





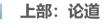
人人都爱tidyverse

艾新波 / 2018 • 北京



课程体系









第3章 格言联璧话学习

■ 第4章 源于数学、归于工程

中部: 执具

- 第5章 工欲善其事必先利其器
- 第6章 基础编程
- 第7章 数据对象





第9章 最美不过数据框





- 第11章 相随相伴、谓之关联

12章 既是世间法、自当有分别

13章 方以类聚、物以群分

第14章 庐山烟雨浙江潮

如果说优雅也有缺点的话,那就是你需要 艰巨的工作才能得到它,需要良好的教育 才能欣赏它

——Edsger Wybe Dijkstra

tidyverse

Tidy datasets are all alike, but every messy dataset is messy in its own way.

— Hadley Wickham

tidy 🔼

英 [ˈtaɪdɪ] 🖒 美 [ˈtaɪdi] 🖒

adj. 整齐的;相当大的

vt. 整理; 收拾; 弄整齐

verse 🗀

英 [V3ːS] 📢 美 [V3·S] 📢

n. 诗 , 诗篇 ; 韵文 ; 诗节

vi. 作诗

tidyverse



tidy = 清爽

tidyxxxx代码界的一股清流

tidyverse

<u>tidybayes</u> I idy Data and 'Geoms' for Bayesian Models <u>tidyboot</u> Tidyverse-Compatible Bootstrapping

tidycensus Load US Census Boundary and Attribute Data as 'tidyverse' and 'sf'-Ready Data Frames

<u>tidygenomics</u> Tidy Verbs for Dealing with Genomic Data Frames

tidygraph A Tidy API for Graph Manipulation

<u>tidyhydat</u> Extract and Tidy Canadian 'Hydrometric' Data

<u>tidyimpute</u> Imputation the Tidyverse Way <u>tidylog</u> Logging for 'dplyr' Functions

tidyLPA Easily Carry Out Latent Profile Analysis (LPA) Using Open-Source or Commercial Software

<u>tidymodels</u> Easily Install and Load the 'Tidymodels' Packages

tidyposterior Bayesian Analysis to Compare Models using Resampling Statistics

tidypredict Run Predictions Inside the Database tidyquant Tidy Quantitative Financial Analysis

<u>tidyqwi</u> A Convenient API for Accessing United States Census Bureau's Quarterly Workforce Indicator

tidyr Easily Tidy Data with 'spread()' and 'gather()' Functions

<u>tidyRSS</u> Tidy RSS for R

<u>tidyselect</u> Select from a Set of Strings tidystats Create a Tidy Statistics Outp

<u>tidystats</u> Create a Tidy Statistics Output File <u>tidystopwords</u> Customizable Lists of Stopwords in 53 Languages

<u>tidytext</u> Text Mining using 'dplyr', 'ggplot2', and Other Tidy Tools

<u>tidytidbits</u> A Collection of Tools and Helpers Extending the Tidyverse

<u>tidytransit</u> Read, Validate, Analyze, and Map Files in the General Transit Feed Specification

<u>tidytree</u> A Tidy Tool for Phylogenetic Tree Data Manipulation

<u>tidyverse</u> Easily Install and Load the 'Tidyverse'

<u>tidyxl</u> Read Untidy Excel Files

tidytext

tidygraph

tidystats

tidyimpute

tidypredict

tidyxl

tidybayes

tidyboot

tidyverse扩展包套装

```
> librarv(tidvverse)
-- Attaching packages ----- tidyverse 1.2.1 -

√ ggplot2 2.2.1

             √ purrr 0.2.4

√ tidyr 0.8.0  
√ stringr 1.3.1

√ readr 1.1.1
 √ forcats 0.3.0

-- Conflicts ----- tidyverse conflicts() -
x dplyr::filter() masks stats::filter()
x dplyr::lag() masks stats::lag()
```

tidyverse扩展包套装

序号	扩展包	功能
1	ggplot2	data visualisation
2	dplyr	data manipulation
3	tidyr	data tidying
4	readr	data import
5	purrr	functional programming
6	tibble	tibbles, a modern re-imagining of data frames
7	stringr	strings
8	forcats	for Categorical Variables (Factors)

