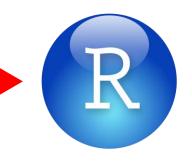
# 有用なR分析技術

① Rの最高のデータ処理のパッケージ"dplyr" 1



# dplyrを使う理由は?

- ①強力な全処理機能!
- →整列、フィルタリング、カラムの選択など簡単にできる
- ②一般的なR関数より直観的で使いやすい!
- →select 、arrange 、filterなど関数名でどのような機能を持っているのか簡単に把握できる
- ③コードの整理が簡単!
- →パイプ記号(%>%)で簡単に整理することもできる(Ctrl+Shift+M)
- 4他のtidyverseパッケージとの連動も簡単!
- →ggplot2などの強力な機能を持つパッケージと同じtidyverseパッケージのため連動が簡単

## 基礎となるdplyrの関数

- 1 filter()
- →提示した条件と一致している資料を表示する
- 2arrange()
- →行(Row)を整列する
- 3select()
- →特定な列(Column)を選択する
- 4mutate()
- →新しい変数を計算・生成する
- **5**summarize()
- →データの統計量を計算し、簡略に表記する

# 練習に使うデータの紹介

#### 以下はとある大学の授業の成績表である

	А	В	С	D	Е
1	Student Number	Name	Major	Midterm	Final
2	1301	Choi	Literature	50	45
3	1302	Park	Engineering	76	42
4	1303	Han	Literature	100	98
5	1304	Jin	Engineering	85	92
6	1305	Liu	Medicine	100	100
7	1306	Li	Engineering	86	100
8	1307	Yamaoka	Social Science	75	90
9	1308	Sirasaki	Natural Science	88	90
10	1309	Honda	Natural Science	92	88
11	1310	Yamada	Social studies	45	0
12	1311	Nguyen	Literature	100	92
13	1312	Lee	Medicine	85	77
14	1313	Hong	Engineering	82	90
15	1314	Hua	Natural Science	56	42
16	1315	Ма	Natural Science	87	78
17	1316	Okusora	Literature	75	62
18	1317	Tu	Engineering	83	55
19	1318	Satou	Social studies	90	82

*	Student <sup>‡</sup> Number	Name ‡	‡ Major	# Midterm	‡ Final
1	1301	Choi	Literature	50	45
2	1302	Park	Engineering	76	42
3	1303	Han	Literature	100	98
4	1304	Jin	Engineering	85	92
5	1305	Liu	Medicine	100	100
6	1306	Li	Engineering	86	100
7	1307	Yamaoka	Social Science	75	90
8	1308	Sirasaki	Natural Science	88	90
9	1309	Honda	Natural Science	92	88
10	1310	Yamada	Social studies	45	0
11	1311	Nguyen	Literature	100	92
12	1312	Lee	Medicine	85	77
13	1313	Hong	Engineering	82	90
14	1314	Hua	Natural Science	56	42
15	1315	Ma	Natural Science	87	78
16	1316	Okusora	Literature	75	62
17	1317	Tu	Engineering	83	55
18	1318	Satou	Social studies	90	82

- ①filter()関数で条件を付けるためには演算子が必要
- →四則演算子 (+, -, x, /)
- →論理演算子(And:&, Or:│, Not:!,
  Same:==, Not same:!=,
  大きさの比較(>, < , >=, <=))
- ②複数の条件を付けて検索することも可能(パイプ記号(%>%)を使うことをおすすめ)

レベル1: MajorがEngineeringの受講生を検索せよ

```
> df %>% dplyr::filter(Major == 'Engineering')
# A tibble: 5 × 5
  `Student Number` Name
                     Major
                                 Midterm Final
           <dbl> <chr> <chr>
                               <dbl> <dbl>
            1302 Park Engineering
                                      76
                                           42
            1304 Jin
                      Engineering
                                  85
                                           92
            1306 Li
                      Engineering
                                  86
                                          100
            1313 Hong Engineering
                                      82
                                           90
                      Engineering
                                      83
            1317 Tu
                                           55
```

