

日本医療検査科学会第52回大会 第2回医療情報技術セミナー 「知っておきたい臨床検査データ解析のテクニック」

2020-09-23\_ver.5 (final)



~Rを用いた解析の基本・応用・発展~

九州大学病院 検査部 瀬戸山大樹









### 本セミナーの内容

- セクション1:Rの基本操作(仮想臨床検査データ)
  - □ ファイルの読み込み、書き出し
  - □ 変数の扱い、挿入、条件抽出など
- セクション2:Rの応用操作(データの連結と縦横変換)
  - □ 列方向の連結
  - □ 横方向の連結
  - 縦持ちと横持ち変換
- セクション3:Rの応用操作(リストとループ処理)
  - □ リストの扱い
  - forループ処理について
- セクション4:病院データの解析前処理(リアル病院検査値データ)
  - □ データの統合、縦持ち・横持ち変換
  - □ 欠損値の少ない検査項目の選別

#### セクション5:臨床検査データの正規化(病院検査値+基準範囲データ)

セクション6:検診データの偏差値表示(検診+基準範囲データ)



	Download RStudio
RStudio	Discover Shiny
pen source and enterprise-ready	Discover RStudio Team
ofessional software for R	Discover RStudio Server Pro Standard and Enterprise

# 本セミナーで使用するRパッケージ

パッケージとはある機能に特化した関数の集合体のことです

現在16,356個のパッケージがCRANに登録されています(2020年9月23日)

- readxl: エクセルファイルを読み込むためのパッケージ(ver. 1.3.1)
- data.table: 大規模データを高速で読み書きするパッケージ(ver. 1.13.0)
- tidyr: データの縦持ち・横持ち形式を変換するパッケージ(ver. 1.1.2)
- dplyr: 数多くの前処理を簡単に行えるパッケージ(ver. 1.0.2)
- ggplot2: データを見栄え良く可視化するためのパッケージ(ver. 3.3.2)

• scales: グラフィックパラメータ調節のためのパッケージ(ver.1.1.1)

スライドでは、(パッケージ名)::(関数名)と表示します

例 readxlパッケージのread\_excel関数を使用する場合

readxl::read\_excel()

5

# Rパッケージのダウンロード



## セクション1:Rの基本

□ ファイルの読み込み、書き出し
 □ 変数の扱い、挿入、条件抽出など

7

# RStudioの立ち上げ



作業ディレクトリを指定しよう

データの読み込み先や保存場所を予め指定しておきます

#現在のディレクトリはgetwd()で把握できます
#macとwindowsでは少し標記が異なります
#デスクトップ上に「seminar\_2020」フォルダを置きます
setwd("/users/xxx/desktop/seminar\_2020/koushukai/")

#windowsの場合 setwd("C:/users/xxx/desktop/seminar\_2020/koushukai/")

## セクション1,2,3で扱うデータセット

#### ✓ 患者の検査値データ(test1.csv, test2.csv, test3.csv, test4.csv)

- 患者のイニシャル、性別sex、各検査値情報を含む仮想データです
- これらのデータを使ってRの基本操作を学びます

11

# csvファイルの読み込み

#csvファイルの読み込み #read.csv関数を使います(標準関数) #第一行目の変数名をheaderとし、第一列目をサンプル名に指定します test1 <- read.csv("test1.csv", header=TRUE, row.names=1)

Rstudioの右上のパネルの下記のように表示されます

			🔳 Project: (None) 👻	
Environment	History		_	
😅 🔒 🖃 📾	mport Dataset 👻 🎻		≣ List - 🥃	
🛑 Global Envir	onment <del>-</del>		(Q	
Data				
🔍 test1	10 obs. of	f 4 variables		
		🛑 Global Envi	ronment 👻	Q,
をクリッ	7	Data		
		🔾 test1	10 obs. of 4 varia	ables
		sex : Fa	ctor w/ 2 levels "F","M":	2 1 1 2 1 2 1 2 2 2
		Tcho : i	nt 220 230 240 240 250 260	260 260 270 280
	亜ボキティか	🛨 🛨 Tg : int	110 150 150 250 200 150 2	50 290 250 290
メッル	女川水小で1	severity	: int 0 1 2 1 3 3 2 1 4 4	
			severity: int	0 1 2 1 3 3 2 1 4 4

#csvファイルの読み込み #headerには変数名、第一行サンプル名を指定する #read.csv 関数を使います(標準関数) test1<-read.csv("test1.csv", header=TRUE, row.names=1)</pre>

#### 190918\_実習用スクリプト.R \* test1 \* 🔷 🔿 🗐 🖓 Filter Environment History 🕣 启 📄 Import Dataset 🗸 🔮 Tcho <sup>‡</sup> Tg <sup>‡</sup> severitŷ sex 🛑 Global Environment 🗸 DS M 220 110 0 Data TH F 230 150 1 🚺 test1 10 obs. of 4 variables 2 MS F 240 150 ここをクリック 1 240 250 MA M 3 TN F 250 200 KG M 260 150 3 YI F 250 2 260 左上のパネルに KI M 260 290 1 データが表示されます KD M 270 250 4 YA M 290 280 4

#### Rstudioの右上のパネルの下記のように表示されます

13

データの表示

#データを眺めてみま test1<-read.csv("te	ເມ st1.α	う csv",	heade	r=TR	UE, row.nan	nes=1)
test1						
変数名に下記の禁則文		-	•	変数		名前はユニーク(一意)で
"\$%^*+-()[] #!?<>		sex	Tcho	Тg	severity	なけれはいりません
	DS	M	220	110	0	
(	TH	F	230	150	1	
	MS	F	240	150	2	
	MA	М	240	250	1	
観測値	ΤN	F	250	200	3	
サンプル名	KG	М	260	150	3	
	ΥI	F	260	250	2	
	ΚI	М	260	290	1	
	KD	М	270	250	4	
Ĺ	YA	М	280	290	4	
	L	:				14

### データの構造を把握する

#str関数(標準関数)を使って構造をみます str(test1)

'data.frame': 10 obs. of 4 variables: \$ sex : Factor w/ 2 levels "F", "M": 2 1 1 2 1 2 1 2 2 2 \$ Tcho : int 220 230 240 240 250 260 260 260 270 280 \$ Tg : int 110 150 150 250 200 150 250 290 250 290 \$ severity: int 0 1 2 1 3 3 2 1 4 4

- test1オブジェクトはデータフレームです
- 10個の観測値(observation)と4つの変数(variables)を含みます
- 変数sexは因子型(factor)でFとMの2つの水準を含みます
- 変数Tcho, Tg, severityは整数型(integer)です

> データフレーム(data frame)は、複数の変数型を持つ2次元配列です
 > 行列(matrix)は、数値型のみの変数を持つ2次元配列です

備考 データの種類 型(class) 変数 スカラ scalar 数値のみ ベクター 単一の型 vector 行列 数値のみ matrix 2次元配列 データフレーム data.frame 2次元配列 複数種 複数(ベクター、 セクション3で 行列、データフ リスト list 説明します レーム) tidyverseパッ 複数(上記に加え、 tibble ケージで扱う リストも可) データ型の一種 #class(data)で確認できます class(test1)

> class(test1)

[1] "data.frame"

# データ型(class)について

15

# 変数の「型(mode)」について

'da	'data.frame': 10 obs. of 4 variables:													
\$	sex	:	Facto	or w/	/ 2 '	level	ls "I	=","N	4": 2	211	L 2 1	121	L 2 2	2 2
\$	Tcho	:	int	220	230	240	240	250	260	260	260	270	280	
\$	Tg	:	int	110	150	150	250	200	150	250	290	250	290	
\$	severity	:	int	0 1	21	33	21	44						

変数の種類	型(mode)	サブタイプ (type)	備考
カテゴリカル	factor		
数值型	numeric	integer	整数值
		double	実数値
文字型	character		
日付型	date		
論理型	boolean		TRUE/FALSE

17

# データの部分表示

#head関数(標準関数)を使って上位5行のみを表示させましょう #デフォルトは上位6行です

#### head(test1, 5)

> ł	<pre>&gt; head(test1,5)</pre>						
	sex	Tcho	Тg	severity			
DS	М	220	110	0			
TH	F	230	150	1			
MS	F	240	150	2			
MA	М	240	250	1			
ΤN	F	250	200	3			

#tail関数(標準関数)を使って下位5行のみを表示させましょう
#デフォルトは上位6行です
tail(test1, 5)

## 変数の表示法

#names関数(標準関数)を使って変数名を表示させます names(test1)

<pre>&gt; names(test1</pre>	.)		
[1] "sex"	"Tcho"	"Tg"	"severity"

19

## 変数のアクセス法

#中性脂肪Tgを取り出す方法 #\$を使用する方法 * <b>「\$」は「~の」の意味</b> だと思ってください					
test1\$1g <b>#[行,列]を使用する方法</b> #指定無しの場合は「すべてを取り出す」の意味です test1[,3] #直接名前を記入する方法 test1[ "Tg"]	"\$"や"["は演算機能を 有する記号=演算子 (オペレーター)と いいます				
<pre>&gt; test1\$Tg   [1] 110 150 150 250 200 150 250 290 250 29</pre>	0				
<pre>&gt; test1[,3] [1] 110 150 150 250 200 150 250 290 250 29</pre>	0				

[1] 110 150 150 250 200 150 250 290 250 290

> test1[,"Tg"]

### 括弧の使い分け

◆ 丸括弧 (parenthesis, "()")と角括弧(square blankets, "[]")の違い



## 複数の変数を取り出す方法 その1

#総コレステロールTcho、中性脂肪Tg、重症度severityの変数を取り出す
 #第2列から第4列までを連続して取り出す
 #「:」は連続値を取り出す演算子です
 test1[,2:4]

> t	est1	,2:4	1]
	Tcho	Тg	severity
DS	220	110	0
TH	230	150	1
MS	240	150	2
MA	240	250	1
ΤN	250	200	3
KG	260	150	3
ΥI	260	250	2
ΚI	260	290	1
KD	270	250	4
YA	280	290	4

## 複数の変数を取り出す方法 その2

#飛び飛びの変数はどうでしょうか?Tchoとseverityを抽出します
#「c」はcombine(結合する)関数です
test1[,c("Tcho", "severity")]
#列数指定でもできます
test1[,c(2,4)]

結果はどちらも一緒です

			_			
> t	est1[,	,c("Tcho","severity")]		> t	est1	[,c(2,4)]
	Tcho s	severity			Tcho	severity
DS	220	0		DS	220	0
TH	230	1		TH	230	1
MS	240	2		MS	240	2
MA	240	1		MA	240	1
ΤN	250	3		ΤN	250	3
KG	260	3		KG	260	3
ΥI	260	2		ΥI	260	2
κı	260	1		ΚI	260	1
KD	270	4		KD	270	4
YA	280	4		YA	280	4

