



Intro to R Programming for Biostatistics

Day 2 - Cleaning and Transforming Data in R

Adam J Sullivan

<http://adamsullivan.com/teaching/biostats/r-cs-633/>

Piping or Chaining Data

Piping or Chaining

- We will discuss a concept that will help us greatly when it comes to working with our data.
- The usual way to perform multiple operations in one line is by nesting.

library(gapminder)											
head(gapminder)											
## # A tibble: 6 x 6											
##	country	continent	year	lifeExp	pop	gdpPercap	<dbl>				
##	<ctr>	<int>	<dbl>	<int>	<dbl>	<dbl>					
##	1 Afghanistan	Asia	1992	28.801	8425333	779.4453					
##	2 Afghanistan	Asia	1997	30.332	9240934	8208.8530					
##	3 Afghanistan	Asia	1962	31.997	10257063	853.1007					
##	4 Afghanistan	Asia	1987	34.609	11371966	836.1971					
##	5 Afghanistan	Asia	1972	36.488	13079469	719.0811					
##	6 Afghanistan	Asia	1977	38.438	14880372	786.1134					

To consider an example we will look at the data provided in the gapminder package:

Nesting vs Chaining

- Let's say that we want to have the GDP per capita and life expectancy Kenya.
- Traditionally speaking we could do this in a nested manner:

filter(select(gapminder, country, lifeExp, gdpPercap), country=="Kenya")
--



# chaining method  
(x1-x2)\*\*2 %>% sum() %>% sqrt()

- If we did it by hand we would perform elementwise subtraction of x2 from x1 **then** we would sum those elementwise values **then** we would take the square root of the sum.

User Defined Function

- However, if we chain this we can see how we would perform this mathematically.

- The first tidyrr function we will look into is the spread() function.
- With spread() it does similar to what you would expect.
- We have a data frame where some of the rows contain information that is really a variable name.
- This means the columns are a combination of variable names as well as some data.

The spread() Function

```
x1 <- 1:5; x2 <- 2:6  
sqrt(sum((x1-x2)**2))
```

- In the nested manner this would be:

User Defined Function

The spread() Function

$$\sqrt{\sum(x1 - x2)^2}$$

- We could use the math formula:
- Let's say that we wish to find the Euclidean distance between two vectors say, x1 and x2.

User Defined Function

# chaining method  
(x1-x2)\*\*2 %>% sum() %>% sqrt()

## [1] 2.236068

User Defined Function

- Many of us have been performing calculations by this type of method for years, so that chaining really is more natural for us.

Notice that in the column of key, instead of there being values we see the following variable names:

- cases
- population

### The spread() Function

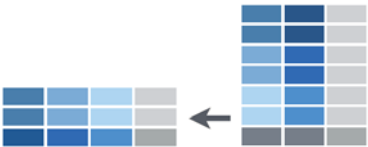
Where

`spread(data, key, value)`

- data is your dataframe of interest.
- key is the column whose values will become variable names.
- value is the column where values will fill in under the new variables created from key.

## # A tibble: 12 x 4	country	year	key	value
## 1	Afghanistan	1999	<ctr>	<int>
## 2	Afghanistan	1999	cases	745
## 3	Afghanistan	2000	cases	2666
## 4	Afghanistan	2000	population	2899360
## 5	Brazil	1999	cases	17717
## 6	Brazil	1999	population	172006362
## 7	Brazil	2000	cases	80488
## 8	Brazil	2000	population	174504898
## 9	China	1999	cases	212258
## 10	China	1999	population	1272915272
## 11	China	2000	cases	213766
## 12	China	2000	population	1280428583

We can consider the following data which is table 2:



The picture below displays this:

## # A tibble: 6 x 4	country	year	cases	population
## 1	Afghanistan	1999	745	19987071
## 2	Afghanistan	2000	2666	2099360
## 3	Brazil	1999	17717	172006362
## 4	Brazil	2000	80488	174504898
## 5	China	1999	212258	1272915272
## 6	China	2000	213766	1280428583

In order to use this data we need to have it so the data frame looks like this instead:

Pipng

If we consider pipng, we can write this as:

data %%
spread(key, value)

25/77

spread() Example

Now if we consider table2 , we can see that we have:

##	country	year	key	value
##	<fctr>	<int>	<fctr>	<int>
## 1	Afghanistan	1999	cases	745
## 2	Afghanistan	1999	population	19987071
## 3	Afghanistan	2000	cases	2666
## 4	Afghanistan	2000	population	2059360
## 5	Brazil	1999	cases	17717
## 6	Brazil	1999	population	17206362
## 7	Brazil	2000	cases	80488
## 8	Brazil	2000	population	174504898
## 9	China	1999	cases	212258
## 10	China	1999	population	1272915272
## 11	China	2000	cases	213766
## 12	China	2000	population	1280428583

26/77

spread() Example

- Now this table was made for this example so key is the key in our spread() function and value is the value in our spread() function.
- We can fix this with the following code:

- We first will load tidyverse.
- If you have not installed it run the following code:
- install.packages("tidyverse")
- Then load this package:
- library(tidyverse)

On Your Own: Rstudio Practice

29/77

spread() Example

- We can now see that we have a variable named cases and a variable named population.
- This is much more tidy.

28/77

spread() Example

table2 %>%
spread(key,value)
## # A tibble: 6 x 4
## country year cases population
## 1 Afghanistan 1999 745 19987071
## 2 Afghanistan 2000 2666 2059360
## 3 Brazil 1999 17717 17206362
## 4 Brazil 2000 80488 174504898
## 5 China 1999 212258 1272915272
## 6 China 2000 213766 1280428583

30/77



### On Your Own: RStudio Practice

```
## # A tibble: 210 x 28
  country
  <chr>
1 Afghanistan 19768973 184871845 199215262 198468616 19967789701
2 Albania 35787588 3410443 3254566 33176294
3 Algeria 2915452 2984280 3004546 3066166 30840825
4 American Samoa 5926 5926 5926 5926 5926
5 Andorra 5354 5354 5354 5354 5354
6 Angola 195252 195252 195252 195252 195252
7 Argentina 19064 19064 19064 19064 19064
8 Aruba 74206 74206 74206 74206 74206
9 Armenia 35670879 35670879 35670879 35670879 35670879
10 Australia 3555555 3555555 3555555 3555555 3555555
  ... with 209 more rows, and 15 unused variables:
  <chr>
1985 1986 1987 1988 1989 1990 1991 1992 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025 2026 2027 2028 2029 2030 2031 2032 2033 2034 2035 2036 2037 2038 2039 2040 2041 2042 2043 2044 2045 2046 2047 2048 2049 2050 2051 2052 2053 2054 2055 2056 2057 2058 2059 2060 2061 2062 2063 2064 2065 2066 2067 2068 2069 2070 2071 2072 2073 2074 2075 2076 2077 2078 2079 2080 2081 2082 2083 2084 2085 2086 2087 2088 2089 2090 2091 2092 2093 2094 2095 2096 2097 2098 2099 2100
```

- You should see the table that we have above, now We have a variable named year, assume that we wish to actually have each year as its own variable.
- Using the `spread()` function, redo this data so that each year is a variable.
- Your data will look like this at the end:

## On Your Own: RStudio Practice

## The gather() Function

- The second tidy<sub>r</sub> function we will look into is the gather<sub>r</sub>() function.
- With gather<sub>r</sub>() it may not be clear what exactly is going on, but in this case we actually have a lot of column names that represent what we would like to have as data values.

### The gather() Function

## On Your Own: RStudio Practice

```
## # A tibble: 1 x 3
##   country year population
##   <chr> <int> <dbl>
## 1 Afghanistan 1995 17586073
```

- In this example we will use the dataset population that is part of tidyverse.
- Print this data:

- This looks similar to the table you created in the spread() practice.
- We now wish to change this data frame so that year is a variable and 1999 and 2000 become values instead of variables.

Table 4

## # A tibble: 3 x 3	##	country	year	cases
## 1 Afghanistan	##	<fctr>	<int>	<int>
## 2 Brazil	##	<fctr>	1999	745
## 3 China	##		2000	212258

Consider table4:

- The gather() Function Example
- For example, in the last spread() practice you created a data frame where variable names were individual years.
- This may not be what you want to have so you can use the gather function.

- You can see that we have created 2 new columns called year and cases.
- We filled these with the previous 2nd and 3rd columns.
- Note that we could have done this in many different ways too.