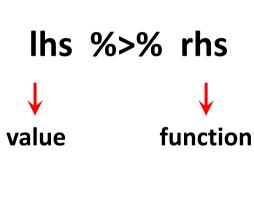
更多内容请参阅: https://www.tidyverse.org/

%>% magrittr

Ceci n'est pas un pipe.

a new "pipe"-like operator, %>%, with which you may pipe a value forward into an expression or function call; something along the lines of x %>% f, rather than f(x)

管道操作符



	表达式	用法
	x %>% f	f(x)
	x %>% f(y)	f(x, y)
)	y %>% f(x, .)	f(x, y)
	z %>% f(x, y, arg = .)	f(x, y, arg = z)

h(g(f(x)))

x %>% f %>% g %>% h

更多用法参见 ?magrittr::`%>%`

查看数据记录

```
cib %>%
 head
#> # A tibble: 6 x 13
#> xm bj xb yw sx wy zz ls dl
#> <chr> <dbl> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <
#> 1 周黎 1101 女
                    94
                         82 96 97
                                      97
                                           98
#> 2 汤海明~ 1101 男
                    87
                         94
                             89 95
                                      94
                                           94
#> 3 舒江辉~ 1101 男
                 92 79
                             86 98
                                      95
                                         96
#> 4 翁柯 1101 女
                       84
                91
                            96 93
                                     97
                                          94
#> 5 祁强
      1101 男
                 85
                       92
                                     87
                            82 93
                                          88
#> 6 淇容 1101 女
                  92
                       82
                           85 91
                                     90
                                          92
#> # ... with 4 more variables: wl <dbl>, hx <dbl>,
#> # sw <dbl>, wlfk <chr>
```

查看数据记录

```
cib %>%
 head(n = 4) #cjb默认为第一个参数
head(cjb, n = 4)
#> # A tibble: 4 x 13
#> xm bj xb yw sx wy zz ls dl
#> <chr> <dbl> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <
#> 1 周黎 1101 女 94
                        82 96 97
                                      97
                                           98
#> 2 汤海明~ 1101 男 87 94 89 95 94 94
#> 3 舒汀辉~ 1101 男
               92 79 86 98 95 96
#> 4 翁柯 1101 女 91 84 96 93 97
                                           94
#> # ... with 4 more variables: wl <dbl>, hx <dbl>,
#> # sw <dbl>, wlfk <chr>
```

dplyr: data manipulation

序号	扩展包	功能
1	select()	picks variables based on their names
2	mutate()	adds new variables that are functions of existing variables
3	filter()	picks cases based on their values
4	arrange()	changes the ordering of the rows
5	summarise()	reduces multiple values down to a single summary

```
cjb %>%
 select(xm, yw, sx) %>%
 head(n = 3)
#> # A tibble: 3 x 3
#> xm
        VW
               SX
#> <chr> <dbl> <dbl>
#> 1 周黎
        94
               82
#> 2 汤海明
         87
              94
#> 3 舒江辉
         92 79
```

```
cjb %>%
 select(xm, yw, sx) %>%
 set names(c("姓名", "语文", "数学")) %>%
 head(n = 3)
#> # A tibble: 3 x 3
#> 姓名 语文 数学
#> <chr> <dbl> <dbl>
#> 1 周黎
        94
                 82
#> 2 汤海明
          87
                  94
#> 3 舒江辉
          92
                 79
```

```
cjb %>%
 select(1, 4:12) %>%
 head(n = 3)
#> # A tibble: 3 x 10
#> xm yw sx wy zz ls dl wl hx
#> <chr> <dbl> <
#> 1 周黎 94 82 96 97 97 98
                                       95
                                           94
#> 2 汤海明~ 87
             94 89 95 94 94 90
                                             90
#> 3 舒江辉~ 92 79 86
                          98 95 96 89 94
#> # ... with 1 more variable: sw <dbl>
```

```
cjb %>%
 select(xm, yw:sw) %>%
 head(n = 3)
#> # A tibble: 3 x 10
#> xm yw sx wy zz ls dl wl hx
#> <chr> <dbl> <
#> 1 周黎 94 82 96
                        97 97 98
                                       95
                                            94
#> 2 汤海明~ 87
             94 89 95 94
                                94 90
                                             90
#> 3 舒江辉~ 92 79 86
                          98 95 96
                                        89
                                             94
#> # ... with 1 more variable: sw <dbl>
```