23

# 条件にあうデータ(行)を抽出する方法

#severityが3の患者だけを抽 #真偽式を利用します test1\$severity==3	出する場合			
<pre>&gt; test1\$severity==3 [1] FALSE FALSE FALSE FALSE</pre>	E TRUE TRUE	FALSE FALSE	FALSE FALS	E
#真偽式がTRUEの行(患者) <b>test1[test1\$severity==3,</b> ]	を抽出します	F		
<pre>&gt; test1[test1\$severity==3,]     sex Tcho Tg severity TN F 250 200 3 KG M 260 150 3</pre>				
#dplyrパッケージの <b>filter()関</b> library(dplyr) #dplyrの実装	<b>数</b> を使うと簡	潔にコードで	きます	APH-

dplyr::filter(test1, severity==3)

Operator	Definition	意味
&	And	かつ
	Or	あるいは
==	Equal to	等しい
!	Not	~ではない
!=	Not equal to	等しくない
>	Bigger than	より大きい
>=	Equal to or bigger than	等しい、もしくは、 より大きい

## 条件抽出で使用するオペレーター(演算子)

# 新たな変数をデータフレームに挿入する方法

```
#新しい変数sumを作りましょう (sum=Tcho+Tg)
sum <- test1$Tcho+test1$Tg
sum
> sum
sum
[1] 330 380 390 490 450 410 510 550 520 570
#test1にsumを加えたデータフレームtest2を作りましょう
#test1をtest2にコピーします
test2<-test1
#test2にtest1がコピーされているか確認しましょう
Test2
#dplyrパッケージのmutate()関数を使うと簡潔にコードできます
test2<-dplyr::mutate(test1, sum=sum)
```

## 新たな変数をデータフレームに挿入する方法

#次に、test2に新しい変数sumを作って、sumを代入しましょう test2\$sum<-sum #test2を確認しましょう test2

>	test2	2<-te	st1		>	test	n ±⊂	新たた恋			
>	testa	2			>	test	利	にな多	〔 父 「		
	sex	Tcho	Тg	severity		sex	Tcho	Тg	severity	sum	
DS	М	220	110	0	DS	5 М	220	110	0	330	
TH	F	230	150	1	TH	l F	230	150	1	380	
MS	F	240	150	2	MS	5 F	240	150	2	390	
MA	М	240	250	1	MA	M	240	250	1	490	
ΤN	F	250	200	3	TN	I F	250	200	3	450	
KG	М	260	150	3	KĢ	i M	260	150	3	410	
ΥI	F	260	250	2	YI	F	260	250	2	510	
ΚI	М	260	290	1	ΚI	M	260	290	1	550	
KD	М	270	250	4	KD	) М	270	250	4	520	
YA	М	280	290	4	YA	М	280	290	4	570	

データフレームから特定の変数を除去する方法

#test2のseverity(4列目)を除きましょう
#いくつかのやり方があります
test2[,-4] #4列目を除去するという意味です
test2[,c(1,2,3,5)] #4列目以外を選択するという意味です

#いずれも列数がわからない場合や変数が多い場合は使いにくいですね
#先程のオペレーターを使ってみましょう
#否定演算子「!」を使うと下記のようになります
test2[,!names(test2)=="severity"]

#dplyrパッケージの**select()関数**を使うと簡潔にコードできます dplyr::select(test2, -severity)



# 条件に合ったラベルを挿入する(ifelse文)

#ある変数に条件をつけて、複数のカテゴリーに分類しましょう
#ifelse関数を使い、ifelse(条件式,真,偽)のように使用します
#「grade」という新しい変数を作り、Tgが160未満であれば0を、160以上で
あれば1をつけましょう
test2\$grade<-ifelse(test2\$Tg<160,0,1)</pre>

#より複雑な条件分岐 #Tgが160未満であれば0を、200未満であれば1を、それ以上であれば2をつ けましょう test2\$grade<-ifelse(test2\$Tg<160, 0,

ifelse(test2\$Tg<200, 1, 2))

#dplyrパッケージのmutate()関数を使うと簡潔にコードできます test2<-dplyr::mutate(test2, grade=ifelse(Tg<160, 0, ifelse(Tg<200, 1, 2)))



# データフレームをcsvファイルに保存する

#データフレームを作業ディレクトリに保存しよう #write.csv関数(標準関数)を使います write.csv(test2, file="test1\_modified.csv")

#### Rの基本操作

- □ csvファイルの読み込み、書き出しができる
- □ データの型(クラス)、変数の型(モード)を理解した
- □ データの構造をつかむことができる
- □ データフレームの特定の変数を抽出できる
- □ データフレームに新たな変数を挿入できる
- □ データフレームの特定の変数を除去できる
- □ 条件分岐により、観察値にラベルを付与することができる

dplyrパッケージの関数	機能
filter()	特定の観測値(行)を抽出する
select()	特定の変数(列)抽出する
mutate()	新たな変数(列)を追加する

31

## セクション2:Rの応用操作 データの連結と縦横変換

- □ 縦方向の連結
- □ 横方向の連結
- 縦持ち横持ち変換

# 縦方向(行)の連結: dplyr::bind\_rows()

#このセクション2では複数のデータを連結する操作について学びます #主にdplyrパッケージの関数を使用します #SQLに精通している方にとってはおなじみの操作ですね #test2.csvファイルを読み込みましょう #test1とほぼ同じフォーマットですが、別の患者のデータが含まれています test3<-read.csv("test2.csv", header=T, row.names=1)

> t	est3			> 1	test1			
	sex	Tcho	Тg		sex	Tcho	Тg	severity
KK	М	230	110	DS	М	220	110	0
KN	F	210	150	TH	F	230	150	1
ΚT	F	220	150	MS	F	240	150	2
MS	М	240	270	MA	М	240	250	1
KS	F	250	200	TN	F	250	200	3
ΤN	М	260	150	KG	М	260	150	3
AY	F	210	230	YI	F	260	250	2
JW	F	260	290	KI	М	260	290	1
NS	М	280	270	KD	М	270	250	4
ТК	М	300	280	YA	М	280	290	4

33

# 縦方向(行)の連結: dplyr::bind\_rows()

	#デー	ータ	を縦	É方 向	」に連	.結する場合はdplyr::bind_rows()関数を使います 🏾 🔬
	test	13<	<-dp	olvr::	bind	rows(test1, test3)
	test	13	•	5	-	- ` ' '
	1031	10				· · · · · · · · · · · · · · · · · · ·
e	> test1	3		<u>.</u>	·····	1
ļ	DC	sex	I Cho	Ig sev	/erity	tost1
ļ		M	220	110	0	
ļ		г с	230	150	1	
ļ	M33	Г	240	250	2	
ļ	TN 5	F	250	200	יד א	
	KG	м	260	150	3	
ļ	YI	F	260	250	2	
ļ	KI	м	260	290	1	
ļ	KD	м	270	250	4	
ļ	YA	М	280	290	4	
	КК	М	230	110	NA	test2
ļ	KN	F	210	150	NA	
ļ	КТ	F	220	150	NA	
ļ	MS14	М	240	270	NA	│ ` 連結するデータの変数が一部異なる場合 │
ļ	KS	F	250	200	NA	は新たな変数が追加され、値を持たな
ļ	TN16	М	260	150	NA	
ļ	AY	F	210	230	NA	い。の力はINA(入損胆)衣小されより
ļ	JW	F	260	290	NA	
ł	NS	М	280	270	NA	

ТК

NA

M 300 280

# 横方向(列)の連結

#次に横方向の連結について学びます
#まずtest3.csvファイルを読み込みましょう
test4<- read.csv("test3.csv", header=T, row.names=1)</pre>

#test1と比べてみましょう test1 test4

						,				
>	te	est1				> t	test4	1		
	S	sex	Tcho	Тg	severity		AST	ALT	GGT	test//-/t_test1の一部
D	S	М	220	110	0	DS	25	22	25	の患者の別の検査デー
Т	Ή	F	230	150	1	TH	26	21	25	タが含まれています
м	IS	F	240	150	2	MS	22	22	39	
м	IA	М	240	250	1	MA	27	23	30	
Т	'N	F	250	200	3	TN	33	40	60	
К	G	М	260	150	3	KG	29	22	32	
Y	Ί	F	260	250	2	YI	30	26	40	
K	I	М	260	290	1	KI	29	25	41	
K	D	М	270	250	4					
Y	Ά	М	280	290	4					

# 横方向(列)の連結を紐付けるプライマリーキー

#2つのデータを横方向に連結させるとき、お互いを紐付ける変数(プライマ リーキーまたは主キーという)を指定します #プライマリーキーは必ず"一意(ユニークである)"でなければいけません

df.x							df.y				
ID	age	sex	ast	alt			ID	ć	alb	HbA1c	glucose
Α							А				
В					>	<b>×</b>	В				
С							С				
D							D				
「ID」をプライマリーキーとした連結 df.z											
	ID	ag	e s	ex	ast	al	t	alb	HbA1	lc glucc	se
	A										
	B										
	C										
	D										

# 横方向(列)の連結: dplyr::join ()関数群

#join()関数には様々なバリエーションがあります #プライマリーキーはby=" "で特定の変数を指定して連結します

+



(1) inner\_join(a, b, by="x1")





(2) left\_join(a, b, by="x1")



	x1	x3		x1	x2	x3
	Α	5		А	1	5
+	В	6	=	В	2	6
	E	7		С	3	NA
	F	8		D	4	NA

(3) full\_join(a, b, by="x1")



37

# 実際の連結

#te #し #引 te: te: te: te:	est1 、 プ なめて st1< st14 st14	とtes ライ -rea -rea .inne	st4を うの: マリ - タ: d.c: d.c: d.c: er<-	き横に連続 ままでは ーキーと を読み込 sv("test sv("test	結さ イニ む え <b>1.cs</b> <b>3.cs</b>	せま こシー て使い sv", sv", r_jo	ミす ャル うこ hea hea	がサンプル名として読み込まれているた とができません ります (row.names=1)のオプション eder=T) eder=T) est1, test4, by="ID")
> te	st14	inner	-					
IC	) sex	Tcho	Тg	severity	AST	ALT	GGT	test1とtest4で共通する串考の
1 DS	5 М	220	110	0	25	22	25	みの検査データを連結しました
2 TH	ł F	230	150	1	26	21	25	
3 MS	5 F	240	150	2	22	22	39	
4 MA	M	240	250	1	27	23	30	
5 TN	I F	250	200	3	33	40	60	
6 KC	i M	260	150	3	29	22	32	
7 YI	F	260	250	2	30	26	40	
8 KI	M	260	290	1	29	25	41	

### 複合プライマリーキーについて

#プライマリーキーは1つの変数である必要はありません #複数の変数の組み合わせが「一意(ユニークである)」であれば、それを 「複合プライマリーキー」として使い、データを紐付ける場合が多いです

#test5にtest4.csvを読み込み、test1と比べてみましょう

test5<-read.csv("test4.csv", header=T)

>	+00+5	
-		

> test1

	ID	sex	AST	ALT	GGT
1	DS	М	25	22	25
2	ΤH	F	26	21	25
3	MS	F	22	22	39
4	MA	М	27	23	30
5	ΤN	F	33	40	60
6	MA	F	29	22	32
7	ΥI	F	30	26	40
8	ΚI	М	29	25	41
9	MS	М	24	22	65