## # A tibble: 6 x 3	##	country	year	cases
## 1 Afghanistan	##	<fctr>	<chr>	<int>
## 2 Brazil	##		1999	745
## 3 China	##		1999	37737
## 4 Afghanistan	##		2000	2666
## 5 Brazil	##		2000	80488
## 6 China	##		2000	212258

table %>%  
gather("year", "cases", 2:3)

In our example here we would do the following:

gather() Example

- In Comes the gather() Function
- We will accomplish this with the gather function:
- gather(data, key, value, ...)
- data is the data frame you are working with.
- key is the name of the key column to create.
- value is the name of the value column to create.
- ... is a way to specify what columns to gather from.
- where

## # A tibble: 2 x 3				
##	country	year	population	
##	<chr>	<dbl>	<dbl>	
##	1 Afghanistan	1995	17586673	
##	2 Afghanistan	1996	18415107	

## On Your Own: Rstudio Practice

In the end your data frame should look like:

population %<- population %/ spread(year, population)	
population %<- population %/ spread(year, population)	

## On Your Own: Rstudio Practice

- Create population2 from last example:

gather("year", "cases", -1)	
gather("year", "cases", -1)	

- For example if we knew the years but not which columns we could do this:

## The dplyr Package

- filter()
- select()
- arrange()
- mutate()
- summarize()

## dplyr Functionality

- With dplyr we have five basic verbs that we will learn to work with:

- Now that we have started to tidy up our data we can see that we have a need to transform this data.
- We may wish to add additional variables.
- The dplyr package allows us to further work with our data.

## The dplyr Package



## # A tibble: 842 x 19											
##	year	month	day	dep_time	sched_dep_time	dep_delay	arr_time	tailnum	carrier	air_time	distance
##	1	2013	1	1	517	515	2	830			
##	2	2013	1	1	533	529	4	850			
##	3	2013	1	1	542	540	2	923			
##	4	2013	1	1	544	545	-1	1004			
##	5	2013	1	1	554	600	-6	812			
##	6	2013	1	1	554	558	-4	740			
##	7	2013	1	1	555	600	-5	913			
##	8	2013	1	1	557	600	-3	700			
##	9	2013	1	1	557	600	-3	838			
##	10	2013	1	1	558	600	-2	753			
## # ... with 832 more rows, and 12 more variables: sched_arr_time <int>, ## # arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>, ## # origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>, ## # minute <dbl>, time_hour <ttm>											

- We could do this without learning a new command and use indexing which we learned yesterday.
- Let's say that we want to look at the flights data but we are only interested in the data from the first day of the year.

### Filtering Example

### nycflights13 Data

library(dplyr) library(nycflights13)
---

- We also will be using the dplyr package from tidyverse:
- This is a dataset that has all flights in and out of NYC in 2013.
- For the purposes of this example we will consider looking at the package nycflights13.

### dplyr Functionality

- We also will consider:
  - joins
  - group\_by()

### Filtering

- At this point we will consider how we pick the rows of the data that we wish to work with.
- If you consider many modern data sets, we have so much information that we may not want to bring it all in at once.
- R brings data into the RMM of your computer. This means you can bring in at once.
  - Very rarely do you need the entire data set.
- We will focus on how to pick the rows or observations we want now.

**Filtering Example**

- Now this is not very difficult to do, what we have is that we are working with `flights` and we only want to keep the rows of data where `month==1` and `day==1`.
- However we could use the `filter()` function to do this in a much easier to read format:

56/77

```
filter(data, ...)
```

where

- . data is a tibble.
- ... is a set of arguments the data you want returned needs to meet.

56/77

## Filtering Example

- This means in our example we could perform the following:

```
flights %>%  
  filter(month==1, day==1)
```

Finally we could also only do one filtering at a time and chain it:

```
flights %>%  
  filter(month==1) %>%  
  filter(day==1)
```

57/77

## Further Filtering

- `filter()` supports the use of multiple conditions where we can use Boolean.
- For example if we wanted to consider only flights that depart between 0600 and 0605 we could do the following:

```
flights %>% filter(dep_time >= 600, dep_time <= 605)
```

58/77

## Further Filtering

```
## # A tibble: 2,468 x 19  
##   year month   day dep_time sched_dep_time arr_time  
##   <int> <int> <int> <int> <int> <int>  
## 1 2013 1 1 600 600 851  
## 2 2013 1 1 600 600 817  
## 3 2013 1 1 601 600 844  
## 4 2013 1 1 602 618 812  
## 5 2013 1 1 602 605 814  
## 6 2013 1 2 600 600 814  
## 7 2013 1 2 600 605 751  
## 8 2013 1 2 600 600 819  
## 9 2013 1 2 600 600 846  
## 10 2013 1 2 600 600 737  
## # ... with 2,458 more rows, and 12 more variables: sched_arr_time <int>,  
##   arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>,  
##   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,  
##   minute <dbl>, time_hour <ctm>
```