レベル2:MajorがEngineering以外の受講生を検索せよ

> 4£ 8>8 4	plyr::fil	ton(Majon	Territoria de la constantidad de			一 しまえし
	: 13 × 5	cer(najor	Tinganeer and	1		→!=と書くと
						自動変換
Studen	t Number`	Name	Major	Midterm	Final	口到久沃
	<dbl></dbl>	<chr></chr>	<chr></chr>	<dbl></dbl>	<dbl></dbl>	
1	<u>1</u> 301	Choi	Literature	50	45	
2	<u>1</u> 303	Han	Literature	100	98	
3	<u>1</u> 305	Liu	Medicine	100	100	
4	<u>1</u> 307	Yamaoka	Social Science	75	90	
5	<u>1</u> 308	Sirasaki	<b>Natural Science</b>	88	90	
6	<u>1</u> 309	Honda	<b>Natural Science</b>	92	88	
7	<u>1</u> 310	Yamada	Social studies	45	0	
8	<u>1</u> 311	Nguyen	Literature	100	92	
9	<u>1</u> 312	Lee	Medicine	85	77	
10	<u>1</u> 314	Hua	<b>Natural Science</b>	56	42	
11	<u>1</u> 315	Ma	<b>Natural Science</b>	87	78	
12	<u>1</u> 316	Okusora	Literature	75	62	
13	<u>1</u> 318	Satou	Social studies	90	82	

レベル3: MajorがLiterature以外で、Midtermが85点以上、Finalが80点以 上の受講生を検索せよ >=と書くと

```
自動変換
 df %>% dplyr::filter(Major ≠ 'Literature' & Midterm ≥ 85 & Final ≥ 80)
  A tibble: 6 \times 5
  Student Number` Name
                          Major
                                          Midterm Final
            <dbl> <chr> <chr>
                                           <dbl> <dbl>
             1304 Jin
                          Engineering
                                              85
                                                    92
                          Medicine
             1305 Liu
                                             100
                                                  100
                          Engineering
             1306 Li
                                              86
                                                  100
             1308 Sirasaki Natural Science
                                              88
                                                    90
5
                          Natural Science
             1309 Honda
                                              92
                                                    88
                          Social studies
             1318 Satou
                                              90
                                                    82
```

# dplyrの関数②:arrange()

- ①arrange()でAscending(昇順)もしくはDescending(降順)にデータを整列することができる
- ②特定なColumnを基準に整列することが可能
- ③重複することも可能

# dplyrの関数②:arrange()

#### 受講生のName基準に昇順整列をせよ

> df %>% dplyr:::	arrange	(Name)	)		
# A tibble: 18 ×	5				
`Student Numb	er` Nam	e	Major	Midterm	Final
<di< td=""><td>bl&gt; <ch< td=""><td>r&gt;</td><td><chr></chr></td><td><dbl></dbl></td><td><dbl></dbl></td></ch<></td></di<>	bl> <ch< td=""><td>r&gt;</td><td><chr></chr></td><td><dbl></dbl></td><td><dbl></dbl></td></ch<>	r>	<chr></chr>	<dbl></dbl>	<dbl></dbl>
1 <u>1</u>	301 Cho	i	Literature	50	45
2 <u>1</u>	303 Han		Literature	100	98
3 <u>1</u>	309 Hon	da	<b>Natural Science</b>	92	88
4 <u>1</u>	313 Hon	g	Engineering	82	90
5 <u>1</u>	314 Hua		<b>Natural Science</b>	56	42
6 <u>1</u>	304 Jin		Engineering	85	92
7 <u>1</u>	312 Lee		Medicine	85	77
8 <u>1</u>	306 Li		Engineering	86	100
9 <u>1</u>	305 Liu		Medicine	100	100
_	315 Ma		Natural Science	87	78
	311 Ngu			100	92
12 <u>1</u>	316 Oku	sora	Literature	75	62
13 <u>1</u>	302 Par	k	Engineering	76	42
	318 Sat			90	82
_	308 Sir	asaki	Natural Science	88	90
_	317 Tu		Engineering	83	55
17 <u>1</u>	310 Yam	ada	Social studies	45	0
18 <u>1</u>	307 Yam	aoka	Social Science	75	90

- ①特定なColumnを選択して表示することができる
- ②複数のColumnを選択したり、特定なカラムを除去して示すことも可能

レベル1:Nameカラムを選択せよ

```
> df %>% dplyr::select(Name)
# A tibble: 18 × 1
   Name
   <chr>
 1 Choi
 2 Park
 3 Han
 4 Jin
 5 Liu
 6 Li
 7 Yamaoka
8 Sirasaki
 9 Honda
10 Yamada
11 Nguyen
12 Lee
13 Hong
14 Hua
15 Ma
16 Okusora
17 Tu
18 Satou
```