	ID	sex	Tcho	Тg	severity
1	DS	М	220	110	0
2	TH	F	230	150	1
3	MS	F	240	150	2
4	MA	М	240	250	1
5	ΤN	F	250	200	3
6	KG	М	260	150	3
7	ΥI	F	260	250	2
8	ΚI	М	260	290	1
9	KD	М	270	250	4
10	YA	М	280	290	4

test5ではMSさんとMA さんが重複しています ので、IDをプライマ リーキーとして使うこ とはできません



39

# 複合プライマリーキーを使った連結

#プライマリーキーをIDに指定した場合とID+sexを指定した場合で比較してみましょう

test15<-dplyr::full\_join(test1, test5, by="ID") test15.unique<-dplyr::full\_join(test1, test5, by=c("ID","sex"))

>	tes	est15										> test15.unique							
	ID	sex.x	Tcho	Тg	severity	sex.y	AST	ALT	GGT			ID	sex	Tcho	Тg	severity	AST	ALT	GGT
1	DS	М	220	110	0	М	25	22	25	-	1	DS	М	220	110	0	25	22	25
2	TH	F	230	150	1	F	26	21	25	Ĩ	2	ΤH	F	230	150	1	26	21	25
3	MS	F	240	150	2	F	22	22	39	3	3	MS	F	240	150	2	22	22	39
4	MS	F	240	150	2	М	24	22	65		4	MA	М	240	250	1	27	23	30
5	MA	М	240	250	1	М	27	23	30		5	ΤN	F	250	200	3	33	40	60
6	MA	М	240	250	1	F	29	22	32	(	6	KG	М	260	150	3	NA	NA	NA
7	ΤN	F	250	200	3	F	33	40	60		7	ΥI	F	260	250	2	30	26	40
8	KG	М	260	150	3	<na></na>	NA	NA	NA	8	8	ΚI	М	260	290	1	29	25	41
9	ΥI	F	260	250	2	F	30	26	40	ç	9	KD	М	270	250	4	NA	NA	NA
10	) KI	М	260	290	1	М	29	25	41	-	10	YA	М	280	290	4	NA	NA	NA
11	l KD	М	270	250	4	<na></na>	NA	NA	NA	:	11	MA	F	NA	NA	NA	29	22	32
17	γ <b>Υ</b> Δ	м	280	290	4	<na></na>	NA	NA	NA		12	MS	М	NA	NA	NA	24	22	65

test15ではMSさんとMAさんの連結がおかしいです test15.uniqueでは正しく連結されています

# データの縦持ちと横持ちについて

#ここではデータの表示形式について学びます #test1のデータ表示方法には2通りあります



41

# tidyrのspread()関数とgather()関数

#縦持ちを横持ちに変換する時は**tidyr::spread()**関数 #横持ちを縦持ちに変換する時は**tidyr::gather()**関数を使います





### 実際の縦持ち変換

#test1のデータは横持ち形式です
#test1を縦持ちに変換しましょう
#test1のIDは変数名として読み込んでおきます
test1<-read.csv("test1.csv", header=T)
library(tidyr) # tidyrパッケージを実装します
test1.long<- tidyr::gather(test1, key="item", value="score",</pre>

Tcho, Tg, severity)

# str(test1.long) test1.long

> 9	str(te	est	:1.lo	ong)									
'da	ata.fr	an	ne':	30	obs	. of	4 ١	/ario	ables	5:			
\$	ID	:	chr	"DS"	' "TI	H" "N	4S" '	'MA''	•••				
\$	sex	:	chr	"M"	"F"	"F"	"M"	•••					
\$	item	:	chr	"Tcł	10"	"Tch	o" "T	[cho	' "То	cho"	•••		
\$	score	e:	int	220	230	240	240	250	260	260	260	270	280
> 1	test1.	.10	ong										
	ID s€	ex		item	sco	re							
1	DS	М		Tcho	27	20							
2	TH	F		Tcho	2	30							
3	MS	F		Tcho	24	40							
4	MA	М		Tcho	24	40							

43

## 実際の横持ち変換

#今度は逆にtest1.longを横持ち変換しましょう test1.wide<-tidyr::spread(test1.long, key="item", value="score") test1.wide

#### > test1.wide

	ID	sex	severity	Tcho	Тg
1	DS	М	0	220	110
2	KD	М	4	270	250
3	KG	М	3	260	150
4	ΚI	М	1	260	290
5	MA	М	1	240	250
6	MS	F	2	240	150
7	ТΗ	F	1	230	150
8	ΤN	F	3	250	200
9	YA	М	4	280	290
10	ΥI	F	2	260	250

\_\_\_\_ 元のtest1に戻りました

#### 解析の目的に応じて、横持ち形式と縦持ち形式を使い分けて、作 業の効率を上げましょう

44

Rの応用操作 データの連結と縦横変換

- □ dplyrパッケージのbind\_rows()関数を使って2つのデータを縦 方向に連結できる
- 2つのデータを横方向に連結するためには一意変数であるプラ イマリーキーで紐付けることを理解した
- □ dplvrパッケージには様々なjoin()関数があり、目的に応じて 使い分けることを理解した
- □ 複数の変数を組み合わた複合プライマリーキーを使って2つの データを連結できる
- tidyrパッケージのgather()関数とspread()関数を使ってデータ の縦持ち・横持ち変換ができる



## セクション3:Rの応用

□ リストの扱い ■ forループ処理について

## リストについて

#### リストは異なるデータを複数格納できるデータ型(class)のことです 行列(matrix)、データフレーム、ベクトルなど何でも入れ込めます

#データフレームtest1とtest2をリストに代入しましょう #list関数 (標準関数)を使います test.list<-list(test1,test2) test.list

>	<u>tes</u> t.	list				
	1]])	-				test.listリストの要素1にtest1データフ
	sex	Tcho	Тg	severity		レーム 要素2にtest2データフレーム
DS	М	220	110	0		が枚納されています
TH	F	230	150	1		
MS	F	240	150	2		
			and the second se	-		
	2]])					
	sex	Tcho	Тg	severity	sum	
DS	Μ	220	110	0	330	
TH	F	230	150	1	380	
MS	F	240	150	2	390	
MA	М	240	250	1	490	

リストの要素(の要素)を取り出す

#リストの要素は二重角括弧[[]]を使って取り出します #リストの要素の要素は[[]][]を使って取り出します #test2の3-5列目Tcho, severity, sumを取り出しましょう test.list[[2]][,3:5]

> 1	cest.	list[[2]][,3:5]
	Тg	severity sum
DS	110	0 330
тн	150	1 380
MS	150	2 390
MA	250	1 490
ΤN	200	3 450
KG	150	3 410
ΥI	250	2 510
κī	290	1 550
KD	250	4 520
YA	290	4 570

### 基本統計量を表示する

#summary関数(標準関数)を使ってデータの基本統計量を表示させます summary(test1)

sex	Tcho	Тд	severity
F:4	Min. :220	) Min. :110	Min. :0.0
M:6	1st Qu.:240	) 1st Qu.:150	1st Qu.:1.0
	Median :255	Median :225	Median :2.0
	Mean :251	. Mean :209	Mean :2.1
	3rd Qu.:260	) 3rd Qu.:250	3rd Qu.:3.0
	Max. :280	) Max. :290	Max. :4.0

表記	意味
Min.	最小値
1st Qu.	第一四分位
Median	中央値
Mean	平均值
3rd Qu.	第三四分位
Max.	最大値

49

# lapply関数を使ったリストの操作

#lapply関数(標準関数)を使って、リストの要素を一括処理しましょう #リストの要素の基本統計量をsummary関数を使って算出します

#### lapply(test.list, summary)

> lapp	ly(test.list,	summary)		
[[1]]				
sex	Tcho	Тд	severity	
F:4	Min. :220	Min. :110	Min. :0.0	
M:6	1st Qu.:240	1st Qu.:150	1st Qu.:1.0	
	Median :255	Median :225	Median :2.0	
	Mean :251	Mean :209	Mean :2.1	
	3rd Qu.:260	3rd Qu.:250	3rd Qu.:3.0	
	Max. :280	Max. :290	Max. :4.0	
[[2]]				
sex	Tcho	Тg	severity	sum
F:4	Min. :220	Min. :110	Min. :0.0	Min. :330.0
M:6	1st Qu.:240	1st Qu.:150	1st Qu.:1.0	1st Qu.:395.0
	Median :255	Median :225	Median :2.0	Median :470.0
	Mean :251	Mean :209	Mean :2.1	Mean :460.0
	3rd Qu.:260	3rd Qu.:250	3rd Qu.:3.0	3rd Qu.:517.5
	Max. :280	Max. :290	Max. :4.0	Max. :570.0

#次に、リストの要素の要素(Tchoとseverity)を取り出しましょう lapply(test.list, "[", c(2,4))

各要素から同じ要素が取り出されていますね この操作は、フォーマットが共通している要素に対して一括処理する際に絶大 な威力を発揮します(詳細はセクション4;スライド61) ここではforループについて学びます。

ループは、他のいろいろなプログラミング言語でも見られる概念で、繰り返 し(反復)処理のことです。しかし実際、Rではforループは計算速度が遅い などの理由から敬遠される傾向にあり、代わりに非常に高速に動く関数群が 用意されています。しかし、どうしてもループを使わざるを得ない場合があ り、その時には絶大な威力を発揮します。後のセクション5,6ではそのような 事例を取り上げます。

#まず、colSums関数(標準関数)を使って、test2の一部の要素を集計しま しょう #colSums関数はcol(列)に対してSums(合計)値を算出する関数です

#colsums)(新知道(19))に対してSums(古計) 値を算出する関数で result1<-colSums(test2[,2:4])

#### > result1<-colSums(test2[,2:4])</pre>

> result1

**result1** 第1要素 第2要素

第3要素

Tcho	Тg	severity
2510	2090	21

51

## colSums関数の処理を解読する

result1<-colSums(test2[,2	2:4])
> test2[,2:4] Tcho Is severity	test2の2列目から4列目までの各行の合計値 を求める処理
TH 230 150 1	1. 処理1:2列目の合計値を求め、第1要素に代入 する
MS 240 150 2 MA 07200 259 1	2. 処理2:3列目の合計値を求め、第2要素に代入 する
TN 250 <b>処理2</b> 3 KG 260 150 <b>処理3</b>	3. 処理3:4列目の合計値を求め、第3要素に代入 する
YI 260 250 2	4. 変数名を代入する
KT 260 290 1 KD 270 250 4	
YA 280 290 4 合計値算出	
2510 2090 21	

# apply関数の処理を解読する

#Rにはapply関数群と呼ばれる一括演算に特化した関数が用意されています #apply関数は**apply(matrix, a, function)**として使用します #「各行」に作用させる場合はa=1,「各列」に作用させる場合はa=2を代入 します #先程のcolSums関数はapply(matrix, 2, sum)と同義です result2<-apply(test2[,2:4], 2, sum)

# colSums()処理をforループで置き換える



### 「臨床検査データ解析」におけるforループ

これまで、私自身が臨床検査データの解析を行ってきた中で、forループを使わざるを得なかった4つの事例について紹介します

事例1: 検査データから、患者ごとに検査値の時系列グラフを作成する場合



<u>事例2: 検査データから、検査値ごとのヒストグラムを一括表示させる場合</u>

<u>事例3: 検査データを検査項目ごとの基準範囲に従って標準化する場合</u>

セクション4,5で説明します

事例4: 検査データの機械学習を自動化する場合

セクション3のまとめ

#### Rの応用

- □ 複数のデータフレームをリスト化できる
- □ リストの要素を取り出すことができる
- □ リストの各要素を一括処理できる
- □ forループで行われている処理を理解できる
- □ forループを使わない関数群を知っている



### セクション4:病院検査値データ解析の実際

### 講習会で扱うデータセット

#### ✓ 病院データベースから抽出したリアルデータ(エクセルファイル)

- 「koushukai」フォルダの「data」に入っています
- 2016年に九州大学病院検査部で測定した生化学検査値のデータです
- 月ごとにエクセルファイル(.xlsx)にまとめられています
- ・ 患者IDは変換済みです
- 主病名に「糖尿病」を含む患者のデータになります
- 169,158件の検査オーダーで3,990,925項目の検査値を含み ます
- ・ 患者21,729人分のデータになります
- 患者ID以外はほぼ原型をとどめています
- ・ 本講習会以外の用途で使用することを禁じます



生データのままでは解析できないので、 解析できるフォーマットに変換します 57

#リアル病院データを使いましょう!