59/77

## Further Filtering

- We can also use the `filter()` function to remove missing data for us.
- Previously we learned about a class of functions called `is.foo()` where `foo` represents a data type.
- We could choose to only use flights that have a departure time.
- That means we wish to not have missing data for departure time:

```
flights %>% filter(!is.na(dep_time))
```

60/77

Selecting

```
flights[, c("dep_time", "arr_time", "flight")]
```

- Let's say that we want to look at the flights data but we are really only interested in the arrival time, departure time and the particular flight number.
- This seems reasonable if we are a customer and wanted to only know these pieces of information. We could do this with indexing:

Select Example

65/77

63/77

## # A tibble: 91,994 x 19											
##	year	month	day	dep_time	sched_dep_time	dep_delay	arr_time	<int>	<int>	<dbl>	<int>
##	1	2013	1	1	517	515	2	830			
##	2	2013	1	1	533	529	4	850			
##	3	2013	1	1	542	540	2	923			
##	4	2013	1	1	554	558	-4	740			
##	5	2013	1	1	558	600	-2	753			
##	6	2013	1	1	558	600	-2	924			
##	7	2013	1	1	558	600	-2	923			
##	8	2013	1	1	559	600	-1	941			
##	9	2013	1	1	559	600	-1	854			
##	10	2013	1	1	606	630	-4	858			
## # ... with 91,984 more rows, and 12 more variables: sched_arr_time <int>,											
##	arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>, hour <dbl>,										
##	origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>										
##	minute <dbl>, time_hour <ctm>										

Your end result should be:

On Your Own: Rstudio Practice

- Choose only rows associated with
  - United Airlines (UA)
  - American Airlines (AA)
- The next logical step would be to select the columns we want as well.
- Many times we have so many columns that we are not interested in for a particular analysis. Instead of slowing down your analysis by continuing to run through extra data, we could just select the columns we care about.

Selecting

64/77

61/77



Removing Columns

- We can also remove columns that contain a certain phrase in the name.
- If we were interested in removing any columns that had to do with time we could search for the word "time" in the data and remove them:

```
flights %>%
  select(-contains("time"))
```

On Your Own: Rstudio Practice

- 1. Select all but the year column.
- 2. Remove the month and day from them.
- 3. Select values which contain "time" in them.
- 4. Chain these together so that you run a command and it does all of these things.

Removing Columns

```
## # A tibble: 336,776 x 13
##   year month day dep_delay arr_delay carrier flight tailnum origin
##   <int> <int> <int> <dbl> <dbl> <chr> <int> <chr> <chr>
## 1 2013 1 1 2 11 UA 1545 N14228 BHM
## 2 2013 1 1 4 20 UA 1714 N14211 LGA
## 3 2013 1 1 2 33 AA 1141 N619AA JFK
## 4 2013 1 1 -1 -18 B6 725 N864JB JFK
## 5 2013 1 1 -6 -25 DL 6611 N666DN LGA
## 6 2013 1 1 -4 12 UA 1696 N19463 EWR
## 7 2013 1 1 -5 -19 B6 587 N516JB EWR
## 8 2013 1 1 -3 -14 EV 5708 N829AS LGA
## 9 2013 1 1 -2 8 AA 301 N361AA LGA
## 10 2013 1 1 -2 8 AA 301 N361AA LGA
## # ... with 336,766 more rows, and 4 more variables: dest <chr>,
##   distance <dbl>, hour <dbl>, minute <dbl>
```

```
## # A tibble: 336,776 x 6
##   dep_time sched_dep_time arr_time sched_arr_time air_time
##   <int> <int> <int> <int> <dbl>
## 1 517 515 529 830 819
## 2 513 529 850 830 830
## 3 542 540 923 850 850
## 4 544 545 1004 850 850
## 5 554 600 812 837 837
## 6 554 740 728 854 854
## 7 555 600 913 854 854
## 8 557 600 709 723 723
## 9 557 600 838 866 866
## 10 558 600 753 745 745
## # ... with 336,766 more rows, and 1 more variable: time_hour <ttm>
##   time_hour
```

Your answer should look like:

On Your Own: Rstudio Practice

Unique Observations

- Many times we have a lot of repeats in our data.
- If we just would like to have an account of all things included then we can use the unique() command.
- Lets assume that we wish to know the origin of a flight and its destination.
- We do not want to have every flight listed over and over again so we ask for unique values:

```
flights %>%
  select(origin, dest) %>%
  unique()
```