#### レベル2:NameとMidterm、Finalカラムを選択せよ

```
df %>% dplyr::select(Name, Midterm, Final)
# A tibble: 18 × 3
   Name
            Midterm Final
            <dbl> <dbl>
   <chr>
 1 Choi
                 50
                       45
 2 Park
                       42
 3 Han
                100
                       98
 4 Jin
                       92
                 85
 5 Liu
                100
                      100
 6 Li
                 86
                      100
 7 Yamaoka
                       90
8 Sirasaki
                       90
9 Honda
                 92
                       88
10 Yamada
11 Nguyen
                100
                       92
12 Lee
                       77
13 Hong
                 82
                       90
14 Hua
                       42
15 Ma
                 87
                       78
16 Okusora
                       62
17 Tu
                       55
                 83
18 Satou
                       82
```

#### レベル3: Student Number以外のすべてのカラムを選択せよ

<pre>&gt; df %&gt;% dplyr::select(!`Student Number`)</pre>								
# /								
	Name	Major	Midterm	Final				
	<chr></chr>	<chr></chr>	<dbl></dbl>	<dbl></dbl>				
1	Choi	Literature	50	45				
2	Park	Engineering	76	42				
3	Han	Literature	100	98				
4	Jin	Engineering	85	92				
5	Liu	Medicine	100	100				
6	Li	Engineering	86	100				
7	Yamaoka	Social Science	75	90				
8	Sirasaki	<b>Natural Science</b>	88	90				
9	Honda	<b>Natural Science</b>	92	88				
10	Yamada	Social studies	45	0				
11	Nguyen	Literature	100	92				
12	Lee	Medicine	85	77				
13	Hong	Engineering	82	90				
14	Hua	Natural Science	56	42				
15	Ma	<b>Natural Science</b>	87	78				
16	Okusora	Literature	75	62				
17	Tu	Engineering	83	55				
18	Satou	Social studies	90	82				

### dplyrの関数④:mutate()

- ①新しいカラムを追加するとき使う
- ②すべて同じ内容を含むカラムを入れることも、存在するカラムをベースに 計算した結果を追加することもできる
- ③ifと論理演算子を応用することもできる
- ④条件が3つ以上の場合はcase\_when()関数が便利

# dplyrの関数④:mutate()

#### レベル2: MidtermとFinalの平均を計算するカラムを作成せよ

•	•	ate(Avera	ge = (Midterm+Fi	nal)/2)		
# A tibble:						
`Student	Number`	Name	Major	Midterm	Final	Average
	<dbl></dbl>	<chr></chr>	<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
1	<u>1</u> 301	Choi	Literature	50	45	47.5
2	<u>1</u> 302	Park	Engineering	76	42	59
3	<u>1</u> 303	Han	Literature	100	98	99
4	<u>1</u> 304	Jin	Engineering	85	92	88.5
5	<u>1</u> 305	Liu	Medicine	100	100	100
6	<u>1</u> 306	Li	Engineering	86	100	93
7	<u>1</u> 307	Yamaoka	Social Science	75	90	82.5
8	<u>1</u> 308	Sirasaki	<b>Natural Science</b>	88	90	89
9	<u>1</u> 309	Honda	<b>Natural Science</b>	92	88	90
10	<u>1</u> 310	Yamada	Social studies	45	0	22.5
11	<u>1</u> 311	Nguyen	Literature	100	92	96
12	<u>1</u> 312	Lee	Medicine	85	77	81
13	<u>1</u> 313	Hong	Engineering	82	90	86
14	<u>1</u> 314	Hua	Natural Science	56	42	49
15	<u>1</u> 315	Ma	<b>Natural Science</b>	87	78	82.5
16	<u>1</u> 316	Okusora	Literature	75	62	68.5
17	<u>1</u> 317	Tu	Engineering	83	55	69
18	<u>1</u> 318	Satou	Social studies	90	82	86

## dplyrの関数④:mutate()

レベル4以上!:MidtermとFinalの平均を示す"Average"カラムを作成した後、Averageの数値が60以上は"Pass"、それ以外は"Non-pass"を示す新しいカラムの"Result"を作成せよ