#扱うデータは、九州大学病院の2016年の生化学検査データ

#まずは2016年1月のエクセルファイルを開いてみましょう

1	A	В	С	D	E	F	G	н	1	J	К	L	м	N C	р	Q	R	S	т	U V
患者	fiD /	患者年齢 患	者性別名称	文書日付 D	オーダ番号	入外名称	診療科コート	診療科名称項	目連番 検	査コード	検査名称	結果	結果(数値)	結果コ;上限	、下『単位	採取日時	検体採取日	分野コ-	分野名	標準材 標準材料(
200	903475	28 女		2016-01-01 06:00:00	27919350	入院	21	第一外科	22 000	00000100	総蛋白	6.3	6.300	L	g/dL	2016-01-01 06:00:0	20160101	00015	化学	000018 全血
200	903475	28 女		2016-01-01 06:00:00	27919350	入院	21	第一外科	23 000	00000200	アルブミン	4.6	4,600		g/dL	2016-01-01 06:00:0	20160101	00015	化学	000018 全血
20'	903475	28 🕏		2016-01-01 06:00:00	27919350	入院	21	第一外科	24 000	00000300	尿素窒素	12	12.000		mg/dL	2016-01-01 06:00:0	20160101	00015	化学	000018 全血
20'	903475	28 😾		2016-01-01 06:00:00	27919350	入院	21	第一外科	25 000	00000400	クレアチニン	0.62	0.620		mg/dL	2016-01-01 06:00:0	20160101	00015	化学	000018 全血
20'	903475	28 🖈		2016-01-01 06:00:00	27919350	入院	21	第一外科	26 000	00000500	尿酸	2.9	2 900		mg/dl	2016-01-01 06:00:0	20160101	00015	化学	000018 全血
20'	903475	28 #	-	2016-01-01 06:00:00	27919350	入院	21	第一外科	27 000	00000800	総ビリルビン	0.9	0.900		mg/dl	2016-01-01 06:00:0	20160101	00015	化学	000018 全血
20'	903475	28 7	-	2016-01-01 06:00:00	27919350	入院	21	第一外科	28 000	00001200	AST	15	15,000		11/1	2016-01-01 06:00:0	20160101	00015	化学	000018 全血
20	903475	28 #	-	2016-01-01 06:00:00	27919350	入院	21	第一外科	29 000	00001300	ALT	10	10.000		U/L	2016-01-01 06:00:0	20160101	00015	化学	000018 全血
20	903475	28 #	-	2016-01-01 06:00:00	27919350	入院	21	第一从利	30 000	00003100	I DH	125	125,000		11/1	2016-01-01 06:00:0	20160101	00015	化学	000018 全血
20	002475	20 5	-	2016-01-01 06:00:00	27010250	入院	21	第二月到	21 000	00003100	ALB	04	84 000	1	U/L	2016-01-01 06:00:0	20160101	00015	化学	000018 全曲
200	003473	20 5	•	2010 01 01 00.00.00	27818550	NOC	21	345 2111-1	31 000	00001300	ALF	04	34.000		0/2	2010 01 01 00.00.0	20100101	00013	10-	000018 ±
200	۰ ۱-	キアウ	<u> </u>	カベ	71-	-1+	1/二/	-	± 1 t	全本	「古日の	· <u>ــــــــــــــــــــــــــــــــــــ</u>		エシーや	一一个	ちゃわ	71	+ -	+	
200	11	内り兀	テー	- ベ ヘー	- 人に	-12	11T (	こうさ	TΤS	史百、	坦日の	ヮ゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚゚	— 'X '	Πク エ\.	- C l 木	14 A 11	、ししい	エ	9	
20:	00 /1	1 3 1 2 0							1.		<u> </u>	-		/// - 0	- 61			0.	-	
20:	- 00	2 10	+ 44	¥ E のニ	* <i>F</i>	- TT	<u>++</u> /+	- 7 -	1+1	L		1-1		+ +						
201	00 (	- X L	~ 前	トラリナ	ーーク	くガク	T1.10	かんし	141	riav	data)	61	いしい	ます	~					
200	00 -		<b>—</b> •••••		-	112	- 0 ( -			<b>,</b>					5					
200	- 00	_~	4	*	- АП ТБ		7 N+	1-1-	7 1	<u>ν π</u>	キーナ	エミギ	Π.	いナ	1 · •					
20:	- 00 -	r -	メヘ	ヽー ス ()	⅀ℳリア	± q	る時	· 1. 12	. 0	) #//-	て(.()) 万	71 ł	ルフレン・	TJ d	しいひと	(`` d				
20:	30 Z		-		· / _ · _	- /	0,000	1-10	<u> </u>	- /// -	-0/5	,2 )	$\sim$	· ·	• • •	~ /				
200	002475	29 +	-	2016-01-01 06:00:00	27010250	3.82	21	第二月刊	42 000	00006209	stêr mîn	(-)		v		2016-01-01 06:00:0	r <sup>7</sup> 20160101	00015	を使	000018 全曲
200	120012	20 5		2016-01-01 06:12:11	20022544	入院	52	か 法 順 倞	24 000	00000203	始星白	44	4 400	î	a / dl	2016-01-01 06:12:1	120160101	00015	ル学	000018 全血
21	100010	30 5		2010-01-01-00.13.11	20022044	八院	50	血液腫病	24 000	00000100	10 風口	0.1	4.400		g/ dL	2010-01-01 00.13.1	20100101	00015	しナ	000010 ± m
21	120012	30 5		2016-01-01 06.13.11	20022044	入院	52	血液腫瘍	20 000	00000200	「ルノミノ	45	45.000		g/uL	2016-01-01 06:13:1	120160101	00015	ル学	000018 主血
21	130813	38 54		2010-01-01 00:13:11	28022544	入院	52	血液腫瘍	26 000	00000300	原 茶 至 茶	40	45.000	H H	mg/ dL	2016-01-01 06:13:1	120160101	00015	16子	000018 主血
3 21	130813	38 54	-	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	27 000	00000400	クレアナーン	1.43	1.430	н	mg/ dL	2016-01-01 06:13:1	120160101	00015	16子	000018 主血
21	130813	38 54	-	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	28 000	00000500	水販	3.5	3.500		mg/dL	2016-01-01 06:13:1	120160101	00015	16字	000018 主血
3 211	130813	38 32		2016-01-01 06:13:11	28022544	人院	52	血液腫瘍	29 000	00000800	総ビリルビン	2.3	2.300	н	mg/dL	2016-01-01 06:13:1	120160101	00015	化学	000018 至血
3 211	130813	38 女		2016-01-01 06:13:11	28022544	人院	52	血液腫瘍	30 000	00006600	直接ヒリルヒン	1.1	1.100	H	mg/dL	2016-01-01 06:13:1	120160101	00015	化字	000018 全面
211	130813	38 女		2016-01-01 06:13:11	28022544	人院	52	血液腫瘍	31 000	00001200	AST	33	33.000	н	U/L	2016-01-01 06:13:1	120160101	00015	化学	000018 全血
211	130813	38 女		2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	32 000	00001300	ALT	26	26.000	н	U/L	2016-01-01 06:13:1	120160101	00015	化学	000018 全血
2 211	130813	38 女	;	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	33 000	00003100	LDH	315	315.000	н	U/L	2016-01-01 06:13:1	120160101	00015	化学	000018 全血
3 2	21130813	38 女		2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	34 000	0001500	ALP	762	762.000	н	U/L	2016-01-01 06:13:11	20160101	00015	化学	000018 全血
* 2	21139813	38 50		ZM10-M1-M1 M01 31	28022566		1.00.0									I AND A REPORT OF A REAL AND A				

立ち上がり遅い	同じ患者IDが多い
日付が細かい	診療科名いろいろ
結果もいろいろ	単位もいろいろ。。

59

### リアルな病院検査データ

A	В	С	D	E	F	G	н	I	J	к	L	м
患者ID	患者年齢	患者性別名称	文書日付_D	オーダ番号	入外名称	診療科コー	ト診療科名称	項目連番	検査コード	検査名称	結果	結果(数値) 結
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	22	0000000100	総蛋白	6.3	6.300
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	23	0000000200	アルブミン	4.6	4.600
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	24	000000300	尿素窒素	12	12.000
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	25	0000000400	クレアチニン	0.62	0.620
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	26	0000000500	尿酸	2.9	2.900
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	27	0000000800	総ビリルビン	0.9	0.900
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	28	0000001200	AST	15	15.000
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	29	0000001300	ALT	10	10.000
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	30	0000003100	LDH	125	125.000
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	31	0000001500	ALP	84	84.000
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	32	0000002700	СК	120	120.000
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	33	0000002000	グルコース	95	95.000
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	34	0000002800	C反応性蛋白	0.05	0.050
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	35	0000002100	ナトリウム	141	141.000
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	36	0000002200	カリウム	3.5	3.500
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	37	0000002300	クロール	107	107.000
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	38	000000600	カルシウム	10.1	10.100
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	39	0000006400	推算糸球体濾過量	93	93.000
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	40	0000006207	A/G比	2.71	2.710
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	41	0000006208	乳濁	(-)	
20903475	28	女	2016-01-01 06:00:00	27919350	入院	21	第一外科	42	0000006209	溶血	(-)	
21130813	38	女	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	24	0000000100	総蛋白	4.4	4.400
21130813	38	女	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	25	000000200	アルブミン	2.1	2.100
21130813	38	女	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	26	000000300	尿素窒素	45	45.000
21130813	38	女	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	27	0000000400	クレアチニン	1.43	1.430
21130813	38	女	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	28	0000000500	尿酸	3.5	3.500
21130813	38	女	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	29	0000000800	総ビリルビン	2.3	2.300
21130813	38	女	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	30	0000006600	直接ビリルビン	1.1	1.100
21130813	38	女	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	31	0000001200	AST	33	33.000
21130813	38	女	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	32	0000001300	ALT	26	26.000
21130813	38	女	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	33	0000003100	LDH	315	315.000
21130813	38	女	2016-01-01 06:13:11	28022544	入院	52	血液腫瘍	34	000001500	ALP	762	762.000
z1130813	38	¥.	2016-01-01 06:13:11	z8022544	入院	52	血液腫瘍	35	000021100	AMY	63	63.000
						•	患者	ťΠ	(第1	列目)		