增加或修改某些列

```
cib %>%
 mutate at(vars(bj, xb, wlfk), factor) %>%
 mutate(zcj = rowSums(.[4:12])) %>%
 arrange(desc(zcj)) %>%
 tail(n = 2)
#> # A tibble: 2 x 14
#> xm bj xb yw sx wy zz ls dl
#> <chr> <fct> <fct> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <
#> 1 滑亚 1113 男 33 46 30 65 82
#> 2 张良平~ 1115 男 0 0 0
#> # ... with 5 more variables: wl <dbl>, hx <dbl>,
#> # sw <dbl>, wlfk <fct>, zcj <dbl>
```

让更改生效

```
#采用%<>%操作符
cjb %<>%
 mutate at (vars (bj, xb, wlfk), factor) %>%
 mutate(zcj = rowSums(.[4:12])) %>%
  arrange(desc(zcj))
cjb <- cjb %>% #和上述语句等价
 mutate at(vars(bj, xb, wlfk), factor) %>%
 mutate(zcj = rowSums(.[4:12])) %>%
  arrange(desc(zcj))
```

选择行

```
cjb %>%
 filter(yw < 60)
#> # A tibble: 2 x 13
#> xm bj xb yw sx wy zz ls dl
#> <chr> <dbl> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <
#> 1 滑亚 1113 男 33 46 30 65 82
                                               76
#> 2 张良平~ 1115 男 0 0
#> # ... with 4 more variables: wl <dbl>, hx <dbl>,
#> # sw <dbl>, wlfk <chr>
```

选择行

```
cjb %>%
 filter at (vars(4:12), any vars(. < 60))
#> # A tibble: 52 x 13
#> xm
       bj xb yw sx wy zz ls
#> <chr> <dbl> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <
#> 1 凌诗雨 1101 女
               84
                         55 95
                                   90
                                        86
#> 2 邓洪涛 1101 男
                    79
                         74 75 95
                                        97
#> 3 罗云 1101 女
                92
                          69
                               85
                                    91 84
#> # ... with 49 more rows, and 5 more variables:
```

xm	xb	yw	SX	wy	ZZ	ls	dl	wl	hx	sw	zcj
汤海明	男	87	94	89	95	94	94	90	90	89	822
舒江辉	男	92	79	86	98	95	96	89	94	87	816
•••	•••		•••	•••	•••	•••		•••	•••		•••
韦永杰	男	81	89	87	97	94	96	81	88	83	796
邰友生	男	88	82	91	89	81	98	89	98	75	791
周黎	女	94	82	96	97	97	98	95	94	88	841
翁柯	女	91	84	96	93	97	94	82	90	83	810
•••		•••	•••	•••	•••	•••	•••	•••			•••
穆伶俐	女	88	72	86	94	87	88	89	98	94	796
龚兰秀	女	88	77	95	94	84	94	87	94	82	795
舒亚	女	94	81	88	91	85	98	81	88	88	794

cjb %>%

```
filter(zcj != 0) %>%
  group by (xb) %>%
                                  韦永杰
                                  邰友生
  summarise(count = n(),
                                  周黎
                                  翁柯
            max = max(zcj),
                                  穆伶俐
            mean = mean(zcj),
                                  龚兰秃
            min = min(zcj)
                                  舒亚
 # A tibble: 2 x 5
#
    хb
          count
                               min
                   max
                        mean
# <fct> <int> <dbl> <dbl> <dbl>
 1 男
                  885 793.
            368
                               523
 2 女
                  879
                       797.
            406
                               647
```

xm

汤海田

舒江辉

94

82 91 89 81 98 89 98 75

72 86 94 87 88 89 98

77

94 82

女 91

 $\boldsymbol{\tau}$

87

816

791

841

810

796

795

794

81

95 94

81 88

90

94 87

97

xm	xb	yw	SX	wy	ZZ	ls	dl	wl	hx	sw	zcj
汤海明	男	87	94	89	95	94	94	90	90	89	822
舒江辉	男	92	79	86	98	95	96	89	94	87	816
祁强	男	85	92	82	93	87	88	95	94	93	809
韦永杰	男	81	89	87	97	94	96	81	88	83	796
邰友生	男	88	82	91	89	81	98	89	98	75	791
周黎	女	94	82	96	97	97	98	95	94	88	841
翁柯	女	91	84	96	93	97	94	82	90	83	810
湛容	女	92	82	85	91	90	92	82	98	90	802
穆伶俐	女	88	72	86	94	87	88	89	98	94	796
龚兰秀	女	88	77	95	94	84	94	87	94	82	795
舒亚	女	94	81	88	91	85	98	81	88	88	794

