107810 -0.7068671
> |
```

Tips: cor, as.numeric

①apply関数を使わずに、遺伝子(つまり行)ごとにcor関数を用いて相関係数を計算する基本手順。②as.numeric関数は、data.clオブジェクトとデータの型を揃える目的で利用している。③as.numericなしの場合と比較すればよい

```
in_f1 <- "sample16_log.txt" #入力ファイル名を指定して
in_f2 <- "sample16_cl.txt" #入力ファイル名を指定して
out_f <- "hoge1.txt" #出力ファイル名を指定して
param <- "pearson" #相関係数の種類を指定("pearson" または "spearman")
```

```
#入力ファイルの読み込みとラベル情報の作成
data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t", quote="") #in_f1で指定したファイルの読み込み
hoge <- read.table(in_f2, sep="\t", quote="") #in_f2で指定したファイルの読み込み
data.cl <- hoge[,2] #テンプレートパターンベクトルdata.clを作成
```

```
#本番
r <- apply(data, 1, cor, y=data.cl, method=param) #各 (行) 遺伝子についてテンプレートパターンdata.clとの相関
```

```
#ファイルに保存
tmp <- cbind(r)
write.table(tmp, "out.txt", sep="\t", quote="")
```

R Console

```
> r
      gene1      gene2      gene3
0.9839815  0.8107810 -0.7068671
```

③

```
> data[1, ]
      G1_1 G1_2 G1_3 G1_4 G1_5 G1_6 G2_1 G2_2 G2_3 G2_4 G2_5
gene1 6.44  6.3 6.51 6.36 6.49 6.39 3.58 4.39 4.25  3.7 4.09
```

```
> data.cl
[1] 1 1 1 1 1 1 0 0 0 0 0
```

①

```
> cor(data[1, ], data.cl, method=param)
以下にエラー cor(data[1, ], data.cl, method = param) : 互換性のない次元です
```

②

```
> as.numeric(data[1, ])
[1] 6.44 6.30 6.51 6.36 6.49 6.39 3.58 4.39 4.25 3.70 4.09
```

```
> cor(as.numeric(data[1, ]), data.cl, method=param)
[1] 0.9839815
```

```
> |
```

①ではエラーが出ているが、②as.numericをつけたことで、無事gene1とテンプレートパターン間の相関係数rを計算できていることがわかる

Tips: cor, as.numeric

```

in_f1 <- "sample16_log.txt" #入力ファイル名を指定してin_f1に格納(発現データ)
in_f2 <- "sample16_cl.txt" #入力ファイル名を指定してin_f2に格納(テンプレート情報)
out_f <- "hoge1.txt" #出力ファイル名を指定してout_fに格納
param <- "pearson" #相関係数の種類を指定("pearson"または"spearman")

#入力ファイルの読み込みとラベル情報の作成
data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t", quote="")#in_f1で指定したファイルの読み込み
hoge <- read.table(in_f2, sep="\t", quote="")#in_f2で指定したファイルの読み込み
data.cl <- hoge[,2] #テンプレートパターンベクトルdata.clを作成

#本番
r <- apply(data, 1, cor, y=data.cl, method=param)#各(行)遺伝子についてテンプレートパターンdata.clとの相関

```

```

R Console
#ファイルに保存
tmp <- cbind(r)
write.table(tmp, "out.txt", sep="\t", quote="")

> r
      gene1      gene2      gene3
0.9839815 0.8107810 -0.7068671
> data[1, ]
      G1_1 G1_2 G1_3 G1_4 G1_5 G1_6 G2_1 G2_2 G2_3 G2_4 G2_5
gene1 6.44  6.3 6.51 6.36 6.49 6.39 3.58 4.39 4.25  3.7 4.09
> data.cl
[1] 1 1 1 1 1 1 0 0 0 0 0
① > cor(data[1, ], data.cl, method=param)
以下にエラー cor(data[1, ], data.cl, method = param) : 互換性のない次元です
② > as.numeric(data[1, ])
[1] 6.44 6.30 6.51 6.36 6.49 6.39 3.58 4.39 4.25 3.70 4.09
> cor(as.numeric(data[1, ]), data.cl, method=param)
[1] 0.9839815
> |

```

Tips: データの型

①「互換性のない次元です」と言われているのは、「相関係数を計算するときの入力データの見栄えが異なる」と文句を言われていると解釈すればよい。この「見栄え」は、Rで「データの型」と呼ばれるものに相当します。②と③が同じ見栄え(型)なのでcor関数の要求を満たしています

```
in_f1 <- "sample16_log.txt" #入力ファイル名を指定
in_f2 <- "sample16_cl.txt" #入力ファイル名を指定
out_f <- "hoge1.txt" #出力ファイル名を指定
param <- "pearson" #相関係数の種類を指定
```

```
#入力ファイルの読み込みとラベル情報の作成
data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t", quote="")#in_f1で指定したファイルの読み込み
hoge <- read.table(in_f2, sep="\t", quote="")#in_f2で指定したファイルの読み込み
data.cl <- hoge[,2] #テンプレートパターンベクトルdata.clを作成
```

```
#本番
r <- apply(data, 1, cor, y=data.cl, method=param)#各(行)遺伝子についてテンプレートパターンdata.clとの相関
```

R Console