本実習では右の項 目を使って解析を 進めます

- ・ 性別(第3列目)
- 日付(第4列目)
- オーダー番号 (第5列目)
- 検査コード(第10列目)
- 結果(数値)(第12列目)

### 病院データの解析前処理



パッケージの実装、関数のヘルプ参照

#必要なパッケージを実装しましょう
#エクセルファイル読み込み用
library(readxl)
#データ前処理用
library(data.table);library(tidyr);library(dplyr)
#可視化ツール
library(ggplot2)

#関数の使い方はヘルプを参照して下さい #?関数名で右下の画面に表示されます #各関数がもつ引数やパラメーター設定など確認して下さい

### ファイルパスの一括取得

#毎月の検査値エクセルデータ12ヶ月分がdataフォルダ内にあります #今からdataフォルダのエクセルファイルを一括で読み込みます #まず、現作業フォルダの下にdataフォルダを指定します path <- "./data/" #次にdataフォルダのファイル名を全て入手します #list.files関数(標準関数)を使います fl<-list.files(path, full.names=T) #確認しましょう fl

> fl

[1] "./data//201601\_chemistry.xlsx" "./data//201602\_chemistry.xlsx"
[3] "./data//201603\_chemistry.xlsx" "./data//201604\_chemistry.xlsx"
[5] "./data//201605\_chemistry.xlsx" "./data//201606\_chemistry.xlsx"
[7] "./data//201607\_chemistry.xlsx" "./data//201608\_chemistry.xlsx"
[9] "./data//201609\_chemistry.xlsx" "./data//201610\_chemistry.xlsx"

# エクセルファイルを一括読み込み

#apply関数のリスト版であるlapply関数を使います #readxl::read\_excel関数と組み合わせることで各月の検査値データをリストの 要素として読み込ませます data <- lapply(fl, readxl::read\_excel)。 。 。 。 #RStudio右上のオブジェクトパネルを確認してください

12ヶ月分の個別データ





## 項目の絞り込み



# リストの各要素を縦方向に連結する処理

#リストに含まれるデータファイルを縦方向に連結します #全ての要素のフォーマットが全く同じであることが必要です #dplyrパッケージのbind\_rows関数を使います

data.bind<-dplyr::bind\_rows(data.select)







65

### 日本語標記名の確認と変換

#### #名前を見ましょう names(data.bind)

> names(data.bind)

[1] "患者 I D" "患者性別名称" "文書日付\_D" "オーダ番号" "検査コード" "結果"

#ちゃんと抽出できていますが、日本語表記は扱いづらいので、英語表記に 変換します

names(data.bind)<-c("id","sexmale", "day","order","item","score") #確認して下さい

# 変数の型(mode)を確認する

#このデータフレームの各変数のデータ型を確認しましょう

#### str(data.bind)

#### > str(data.bind)

tibble [3,990,925  $\times$  6] (S3: tbl\_df/tbl/data.frame)

- \$ id : num [1:3990925] 20903475 20903475 20903475 20903475 ...
- \$ sexmale: chr [1:3990925] "女" "女" "女" "女" ...
- \$ day : chr [1:3990925] "2016-01-01 06:00:00" "2016-01-01 06:00:00" "2016-01-01 06:00:00" "2016-01-01 06:00:00" ...
- \$ order : num [1:3990925] 27919350 27919350 27919350 27919350 ...
- \$ item : chr [1:3990925] "0000000100" "0000000200" "0000000300" "0000000400" ...
- \$ score : chr [1:3990925] "6.3" "4.6" "12" "0.62" ...



### 変数の型を変換する

#idと検査オーダーを文字化(character)します
#オーダー日を日付型化(Date)します(\*年月日に短縮されます)
#検査コードと検査値は数値化(numeric)します
data.bind\$id<-as.character(data.bind\$id)
data.bind\$sexmale<-factor(data.bind\$sexmale,levels=c("女","
男"),labels=c(0,1))
data.bind\$order<-as.character(data.bind\$order)
data.bind\$item<-as.numeric(data.bind\$item)
data.bind\$item<-as.Date(data.bind\$item)
data.bind\$score<-as.numeric(data.bind\$score)
#確認しましょう
str(data.bind)</pre>

69

## 縦持ち->横持ち変換

#縦持ち(long型)から横持ち(wide型) に変換します #検査オーダー(ID,日付) ごとのデータに変形します spread.data<-tidyr::spread(data.bind, key=item, value=score) spread.data<-as.data.frame(spread.data)



#右上のパネルでspread.dataオブジェクトを確認してください

### リアル検査データの観察

#spread.dataの中身を観察してみましょう str(spread.data)

> str(spr	ead.do	ata)	
'data.fra	me':	169158 obs. of 153 variables:	― 169,158の観測値(検
\$ id	: chi	r "20903475" "21130813" "13279003" "8906411"	査オーター)があり、
\$ sexmal	e: Fa	ctor w/ 2 levels "0","1": 1 1 2 2 1 2 2 1 1 1	153種類の検査項目
\$ day	: Dat	te, format: "2016-01-01" "2016-01-01" "2016-01-01"	… し かあります
\$ order	: chi	r "27919350" "28022544" "28019770" "28011398"	
\$ 100	: nur	m 6.3 4.4 4.6 4.5 6.2 8.7 NA 4.2 3.3 4	
\$ 200	: nur	m 4.6 2.1 2.3 2.7 3.4 1.9 NA 2.9 2.6 1.2	
\$2900 \$3100	: nur : nur	m NA	「NA」は欠損値です ) 検索頂日によってけば
\$ 3300	: nur	m NA NA NA 38 NA NA NA NA NA NA	とんどが欠損値です
\$ 3900	: nur	m NA	
\$ 3901	: nur	m NA NA NA NA NA NA NA NA NA	
\$ 3902	: nur	m NA NA NA NA NA NA NA NA NA	

検査項目ごとに、どれだけ欠損値を含むのかを調べる必要があります

71

# 検査値の欠損率の算出

#各検査値の欠損値を把握しましょう #各検査値の欠損値を算出する関数を作りました #下記スクリプトをコピペして実行して下さい na.summary<-function(data){ name<-names(data) #nameは各検査値の名前(検査コード) Value<-colSums(!is.na(data)) #Nは値を含むデータ数 Missing<-colSums(is.na(data)) #Missingは欠損値を含むデータ数 NA.ratio<-round((Missing/nrow(data))\*100,2) #NA.ratioは欠損率 df<-as.data.frame(cbind(name,Value,Missing,NA.ratio)) write.csv(df,file="summary of NA.csv") } #使い方はna.summary(データフレーム)です

#作業フォルダに結果が保存されます

na.summary(spread.data)

### 病院リアルデータの特徴把握

#summary of NA.csvファイルをエクセルで開きましょう

			値を含む	む数 欠損値	直を含む数	欠損率		
検	査値コード	В		D	E		-	
1		name	Value	Missing	NA.ratio			
2	id	b,	169158	0	(	) 7	検査値コート	、対応表
3	day	da	169158	0	(	)	(エクセル) も	ファイル)
4	order	order	169158	0	(	)	を参照して	1211
5	100	100	143359	25799	15.25	5 00	00000100	総蛋白
6	200	200	146669	22489	13.29	9 00	00000200	アルブミン
7	300	300	157149	12009	7.3	1 00	00000300	尿素窒素
8	400	400	157817	11341	6.7	7 0	00000400	クレアチニン
9	500	500	118947	50211	29.68	3 0	00000500	尿酸
.0	600	600	113614	55544	32.84	4 00	000000600	カルシウム
.1	700	700	58519	110639	65.43	1 00	00000700	無機リン
.2	800	800	150520	18638	11.02	2 0	00800000	総ビリルビン
.3	1000	1000	2638	166520	98.44	4 00	00001000	ТТТ

検査項目ごとに欠損値が異なる

73

# 検査項目の絞り込み

#欠損率を多く含む検査値を除外します

#充足率が70%以上の検査値だけをspread.dataから抽出します

select.70<-

spread.data[ ,colSums(!is.na(spread.data))>nrow(spread.data)\*0.7]

真偽式になっています

翻訳すると。。

spread.dataの中から、欠損値ではない(!is.na(spread.data))数を各検査項目 ごとにカウントして(colSums)、それがデータ数70%(nrow(spread.data)\*0.7) より大きい(>)場合を判定し、TRUEの検査項目(列)を抽出します

names(select.70)

> names(select.70) [1] "id" "day" "order" "100" "200" "300" "400" "500" "800" "1200" "1300" [12] "1500" "1600" "2000" "2100" "2800" "3100" "2200" "2300" "6207" "6400" これらの21種類の変数(検査項目)が該当するそうです

### 検査値名への変換

#正式な検査値名(略称)に変換しましょう #name70.csvファイルに21種類の名前が入っています #このファイルを読み込んで、select.70の名前を置換しましょう name70<-read.csv("name70.csv",header=T) names(select.70)<-name70\$name names(select.70)

ようやくおなじみの名前が出てきました。。

#### > names(select.70)

] [1	[1] [2]	"id" "ALP"	"day" "GGT"	"ord "GLU	er" "	"TP" "Na"	"ALB" "K"	"BUN" "Cl"	"CRE" "CRP"	"UA" "LDH"	"Tbil" "A/G"	"AST" "eGFR"	"ALT"
	1	TP	総蛋白		10	GGT	gamma-GT						
	2	ALB	アルブミン		11	GLU	グルコース						
	3	BUN	尿素窒素		12	Na	Na+						
	4	CRE	クレアチニン	/	13	K	K+						
	5	UA	尿酸		14	CI	CI-						
	6	Tbil	総ビリルビン	/	15	CRP	CRP						
	7	AST	AST		16	LDH	LDH						
	8	ALT	ALT		17	A/G	AG比						
	9	ALP	ALP		18	eGFR	推算糸球体》	慮過量					75
													(')

データの確認

#データの構造を確認しておきましょう str(select.70)