tidyr:长宽变换

spread(., key, value)

xm	хb	yw	SX	wy
汤海明	男	87	94	89
舒江辉	男	92	79	86
祁强	男	85	92	82
\	,			/
宽的	数据			

xm	хb	key	value
汤海明	男	yw	87
汤海明	男	SX	94
汤海明	男	wy	89
舒江辉	男	yw	92
舒江辉	男	SX	79
舒江辉	男	wy	86
祁强	男	yw	85
祁强	男	SX	92
祁强	男	wy	82

gather(., key, value, yw:wy)

xm	хb	key	value	xm	xb	key	val
汤海明	男	yw	87	汤海明	男	yw	8
汤海明	男	SX	94	舒江辉	男	yw	9
汤海明	男	wy	89	祁强	男	yw	8
舒江辉	男	yw	92	汤海明	男	SX	9
舒江辉	男	SX	79	舒江辉	男	SX	7
舒江辉	男	wy	86	祁强	男	SX	9
祁强	男	yw	85	汤海明	男	wy	8
祁强	男	SX	92	舒江辉	男	wy	8
祁强	男	wy	82	祁强	男	wy	8

```
cjb %>%
 gather(key = ke mu, value = cheng ji, yw:sw) %>%
 arrange (xm)
#> # A tibble: 6,975 x 7
#> xm
           bj
                 \xb
                      wlfk zcj ke mu cheng ji
#> <chr> <fct> <fct> <fct> <dbl> <chr>
                                          <db1>
#> 1 艾春莲
          1103
                       文科
          1103 女
1103 女
                                             86
                           713 yw
                       文科
#> 2 艾春莲
                                             59
                           713 sx
          1103 女
                     文科
#> 3 艾春莲
                              713 wy
                                             87
            1103 女
#> 4 艾春莲
                       文科
                            713 zz
                                             89
            1103 女
#> 5 艾春莲
                       文科
                              713 ls
                                             85
#> 6 艾春莲
            1103
                  女
                       文科
                              713 dl
                                             92
#> # ... with more rows
                  6975 = 775 \times 9
```

```
#按科目讲行汇总统计
cjb %>%
  filter(zcj != 0) %>%
  gather(key = ke mu, value = cheng ji, yw:sw) %>%
  group by (ke mu) %>%
  summarise(max = max(cheng ji),
            mean = mean(cheng ji),
            median = median(cheng ji),
            min = min(cheng ji)) %>%
  arrange (desc (mean) )
```

```
\# # A tibble: 9 x 6
#> ke mu
                    mean median
                                   min
                                            sd
              max
             <dbl> <dbl> <dbl> <dbl> <
#> <chr>
                                          <db1>
              100
                                     70
                                           4.87
#> 1 dl
                    93.0
                              94
#> 2 zz
              100
                    92.3
                              93
                                     65
                                           4.56
#> 3 hx
                                     52
                                          7.60
              100
                    91.7
                              94
                                      0
                                           7.66
#> 4 ls
              100
                    89.1
                              90
                                     30
                                          7.00
#> 5 wy
               99
                    87.5
                              88
#> 6 vw
               96
                    87.4
                              88
                                     33
                                           4.94
#> 7 sw
              100
                    86.4
                              88
                                     55
                                           8.26
#> 8 sx
                                     26
                                           10.5
              100
                    86.2
                              89
#> 9 wl
                              83
                                     21
              100
                    81.2
                                           12.1
```

后话: 优雅的R

hipsteR

hipsteR re-educating people who learned R before it was cool

- 1. Switch to knitr
- 2. Learn Hadley Wickham's packages
- 3. Adopt the pipe operator
- 4. Consider RStudio
- 5. CRAN is huge, and there's also GitHub
- 6.

http://kbroman.org/hipsteR

謝謝聆听 Thank you

教师个人联系方式

艾新波

手机: 13641159546

QQ: 23127789

微信: 13641159546

E-mail: 13641159546@126.com

axb@bupt.edu.cn

地址:北京邮电大学科研楼917室

课程 网址: https://github.com/byaxb/RDataAnalytics



