

Linear Regression Models

P8111

Lecture 02

Jeff Goldsmith
January 21, 2016



THE DEPARTMENT OF
BIostatISTICS




Columbia University
**MAILMAN SCHOOL
OF PUBLIC HEALTH**

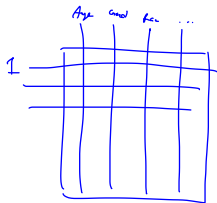
Today's Lecture

- dplyr

Data organization and manipulation

- You're going to spend a lot of time doing this
 - Being better will make your life easier and happier
- 

Data frames



- Most common way of storing dataset in R
- 2D array consisting of named vectors
- Can include variables of many types (numeric, logical, factor, string)
- Used consistently, data frames make analysis easier

Data frames

```
data = data.frame(  
  seq = 1:10,  
  let = letters[1:10],  
  bool = 1:10 < 5  
)
```




Data frames: exploration

```
> dim(data)
```

```
[1] 10  3
```

```
> head(data)
```

```
Source: local data frame [6 x 3]
```

	 seq	 let	 bool
	(int)	<u>(fctr)</u>	<u>(lg1)</u>

1	1	a	TRUE
2	2	b	TRUE
3	3	c	TRUE
4	4	d	TRUE
5	5	e	FALSE
6	6	f	FALSE

Data frames: exploration

```
> summary(data)
```

seq

Min. : 1.00

1st Qu.: 3.25

Median : 5.50

Mean : 5.50

3rd Qu.: 7.75

Max. : 10.00

a

b

c

d

e

f

(Other) : 4

let

:1

:1

:1

:1

:1

:1

bool

Mode :logical

FALSE:6

TRUE :4

NA's :0

Data frames: exploration

↓
data\$seq
summary(data\$seq)

tbl_df: an upgrade to data frames

> data = tbl_df(data) 

> data

Source: local data frame [10 x 3]

	seq (int)	let (fctr)	bool (lg1)
1	1	a	TRUE
2	2	b	TRUE
3	3	c	TRUE
4	4	d	TRUE
5	5	e	FALSE
6	6	f	FALSE
7	7	g	FALSE
8	8	h	FALSE
9	9	i	FALSE
10	10	j	FALSE

tbl_df: an upgrade to data frames

```
> glimpse(data)
```

```
Observations: 10
```


```
Variables: 3
```

```
$ seq  (int) 1, 2, 3, 4, 5, 6, 7, 8, 9, 10
```

```
$ let  (fctr) a, b, c, d, e, f, g, h, i, j
```

```
$ bool (lgl) TRUE, TRUE, TRUE, TRUE, FALSE, FALSE, FALSE
```

dplyr

- dplyr is a new(ish) package for the management and manipulation of data frames built by Hadley Wickham and Romain Francois
 - The package contains several functions focused on the most common tasks – these functions each do one thing, and do it extremely well
 - Functions are designed in a uniformly sensible way: the first argument is a dataframe, and the output is a dataframe
- 

dplyr

- Think of dplyr's functions as verbs: they're actions you want to take on the data
- Arguments to functions clarify the action to take
- Verbs include filter(), arrange(), select(), rename(), mutate(), summarize(), sample_n()
- You should absolutely become fluent in these actions

dplyr : Two Other Things

- Grouping (group_by()) can make some tasks infinitely easier
- The pipe operator (%>%) will change your life

Live coding

Some final thoughts

- You don't know how good you have it
- You can (and should) limit the amount of subsets you save to your workspace
- Many R functions (like `lm`) have `data` and `subset` options, allowing you to pass datasets and trim, or to have these as the last step in a pipe

Today's big ideas

- Intro to `dplyr`
 - Intro to coding
-

- Introduction to `dplyr` (on CRAN)
- Data Wrangling Cheat Sheet
- `swirl` (Getting and Cleaning Data)
- STAT 545 "Basic care and feeding... ", "dplyr: ..."
- Exploratory Data Analysis with R (Managing Data)
- R Programming for Data Science (Ch 13)