#### > str(select.70)

'da	ata.fran	ne	': :	169158 obs. of 22 variables:
\$	id	:	chr	"20903475" "21130813" "13279003" "8906411"
\$	sexmale	e:	Fact	or w/ 2 levels "0","1": 1 1 2 2 1 2 2 1 1 1
\$	day	:	Date	, format: "2016-01-01" "2016-01-01" "2016-01-01"
\$	order	:	chr	"27919350" "28022544" "28019770" "28011398"
\$	ТР	:	num	6.3 4.4 4.6 4.5 6.2 8.7 NA 4.2 3.3 4
\$	ALB	:	num	4.6 2.1 2.3 2.7 3.4 1.9 NA 2.9 2.6 1.2
\$	BUN	:	num	12 45 35 7 24 22 NA 27 35 7
\$	CRE	:	num	0.62 1.43 2.21 0.58 0.59 3.09 NA 0.45 0.7 0.54 .
\$	UA	:	num	2.9 3.5 5 2.3 NA 7.1 NA 4.7 3.4 2.4
\$	Tbil	:	num	0.9 2.3 1.1 1.8 0.7 2.1 NA 3.5 1.5 0.8
\$	AST	:	num	15 33 52 148 17 61 NA 28 46 32
\$	ALT	:	num	10 26 26 561 25 24 NA 36 75 12
\$	ALP	:	num	84 762 242 189 207 375 NA 259 65 577
\$	GGT	:	num	NA NA NA 232 NA NA 233 NA NA NA
\$	GLU	:	num	95 118 195 120 NA 119 NA 143 125 105
\$	Na	:	num	141 139 135 143 134 142 NA 147 135 138
\$	К	:	num	3.5 3.2 3.5 3.9 4.2 3.4 NA 4.9 4.5 4.4
\$	C1	:	num	107 105 102 106 101 93 NA 118 108 105
\$	CRP	:	num	0.05 1.86 21.98 0.1 0.03
\$	LDH	:	num	125 315 383 130 133 304 NA 218 207 227
\$	AG	:	num	2.71 0.91 1 1.5 1.21 0.28 NA 2.23 3.71 0.43
\$	eGFR	:	num	93 34 23 115 75 18 NA 103 66 86

- 169,158の観測値(18種類の検 査値プロファイル)に絞り込む ことができました。
- このデータセットを使って、検 査値の相関解析やプロファイル 比較、あるいは時系列変化等を 解析することができます。
- また、患者の層別化解析なども 可能です。
- 患者の病歴や投薬情報等を付与 することができれば、より病態 に即した検査値変化の解析が可 能になります。

### 患者数の把握

#患者数は何人ですか? #unique関数(標準関数)を使って重複のないidを求めます #更に、length関数(標準関数)を使って数をカウントします #つまり、「重複無しのIDの数=患者数」です length(unique(select.70\$id))

#### > length(unique(select.70\$id))

[1] 21789

### セクション4のまとめ

#### 病院検査値データ解析の実際

- □ 病院から取り出した検査値データの構造を理解できる
- □ 共通したフォーマットファイルをリストとして読み込ませる ことができる
- □ リスト要素の連結ができる
- □ データフレームの変数の型を変換できる
- □ 縦持ち、横持ち変換ができる
- □ 変数ごとの欠損率を把握することができる
- □ 欠損率(またはデータの充足率)に従って変数を絞り込むことができる

### セクション5:病院検査値データの正規化

### このセクションの目的と免責事項

- これ以降のセクションでは、「検査値の基準範囲(reference interval」を 使って検査値データを正規化する方法を解説します。
- ただし、本手法は、私が独自に考案したもので、<u>本学会、本委員会および</u> 臨床検査学・医療情報学の分野で公認されているものではありません。
- 本手法に関するご質問・ご意見はいつでも歓迎しますが、本手法を用いて 生じた結果については、私は責任を負いかねますのでご了承ください。
- あくまでも、「<u>リファレンス情報を使って変数(検査値)ごとに値を正規</u> <u>化する技術を習得する</u>」という位置づけで臨んでください。



79

## 各検査値を可視化する

#各検査値の分布をヒストグラムを描いて眺めましょう #hist()関数を使い、hist(data)で可視化します #まず、変数名をnameオブジェクトに取り出しましょう #select.70の第5-22列目に検査値が含まれています name<-names(select.70[,5:22]) name > name [1] "TP" "ALB" "BUN" "CRE" "UA" "Tbil" "AST" "ALT" "ALP" "GGT" "GLU" "Na" "К" "Cl" [15] "CRP" "LDH" "AG" "eGFR" #尿酸値(UA)のヒストグラムは下記のように描けます hist(select.70[,name[5]], main=name[5]) #mainはタイトルを表示させるオプションです UA 50000 尿酸値は正規 Frequency 30000 分布をしてい  $\bigcirc$ ますね 10000

81

# 各検査値を可視化する

#forループ関数を使ってヒストグラムを一括表示させましょう #検査値は18種類ありますので、18枚の図を同じ画面に表示させます #par()関数のmfrowオプションで指定します

**par(mfrow=c(4,5))** #図を4行5列に配置する設定です

10

15 select.70[, name[5]]





検査値の対数変換値を可視化する

#検査値を対数変換(eを底とする) してヒストグラムを描いてみましょう
#log()関数を使って変換します
#ALTを例にとり、元の検査値と対数値のヒストグラムを比較しましょう
par(mfrow=c(1,2))
hist(select.70[,name[8]], main=name[8]) #元の検査値のグラフ
hist(log(select.70[,name[8]]), main=name[8]) #対数変換値のグラフ





# 臨床検査値の基準範囲の算出法



# 検査データの正規化の意味と算出法

リアル検査値 正規化検査値 値の分布や単位が異なる 左の制約が無くなる • 基準範囲は年齢や性別に 全ての検査値を同じ尺度 よって異なる場合がある (Zスコア)で表示させる 対数変換は全ての検査値を、厳密な意味において、正規分布化させるわけではない点に ご留意ください。あくまで、正規分布に「近似」させて処理しているに過ぎません <u>対数正</u>規化 x  $x^p$  $LL^T$ ULLL $UL^T$ Me m  $=(LL^{T})^{1/p} =m^{1/p}$  $=(UL^T)^{1/p}$ =m-1.96s=m+1.96s• 正規化検査值 = (ln(x)-m)/s • 正規化平均值m=(ln(下限值) + ln(上限值))/2 • 正規化標準偏差s = (ln(上限值) - ln(下限值))/3.92 \*In: eを底とするlogのこと

正規化データ(Zスコア)の解釈 Zスコア化された検査値 The Normal Distribution <sup>2</sup>robability 基準範囲 alue Values -2.5 99% of values Probability of Cases in portions of the curve ≈ 0.02 4 = 0.1359 ≈ 0.1359 ≈ 0.0013 ≈ 0.3413 ≈ 0.3413 .0214 0.0013 Standard Deviations From The Mean -4σ -2σ +1σ +2σ +3σ +4σ -3σ -1σ Cumulative % 0.1% 2.3% 15.9% 50% 84.1% 97.7% 99.9% Z Scores -4.0 . -3.0 . -2.0 -1.0 0 +1.0 +2.0 +3.0 +4.0 Xスコアの値 解釈 Z <- 2 異常低値(極端値) です -2 <= Z <= +2 基準範囲内です Z > +2 異常高値(極端値)です

87

#### ✓ 検査値の正規化を行うロジック

 1. 検査値の基準範囲(reference interval)は、基準個体の基準値 95%区間の上限・下限値で定義される
 2. 一部の検査値の基準範囲は性別や年齢によって異なる
 3. 検査値は固有の分布を持つが、対数変換(またはべき乗変換) すると正規分布(対数正規分布)に近似できる(経験則)
 4. 正規分布では平均値±2SD値の範囲に全体の約95%が含まれる
 対数変換した基準範囲の上・下限値を使った対数正規分布に対し、 対数変換した検査値を当てはめることで、基準範囲に基づく正規 化された値(Zスコア)を算出できる
 ・ 検査値同士をZスコアで比較することができる
 ・ 検査値を性別や基準範囲に照合させることなく、正常高・低値 または異常値の判断が可能となる
 ・ 正規化データを使った様々な解析(機械学習など)が可能となる

単位 検査項目 男性 女性 略称 Serum chemistry Serum creatinine CRE mg/dL 0.65~1.07 0.46~0.79 Uric acid UA mg/dL 3.7~7.8  $2.6 \sim 5.5$ Triglyceride ΤG mg/dL 40~234 30~117 High-density lipoprotein mg/dL 38~90 48~103 HDL Alanine aminotransferase U/L ALT 10~42 7~23 gamma-Glutamvl GGT U/L 13~64 9~32 transferase Creatine kinase СК U/L 59~248 41~153 Complete blood count (CBC) Red blood cell  $\times 10^{6} / \mu L$ RBC 4.35~5.55 3.86~4.92 Hemoglobin g/dL 13.7~16.8 11.6~14.8 Hb Ηt % Hematocrit 40.7~50.1 35.1~44.4

JCCLSの共用基準範囲による

## 性別で基準範囲が異なる検査項目

89

#### 男女別の検査値の共用基準範囲データ

- 「reference interval.m.csv」
- [reference interval.f.csv]
- 全43項目の検査値名(description)、略称(item)、単位(unit)、基準 範囲の下限値(min)、上限値(max)、由来(source)に関する情報をま とめています
- セクション4で作成したデータセット「select.70」の18種類の検査 項目の基準範囲も含まれています

desctiption	item	unit	min	max	source
Total protein	TP	g/dL	6.6	8.1	JCCLS
Albumin	ALB	g/dL	4.1	5.1	JCCLS
Albumin/IgG ratio	AG	nd	1.32	2.23	JCCLS
Blood urea nitrogen	BUN	mg/dL	8	20	JCCLS
Serum creatinine	CRE	mg/dL	0.65	1.07	JCCLS
Uric acid	UA	mg/dL	3.7	7.8	JCCLS
Na	Na	mmol/L	138	145	JCCLS
К	K	mmol/L	3.6	4.8	JCCLS
CI	CI	mmol/L	101	108	JCCLS
•	•	•	•	•	•
:		•	•		•

91

# 検査値データの正規化手順



## STEP1: 検査値データを男女別に分割する

##Step1: 検査データを男女別に分割する data.m<-select.70[select.70\$sexmale==1,] data.f<-select.70[select.70\$sexmale==0,]

# STEP2: 男女別の基準範囲データを読み込む

#Step2: 基準範囲データの読み込み #男女別の検査値ごとの基準範囲(95%区間の上限値、下限値)を記述 ri.m<-read.csv("reference interval.m.csv",header=T) #男性の基準範囲 ri.f<-read.csv("reference interval.f.csv",header=T) #女性の基準範囲								
#中身を確認しましょう #検査値表記, 略称, 単位, 下限値, 上限値, 参照値が含まれています str(ri.m)								
<pre>&gt; str(ri.m)</pre>								
'data.frame': 43 obs. of 6 variables:								
<pre>\$ desctiption: chr "Total protein" "Albumin" "Albumin/IgG ratio" "Blood urea nitrogen"</pre>								
\$ item : chr "TP" "ALB" "AG" "BUN"								
\$ unit : chr "g/dL" "g/dL" "nd" "mg/dL"								
\$ min : num 6.6 4.1 1.32 8 0.65 3.7 138 3.6 101 8.8								
\$ max : num 8.1 5.1 2.23 20 1.07 7.8 145 4.8 108 10.1								
\$ source : chr "JCCLS" "JCCLS" "JCCLS"								
\ 基準範囲の下限値(min)と上限値(max)が含まれています								