> 0	<pre>&gt; df %&gt;% dplyr::mutate(Average = (Midterm+Final)/2,</pre>								
+		Resul <sup>*</sup>	t = ifelse(Avera	ge ≥ 60	, 'Pass	s', 'Non-	-pass'))		
# #	tibble: 18 × 7								
	`Student Number	Name	Major	Midterm	Final	Average	Result		
	<dbl< td=""><td>&gt; <chr></chr></td><td><chr></chr></td><td><dbl></dbl></td><td><dbl></dbl></td><td><dbl></dbl></td><td><chr></chr></td></dbl<>	> <chr></chr>	<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<chr></chr>		
1	<u>1</u> 30	L Choi	Literature	50	45	47.5	Non-pass		
2	<u>1</u> 30	2 Park	Engineering	76	42	59	Non-pass		
3	<u>1</u> 30	3 Han	Literature	100	98	99	Pass		
4	<u>1</u> 30	l Jin	Engineering	85	92	88.5	Pass		
5	<u>1</u> 30	5 Liu	Medicine	100	100	100	Pass		
6	<u>1</u> 30	5 Li	Engineering	86	100	93	Pass		
7	_	7 Yamaoka			90	82.5	Pass		
8	<u>1</u> 30	3 Sirasaki	Natural Science	88	90	89	Pass		
9	<u>1</u> 30	Honda	<b>Natural Science</b>	92	88	90	Pass		
10	<u>1</u> 31	9 Yamada	Social studies	45	0	22.5	Non-pass		
11	<u>1</u> 31	l Nguyen	Literature	100	92	96	Pass		
12	<u>1</u> 31	2 Lee	Medicine	85	77	81	Pass		
13	<u>1</u> 31	3 Hong	Engineering	82	90	86	Pass		
14	<u>1</u> 31	l Hua	Natural Science	56	42	49	Non-pass		
15	<u>1</u> 31	5 Ma	Natural Science	87	78	82.5	Pass		
16	<u>1</u> 31	6 Okusora		75	62	68.5	Pass		
17	<u>1</u> 31	7 Tu	Engineering	83	55	69	Pass		
18	<u>1</u> 31	3 Satou	Social studies	90	82	86	Pass		

## dplyrの関数⑤:summarize()

- ①数値型データの統計量を計算する
- ②様々な統計オプションがあるが、以下のオプションが主に使われている
  →mean(x, na.rm = TRUE): 平均値の計算(欠損値を含む場合はFALSEに変更)
- **→n():データの個数**
- ③group\_byを使うことでグループを基準として計算することもできる
- ④"summarise"と書いてもOK

# dplyrの関数⑤:summarize()

レベル1:Midtermの平均値をsummarize関数を用いて計算せよ

### dplyrの関数⑤:summarize()

#### レベル3:専攻別MidtermとFinalの平均値を計算せよ

```
dplyr::summarize(group_by(df, Major),
                   Midterm_average = mean(Midterm),
                   Final_average = mean(Final))
 A tibble: 6 \times 3
                  Midterm_average Final_average
 Major
                             <dbl>
                                           <dbl>
  <chr>>
1 Engineering
                             82.4
                                            75.8
2 Literature
                                            74.2
                             81.2
3 Medicine
                                            88.5
                             92.5
4 Natural Science
                                            74.5
                             80.8
5 Social Science
                                            90
                             75
6 Social studies
                              67.5
                                            41
```

### 演習問題(挑戦)

以下の条件に合わせてデータを処理しなさい

Medicine専攻以外の理系(Natural ScienceとEngineering)の場合、平均値が85点以上の場合は奨学金の対象になる。 MidtermとFinalの平均値を示すAverageカラムを生成し、該当カラムも示すこと。また、対象学生を示すScholarshipカラムを生成し、対象学生には"S"を、対象外の学生には"一"を表記すること。最後に、対象学生のStudent number、Name、Major、Averageのみ示すこと。

### 演習問題(挑戦)

```
> df2 ← df %>% dplyr::mutate(Average = (Midterm + Final) / 2) %>%
   dplyr::mutate(Scholarship = ifelse(Average ≥ 85 & Major ≠ 'Medicine', "S", "-")) %>%
   dplyr::filter(Scholarship == 'S') %>%
   dplyr::select(`Student Number`, Name, Major, Average)
 A tibble: 8 × 4
  `Student Number` Name
                        Major
                                          Average
            <dbl> <chr> <chr>
                                            <dbl>
             1303 Han Literature
                                             99
                          Engineering
             <u>1</u>304 Jin
                                          88.5
                          Engineering
             1306 Li
                                             93
             1308 Sirasaki Natural Science
                                             89
             1309 Honda Natural Science
                                             90
                          Literature
                                             96
             <u>1</u>311 Nguyen
             1313 Hong Engineering
                                             86
                          Social studies
             1318 Satou
                                             86
```