```
#ファイルに保存
tmp <- cbind(r)
write.table(tmp, "r.txt")
```

```
> r
      gene1      gene2      gene3
0.9839815 0.8107810 -0.7068671
> data[1, ]
      G1_1 G1_2 G1_3 G1_4 G1_5 G1_6 G2_1 G2_2 G2_3 G2_4 G2_5
gene1 6.44  6.3  6.51 6.36 6.49 6.39 3.58 4.39 4.25  3.7  4.09
> data.cl
[1] 1 1 1 1 1 1 0 0 0 0 0
> cor(data[1, ], data.cl, method=param)
以下にエラー cor(data[1, ], data.cl, method = param) : 互換性のない次元です
> as.numeric(data[1, ])
[1] 6.44 6.30 6.51 6.36 6.49 6.39 3.58 4.39 4.25 3.70 4.09
> cor(as.numeric(data[1, ]), data.cl, method=param)
[1] 0.9839815
> |
```



Tips: as.vector, mode

慣れてくるとas.numericを使えばいいとすぐにわかるが、はじめのうちはas.character, as.integer, ①as.vectorなど既知のas…関数候補の中からそれっぽいものを試して型(見栄え)が揃うものを探すのが一般的かも…。②modeというデータの型を表示する関数もあるが…私は試行錯誤派です

```
in_f1 <- "sample16_log.txt" #入力ファイル名を指定してin_f1
in_f2 <- "sample16_cl.txt" #入力ファイル名を指定してin_f2
out_f <- "hoge1.txt" #出力ファイル名を指定してout_f
param <- "pearson" #相関係数の種類を指定("pearson")

#入力ファイルの読み込みとラベル情報の作成
data <- read.table(in_f1, header=TRUE, row.names=1, sep="\t", quote="")
hoge <- read.table(in_f2, sep="\t", quote="") #in_f2で指定したファイルの読み込み
data.cl <- hoge[,2] #テンプレートパターンベクトルdata.clを作成

#本番
r <- apply(data, 1, cor, y=data.cl, method=param) #各(行)遺伝子についてテンプレートパターンdata.clとの相関
```

```
R Console
#ファイルに保存
tmp <- cbind(r, data.cl)
write.table(tmp, "hoge1.txt", sep="\t", quote="")

> data.cl
[1] 1 1 1 1 1 1 0 0 0 0 0

> cor(data[1, ], data.cl, method=param)
以下にエラー cor(data[1, ], data.cl, method = param) : 互換性のない次元です

> as.numeric(data[1, ])
[1] 6.44 6.30 6.51 6.36 6.49 6.39 3.58 4.39 4.25 3.70 4.09

> cor(as.numeric(data[1, ]), data.cl, method=param)
[1] 0.9839815

① > as.vector(data[1, ])
      G1_1 G1_2 G1_3 G1_4 G1_5 G1_6 G2_1 G2_2 G2_3 G2_4 G2_5
gene1 6.44 6.3 6.51 6.36 6.49 6.39 3.58 4.39 4.25 3.7 4.09

② > mode(data[1, ])
[1] "list"
> |
```

コードの解説: cbind

出力部分。cbind関数を用いて出力させたい順番で列(column)方向で結合(bind)したものがtmpオブジェクト。尚、行(row)方向で結合させたい場合は、rbind関数を用いる

```
#本番  
r <- apply(data, 1, cor, y=data.cl, method=param)#各 (行) 遺伝子に
```

```
#ファイルに保存  
tmp <- cbind(rownames(data), data, r) #入力データの右側に相関係数rのベクトルを結合した結果をtmpに格納。  
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
```

```
R Console  
> #ファイルに保存  
> tmp <- cbind(rownames(data), data, r) #入力データの右側に相関係数rのベクトルを結合$  
> write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F)#tmpの中身を指定し$  
> rownames(data)  
[1] "gene1" "gene2" "gene3"  
> data  
      G1_1 G1_2 G1_3 G1_4 G1_5 G1_6 G2_1 G2_2 G2_3 G2_4 G2_5  
gene1 6.44 6.30 6.51 6.36 6.49 6.39 3.58 4.39 4.25 3.70 4.09  
gene2 5.81 6.93 6.73 5.55 6.39 6.61 2.81 5.46 1.00 3.46 4.17  
gene3 3.91 4.81 5.04 3.17 4.75 5.36 5.58 5.52 5.70 5.64 5.61  
> r  
      gene1      gene2      gene3  
0.9839815 0.8107810 -0.7068671  
> tmp  
      rownames(data) G1_1 G1_2 G1_3 G1_4 G1_5 G1_6 G2_1 G2_2 G2_3 G2_4 G2_5      r  
gene1      gene1 6.44 6.30 6.51 6.36 6.49 6.39 3.58 4.39 4.25 3.70 4.09 0.9839815  
gene2      gene2 5.81 6.93 6.73 5.55 6.39 6.61 2.81 5.46 1.00 3.46 4.17 0.8107810  
gene3      gene3 3.91 4.81 5.04 3.17 4.75 5.36 5.58 5.52 5.70 5.64 5.61 -0.7068671  
> |
```