# STEP3: 対数正規分布のパラメータを逆算する

#Step3: 対数正 #対数正規分布の #男性の基準範囲 ri.m\$avr<-(log ri.m\$sd<-(log( #女性の基準範囲 ri.f\$avr<-(log( ri.f\$sd<-(log(ri str(ri.m)	規分布のパラメータを逆算する D平均値とSD値を求め、新たな変数として格納する 国データを変換する (ri.m\$min)+log(ri.m\$max))/2) #log()は自然対数に変換します ri.m\$max)-log(ri.m\$min))/3.92) 国データを変換する ri.f\$min)+log(ri.f\$max))/2) i.f\$max)-log(ri.f\$min))/3.92)
• str(ri.m)	
data.frame': 43 ol	ps. of 8 variables:
<pre>\$ desctiption: chr</pre>	"Total protein" "Albumin" "Albumin/IgG ratio" "Blood urea nitrogen"
\$ item : chr	"TP" "ALB" "AG" "BUN"
\$ unit : chr	"g/dL" "g/dL" "nd" "mg/dL"
\$ min : num	6.6 4.1 1.32 8 0.65 3.7 138 3.6 101 8.8 対数正規分布の平均値avr
\$ max : num	8.1 5.1 2.23 20 1.07 7.8 145 4.8 108 10.1/ 外致正況分布の標準備差別
\$ source : chr	"JCCLS" "JCCLS" "JCCLS"
\$ avr : num	1.989 1.52 0.54 2.538 -0.182
\$ sd : num	0.0512 0.0546 0.1311 0.2291 0.1246
	95





"AST"の部分を name[*i*]を使って書き換えるとforループで処理できる

# STEP4: forループを使って検査データを正規化する

##Step4: 基準範囲による正規化 #新しい男女別データフレームを用意しましょう std.m<-data.m std.f<-data.f

#forループで検査値ごとの対数正規分布に従って正規化します for(i in 1:length(name)){

#男性データの変換 std.m[,name[i]]<-((log(std.m[,name[i]])-ri.m[ri.m\$item==name[i],"avr"])/ri.m[ri.m\$item==name[i],"sd"]))

#女性データの変換 std.f[,name[i]]<-((<mark>log</mark>(std.f[,name[i]])-ri.f[ri.f\$item==name[i],"avr"])/ri.f[ri.f\$item==name[i],"sd"]))

}

#解読できましたか? #今まで学んだことを総動員したら理解できるはずです! #検査データの検査値名と基準範囲データのitemの名前は完全一致させる必要がありますので、注意してください

97

# STEP5: データを統合し、並び替える

##Step5

#dplyr::bind\_rows関数を使い、男女別の標準化データを統合する std <- dplyr::bind\_rows(std.m, std.f)

#dplyr::arrange関数を使い、日付の昇順に並び替える test <- dplyr::arrange(std, day) #以上の操作により元のデータを検査値のみ標準化したデータセットを得るこ とができる 結果の観察と考察

_							$\pi c$	った	杏	はう	<u> </u>	内							7
>	head	l(se	lect	.70[,	5:22	],6)	760		ш	但 /		^							
	ΤP	ALB	BUN	CRE	UA	Tbil	AST	ALT	ALP	GGT	GLU	Na	К	C1	CR	P LDF	I AG	eGFR	
1	6.3	4.6	12	0.62	2.9	0.9	15	10	84	NA	95	141	3.5	107	0.0	5 125	2.71	93	
2	4.4	2.1	45	1.43	3.5	2.3	33	26	762	NA	118	139	3.2	105	1.8	6 315	0.91	34	
3	4.6	2.3	35	2.21	5.0	1.1	52	26	242	NA	195	135	3.5	102	21.9	8 383	1.00	23	
4	4.5	2.7	7	0.58	2.3	1.8	148	561	189	232	120	143	3.9	106	0.1	0 130	1.50	115	
5	6.2	3.4	24	0.59	NA	0.7	17	25	207	NA	NA	134	4.2	101	0.0	3 133	1.21	75	
6	8.7	1.9	22	3.09	7.1	2.1	61	24	375	NA	119	142	3.4	93	7.1	1 304	0.28	18	正担化データ
>	head	l(st	d.se]	lect.	70[,5	5:22]	,6)												
	٦	ГР	ALB	BUI	N (	CRE	UA	Tbil	. AS	ST	ALT	ALF	P G(	GT (	GLU	Na	К	C1	CRP LDH AG eGFR
1	-0.2	25 -0	0.35	1.3	51	.46	0.17	0.47	1.3	30 (	0.38	0.42	2 1	NA 1	.12 -	0.02	-0.02	0.01	157.00 1.17 -0.14 -0.52
2	-0.2	26 -0	0.27	-0.0	5 -0	.07 -	0.18	0.93	4.5	50 1	3.12	0.26	5 3.4	42 0	.43	0.03	0.06	0.05	0.71 0.03 0.08 0.50
3	0.2	26 -	0.43	0.70	ð 2	.28	0.44	1.13	3 1.6	50 (	0.33	0.84	1_1	NA Ø	.42	0.03	-0.04	-0.07	50.79 0.81 -0.47 -0.58
4	1	١A	NA	N	4	NA	NA	NA	1 1	٨	NA	NA	3.4	44	NA	NA	NA	NA	NA NA NA NA
5	-0.0	)5 -	0.24	0.50	0 0	.42	0.17	0.07	0.7	77 (	0.36	NA	1 1	٨N	NA	0.01	0.19	0.05	79.29 1.40 -0.20 -0.19
6	-0.1	LØ -	0.24	0.1	50	.16	0.13	0.47	0.2	27 -0	0.02	1.00	0 0.7	73	NA -	0.02	0.15	0.01	39.79 1.64 -0.14 0.08

元の検査値データは、基準範囲を見比べないと異常値かどうか判断できません。それに対して、正規化データでは、+2よりも離れた値が異常高値ですので、ひと目見てでわかります。と同時に、異常の程度もわかるはずです。

99

### セクション5のまとめ

#### 病院検査値データの正規化

- □ 個々の検査値の分布は正規分布している場合とそうでない場合 があることを理解した
- □ 検査値をべき乗(対数)変換すると正規分布することを理解した
- □ 検査値の基準範囲 (reference interval)の成り立ちを理解した
- □ 基準範囲の上限値・下限値に基づき、検査値を正規化する手順 について理解した
- □ 正規化した検査値の解釈を理解した



# セクション6:検診データの偏差値化

101

### このセクションの目的

このセクションでは、検診データの検査値を偏差値を使って表示させる 方法を解説します

#### 健康診断の結果

偏差值?





# このセクションで使用するデータ

- 1. ある男性の2011年から2020年までの健康診断結果(生化学検査 と血球検査のみを抜粋)
  - 「laboratory test.csv」ファイル
  - 空白期間がありますが、基本的に年2回の健康診断結果です
  - 19項目の臨床検査項目を含みます
  - 空白は欠損値(未測定のため)
- 2. 検診結果に記載されている基準範囲?のデータ(一部改変)
   \*健診の「基準範囲」は臨床検査分野で使われている「共用基準範囲」とは 少し異なりますので、注意してください。
  - 「reference interval.kenko zaidan.csv」ファイル
  - セクション5で使用した男性用基準範囲データの一部を検診 データに置き換えています(source=kenko zaidan)
  - CRNの下限値はJCCLSの下限値に変更しました
  - eGFRの基準範囲は70-90に設定しました(便宜上)
  - UAの基準範囲は7.0以下ですが、下限値をJCCLS値に設定しました

検診データの偏差値化手順



105

# STEP1: 検診データと基準範囲データを読み込む

#Step1:検診データと基準範囲データの読み込み
data<-read.csv("laboratory test.csv", header=T)
ri<-read.csv("reference interval.kenko zaidan.csv", header=T)
head(data, 6)
head(ri, 6)

>	data	

	10	order	aate	sexmale	age	RBC	HD	Hτ	MRC	PLI	AST	ALI	661	IDIL	ALP	IG	HUL
1	61748341	202082	2011/3/10	1	32	524	14.9	47.3	4.710	NA	25	22	25	NA	NA	78	61
2	61748341	106065	2012/6/19	1	33	525	15.6	47.7	5.000	NA	26	21	25	NA	NA	79	53
3	61748341	153980	2012/12/5	1	33	540	16.1	49.1	4.340	NA	22	22	39	NA	NA	110	59
4	61748341	478382	2013/6/12	1	34	503	15.1	46.5	5.480	NA	27	23	30	NA	NA	122	51
5	61748341	482283	2013/12/2	1	34	523	15.6	47.5	5.650	NA	33	40	60	NA	NA	104	58
6	61748341	714751	2016/6/17	1	37	516	15.4	47.8	6.340	19.9	29	22	32	NA	NA	84	58
> ri																	
				descti	ptic	on i	tem c	lass			un	it	min	m	ax		source
1				Total pr	otei	.n	TΡ	1			g/	dL	6.60	8.	10		JCCLS
2				A	bumi	.n	ALB	2			g/	dL	4.10	5.	10		JCCLS
3			Alb	umin/IgG	rati	.0	AG	3			1	nd	1.32	2.	23		JCCLS
4		As	spartate am	inotransf	eras	se	AST	4			U,	/L :	10.00	40.	00 k	enko	zaidan
5			Alanine am	inotransf	eras	se	ALT	5			U,	/L	5.00	45.	00 k	enko	zaidan
6		gc	amma-Glutam	yl transf	eras	se	GGT	6			U.	/L	5.00	79.	00 k	enko	zaidan

# STEP2:対数正規分布のパラメータを逆算する

#Step2: 対 ri\$avr<-(( ri\$sd<-((  str(ri)	数正 (log( log(r	規分布のmとSDを逆算する(スライド95と同様です) ri\$min)+log(ri\$max))/2) i\$max)-log(ri\$min))/3.92)
> str(ri)		
'data.frame':	43 o	bs. of 9 variables:
<pre>\$ desctiption</pre>	: chr	"Total protein" "Albumin" "Albumin/IgG ratio" "Aspartate aminotransferase"
\$ item	: chr	"TP" "ALB" "AG" "AST"
<pre>\$ class</pre>	: int	1 2 3 4 5 6 7 8 9 10
\$ unit	: chr	"g/dL" "g/dL" "nd" "U/L"
\$ min	: num	6.6 4.1 1.32 10 5 5 0.4 106 8 3.7
\$ max	: num	8.1 5.1 2.23 40 45 79 1.5 322 20 7
\$ source	: chr	"JCCLS" "JCCLS" "JCCLS" "kenko zaidan"
\$ avr	: num	1.99 1.52 0.54 3 2.71
\$ sd	: num	0.0522 0.0557 0.1338 0.3536 0.5605
` <b>\</b>		
計算結果	果が打	<b>挿入されています</b>

107

# STEP3: forループ処理で検診データを偏差値化する

# 結果の観察

	#偏差( head(	值化; <b>std,</b>	された <sup>7</sup> 6)	検診テ	<u>-^</u>	タ	を着	観察	察し	_ <del>11</del> 0	εL	, L	う											
>	head(std	,6)																						
	id	order	date	sexmale	age		RBC		Н	b		Ht		WBC	PL	T	Д	ST		ALT		GGT	Tbil	ł
1	61748341	202082	2011/3/10	1	32	55.89	9672	41.	7692	0 55	5.070	604	41.8	1836	Ν	NA 5	56.309	79	56.	83285	53.25	5860	NA	
2	61748341	106065	2012/6/19	1	33	56.10	5788	47.	9855	3 56	5.393	389	44.1	1607	١	NA 5	57.418	83	56.0	00290	53.25	5860	NA	
3	61748341	153980	2012/12/5	1	33	60.17	7443	52.	2572	9 60	0.920	089	38.6	7219	١	NA 5	52.695	07	56.	83285	59.57	7441	NA	
4	61748341	478382	2013/6/12	1	34	50.07	7959	43.	5746	1 52	2.40	657	47.6	4117	١	NA 5	58.486	00	57.0	62590	55.84	1809	NA	
5	61748341	482283	2013/12/2	1	34	55.62	2505	47.	9855	3 55	5.730	635	48.8	1600	١	NA 6	54.160	33	67.4	49867	65.69	9275	NA	
6	61748341	714751	2016/6/17	1	37	53.70	0863	46.	2383	7 56	5.72	162	53.2	4696	44.24	12 6	50.506	64	56.	83285	56.76	6472	NA	
	#偏差 <sup>4</sup> #apply std[, 6 head(	値の y関数 6:24] (std)	小数部 友を使い  <-app	分はa ヽ、偏 oly (s <sup>.</sup>	るま 差( <b>td</b> [	い 値音 <b>,6</b>	意「 『分 <b>:24</b>	味: }に <b>!],</b>	かれ Ero <b>2,</b>	นท ro	, \∂ Id ₪ <b>un</b>	りて 関数 nd)	ご、 次を	丸	めま	E ( さ-	ンよせま	つす	_					
>	head(std	,6)																						
	id	order	date	sexmale	age	RBC	Hb	Ht V	WBC I	PLT	AST	ALT	GGT	Tbil	ALP	ТG	HDL	LDL	тс	GLU	HbA1c	CRE	eGFR	UA
1	61748341	202082	2011/3/10	) 1	. 32	56	42	55	42	NA	56	57	53	NA	NA NA	46	54	42	39	NA	41	NA	NA	66
2	61748341	106065	2012/6/19	) 1	. 33	56	48	56	44	NA	57	56	53	NA	NA NA	47	46	39	32	NA	44	NA	NA	51
3	61748341	153980	2012/12/5	5 1	. 33	60	52	61	39	NA	53	57	60	NA	NA NA	59	52	47	49	NA	47	NA	NA	70
4	61748341	478382	2013/6/12	2 1	. 34	50	44	52	48	NA	58	58	56	NA	NA	62	44	43	34	NA	49	NA	NA	65
5	61748341	482283	2013/12/2	: 1	. 34	56	48	56	49	NA	64	67	66	NA	NA NA	57	51	58	54	NA	47	NA	NA	70
6	61748341	714751	2016/6/17	' 1	. 37	54	46	57	53	44	61	57	57	NA	NA NA	49	51	55	50	NA	54	NA	NA	55

109

# STEP4: 値に応じてグレードを付与する

#偱	幕差値の値に応じ	て下記のよう	うにグレ	ードを付与します	
	偏差値の範囲	グレード	割合(	6)	
	<30	Low	2.3	<u>スライド42, 70</u>	
	30 <= 偏差値 < 35	Normal Low	4.4	tidyr::spread	
	35 <= 偏差値 =< 65	Normal	86.6		
	65<偏差値 =< 70	Normal High	4.4	tidyr::gather	
	>70	High	2.3		
# C libr std std hea	このため、横持ち rary(tidyr) l.ts<-tidyr::gather l.ts\$date<-as.Dat ad(std.ts,6)	から縦持ちう (std, key="if æ(std.ts\$dat	グレード       割合(%)         Low       2.3         Normal Low       4.4         Normal B6.6       Image: Second State of the second state of t		
> he <ol> <li>1 61</li> <li>2 61</li> <li>3 61</li> <li>4 61</li> <li>5 61</li> <li>6 61</li> </ol>	ad(std.ts) id order dat 748341 202082 2011-03- 748341 106065 2012-06- 748341 153980 2012-12-0 748341 478382 2013-06- 748341 482283 2013-12-0	e sexmale age it 10 1 32 F 19 1 33 F 15 1 33 F 12 1 34 F 12 1 34 F	em score BC 56 BC 56 BC 60 BC 50 BC 56 BC 54	この偏差値に対してグレードを付与します	

# STEP4: 値に応じてグレードを付与する

#ifelse関数を使ってラベル化します(スライド29を参照) #dplyrパッケージのmutate関数を使って新たな変数gradeを挿入します library(dplyr) std.ts<-dplyr::mutate(std.ts, grade=ifelse(score<30, "Low (<30)", ifelse(score<35, "Normal Low (30-35)", ifelse(score<65, "Normal (35-65)", ifelse(score<70, "Normal High (65-70)", "High (>70)"))))) head(std.ts, 6)

>	head	(std	+s.	6)

		,.,								
	id	order	date	sexmale	age	item	score		grade	
1	61748341	202082	2011-03-10	1	32	RBC	56	Normal	(35-65)	らやんと竹与されています
2	61748341	106065	2012-06-19	1	33	RBC	56	Normal	(35-65)	Γ
3	61748341	153980	2012-12-05	1	33	RBC	60	Normal	(35-65)	
4	61748341	478382	2013-06-12	1	34	RBC	50	Normal	(35-65)	
5	61748341	482283	2013-12-02	1	34	RBC	56	Normal	(35-65)	
6	61748341	714751	2016-06-17	1	37	RBC	54	Normal	(35-65)	

111

# STEP5: グラフ化

#最後はggplot2を使ったグラフ化ですが、詳細は省略します \*ggplot2を使ったグラフ化はスライド118の参考文献(特に下2冊)をご参照ください

#コメントのみ残しておきます
#検査値のグルーピングに関するラベルを読み込む
ri.class<-ri[, 2:3]
#縦持ちの偏差値データと連結する
std.join<-dplyr::left\_join(std.ts, ri.class, by="item")</pre>

#ggplot2のgeom\_tile関数でヒートマップ化する library(ggplot2); library(scales)

 $\begin{array}{l} graph <-ggplot(std.join, aes(x=date, y=reorder(item, -class), fill=grade)) \\ graph <-graph + \\ scale_x_date(date_breaks="6 month", labels=date_format("%Y-%b")) + \\ geom_tile(color="white", size=.4) + \\ theme(axis.text.x=element_text(angle=60, hjust=1, vjust=1)) + \\ scale_fill_discrete(limits=c("High (>70)", "Normal High (65-70)", "Normal (35-65)", "Normal Low (30-35)", "Low")) \\ \end{array}$ 

graph





パッと見て、正常高値、正常低値、異常値の判別ができる! "患者フレンドリー"な検査結果の表現方法といえる!

### セクション6のまとめ

#### 検診データの偏差値化

- 検診結果をよりわかりやすい形に変換したいという講演者の意図を理解した
- □ 検査値を正規化し、偏差値化することができる
- □ 偏差値結果を縦持ちに変換し、それぞれの値にグレードを付与 することができる



## 臨床検査データ解析を成功させるための3箇条





117

# もっと知りたい人のために



9年前に読んだR本です。 著者は「**Rはあなた方** が考えつくことは何で も実現できます」と述 べています。 私も同感です。



辞書代わりに重宝します



1,3章は本講習会にピッ タリ グラフ化は4章を参照 してください



ggplot2を使う際、 辞書代わりに重宝します

# 本セミナーで使用した関数一覧(A-Z)

関数	由来	用途	出現スライド				
!is.na()	標準	欠損値ではない(!=notの意)判定	74				
apply	標準	並列演算処理	53,109	left_join()	dplyr	データの連結	112
arrange()	dplyr	並び替える	98	length()	標準	ベクトルの長さを求める	77,97,108
as.character()	標準	文字化	69	library()	標準	パッケージの実装	7,35,58
as.data.frame()	標準	データフレーム化	70,72	list.files()	標準	ファイル名の取得	63
as.Date()	標準	日付化(年月日)	69,110	list()	標準	リストの作成	47
as.numeric()	標準	数値化	69	log()	標準	自然対数を算出する	84,95,97,107,108
bind_rows()	dplyr	リストの要素を縦方向に連結する	34,66,98	mutate()	dplyr	列の挿入	26,29,111
class()	標準	オブジェクトの型を確認する	17	na.summary()	自作	変数の欠損率を算出	72
colSums()	標準	列の総和を算出	51,52,72,74	names()	標準	データの変数名を表示	19,54,67,72,74,75,81,1
factor()	標準	因子化	55	nrow()	標準	行数をカウントする	74
filter()	dplyr	特定の行を抽出する	24	par()	標準	グラフィック環境調整パラメータ	82,84
for()	標準	ループ処理	54,82,97,108	path.packages()	標準	パッケージの保存場所の確認	7
full_join()	dplyr	データの連結	40	read_excel()	readxl	エクセルファイルの読み込み	64
gather()	tidyr	縦持ち変換	42,43,110	read.csv()	標準	csvファイルの読み込み	12,33,35,38,43,75,94,1
getwd()	標準	現在の作業ディレクトリを表示	11	round()	標準	四捨五入する	72,109
ggplot()	ggplot2	データの可視化	112	select()	dplyr	変数の選択、削除	28
head()	標準	データの上位一部表示	18,99,106,109,110,111	setwd()	標準	作業ディレクトリの指定	11
hist()	標準	ヒストグラムを描画	81,82,84	spread()	tidyr	縦持ちデータを横持ち化	42,44,70
ifelse()	標準	条件分岐	29,111	str()	標準	データ構造を表示	15,43,68,69,71,76,94,95,
inner_join()	dplyr	データの連結	38	sum()	標準	合計値の算出	53,54
install.packages()	標準	バッケージのダウンロード	7	summary()	標準	データの基本統計量を表示	49,50
is.na()	標準	欠損値の判定	47	tail()	標準	データの下位一部表示	18
laash 0	+m-3#	関数を複数のオブジェクトに適用、結	F0 C4 CF	unique()	標準	ユニーク(一意)の値を返す	77
labbia	惊华	果をリストで返す	50,04,05	write.csv()	標準	csvファイルの書き出し	30,72

#関数の使い方は"**?関数名"**で調べてください #tidyr, dplyr, ggplot2などはRStudioのHPに**Cheat Sheet**があります。 非常に便利ですので、ぜひ参考にしてください https://www.rstudio.com/resources/cheatsheets/

119



# Enjoy !

### ご質問等有りましたらいつでもどうぞ

#### dseto@cclm.med.kyushu-u.ac.jp

