

SM-2302 Software for Mathematicians

R3: The tidyverse [handout version]

Dr. Haziq Jamil Mathematical Sciences, Faculty of Science, UBD https://haziqj.ml

Semester I 2022/23

Overview

Tidy data

Tibbles

Pipeline

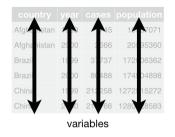
```
Data manipulation
filter
slice
select
relocate
pull
arrange
Others
```

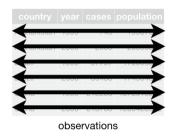
Overview

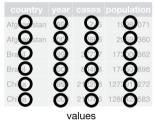
Reshaping data (Wide vs. Long)



Tidy data







From R4DS, tidy data.

Ubd

Tidy vs Untidy

Happy families are all alike; every unhappy family is unhappy in its own way.

—Leo Tolstoy, Anna Karenina

```
tidyr::billboard[, 1:7]
## # A tibble: 317 x 7
##
      artist
                                                date.entered
                                                                wk1
                                                                      wk2
                                                                             wk3
                      track
                                                                                   wk4
      <chr>
                      <chr>>
                                                              <dbl> <dbl> <dbl> <dbl> <dbl>
##
                                                <date>
    1 2 Pac
                      Baby Don't Cry (Keep...
                                                2000-02-26
                                                                 87
                                                                       82
                                                                              72
                                                                                    77
##
##
    2 2Ge+her
                      The Hardest Part Of ... 2000-09-02
                                                                 91
                                                                       87
                                                                              92
                                                                                    NA
##
    3 3 Doors Down
                      Kryptonite
                                                2000-04-08
                                                                 81
                                                                       70
                                                                              68
                                                                                    67
##
    4 3 Doors Down
                                                2000-10-21
                                                                 76
                                                                       76
                                                                                    69
                      Loser
##
                      Wobble Wobble
                                                2000-04-15
                                                                 57
                                                                       34
                                                                              25
                                                                                    17
    5 504 Bovz
##
    6 98^0
                                                2000-08-19
                                                                 51
                                                                       39
                                                                              34
                                                                                    26
                      Give Me Just One Nig...
##
    7 A*Teens
                      Dancing Queen
                                                2000-07-08
                                                                 97
                                                                       97
                                                                              96
                                                                                    95
##
    8 Aaliyah
                      I Don't Wanna
                                                2000-01-29
                                                                 84
                                                                       62
                                                                              51
                                                                                    41
##
    9 Aaliyah
                      Try Again
                                                2000-03-18
                                                                 59
                                                                       53
                                                                              38
                                                                                    28
   10 Adams, Yolanda Open My Heart
                                                2000-08-26
                                                                 76
                                                                       76
                                                                              74
                                                                                    69
     with 307 more rows
```

Ubd

More tidy vs untidy

A			Untid	y Data	3			
species		habitat	weight	ight length latitude/longitude		ongitude	date	
Alligator mississippiensis		swamp	431 lb	4 ft 2	29.531,-82.184		Sept 15, 2015	
Puma concolor		forest	125 lb	2.2m	29.125,-81.682		08/10/2015	
Ursus americanus		forest	88 kg	133 cm	N29°7'30"/W81°40'55.2"		07-13-2015	
ı			Tidy	Data				
	meta	a-data			a	ata	_	
species_code	date		station_code		weight_kg length_cm		7	
TSN 551771	2015-	09-15	1		196	127		
TSN 55247	2015-	08-10	2		57	220		
TSN 180544	2015-	07-13	2		88	133	_	
station_code	habitat	latitude	longitude					
1	swamp	29.531	-82.184					
2	forest	29.125	-81.682					
species_code	class	genus	species					
TSN 551771	Reptilia	Alligato	or mississi	ppiensis				
TSN 55247	Mammalia	a Puma	concolo	r				
TSN 180544	Mammali:	a Ursus	america	nus				

Hart EM, Barmby P, LeBauer D, Michonneau F, Mount S, et al. (2016) Ten Simple Rules for Digital Data Storage. PLOS Computational Biology 12(10): e1005097. DOI: 10.1371/journal.pcbi.1005097

Tidy data

Tibbles

Pipeline

Data manipulation

Reshaping data (Wide vs. Long

Modern data frames

4.7

4.6

5.0

5.4

4.6

5.0

4.4

4.9

5.4

4.8

4.8

4.3

3.2

3.1

3.6

3.9

3.4

3.4

2.9

3.1

3.7

3.4

3.0

3.0

3

4

5

6

7

8

9

10

11

12

13

6##014

Hadley Wickham / RStudio have a package that modifies data frames to be a bit more modern.

```
iris
                                                   (tbl iris = as tibble(iris))
                                                                                                   3 S
```

				_			
##	Sepal.Length Sepal	.Width Peta	1.Length	## #	A tibble: 150	x 5	
## 1	5.1	3.5	1.4	##	Sepal.Length	Sepal.Wi~1	Petal~2 Pet

1.3

1.5

1.4

1.7

1.4

1.5

1.4

1.5

1.5

1.6

1.4

1.1

##	Sepal.Length	Sepal.Width	Petal.Length ##	# A tibble: 150	x 5		
## 1	5.1	3.5	1.4 ##	Sepal.Length	Sepal.Wi~1	Petal~2	Petal~3
## 2	4.9	3.0	1.4 ##	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>

> <

##

##

##

##

##

##

##

##

##

10

#

5.1

4.9

4.7

4.6

5.4

4.6

4.4

4.9

5

5

3.5

3.2

3.1

3.6

3.9

3.4

3.4

2.9

3.1

... with 140 more rows, and abbreviated

variable names 1: Sepal.Width,

3

1.4

1.4

1.3

1.5

1.4

1.7

1.4

1.5

1.4

1.5

0.2 s

0.2 s

0.2 s

0.2 s

0.2 s

0.4 s

0.3 s

0.2 s

0.2 s

0.1 s

Ubd

library(tibble)

Tibbles are lazy

By default, subsetting tibbles always results in another tibble (\$ or [[can still be used to subset for a specific column).

```
tbl_iris[1, ]
                                                            tbl_iris[, 1]
     A tibble: 1 x 5
                                                            ## # A tibble: 150 x 1
 ##
      Sepal.Length Sepal.Width Petal~1 Petal~2 Species
                                                            ##
                                                                  Sepal.Length
             <dbl>
                          <dbl>
                                  <dbl>
                                          <dbl> <fct>
                                                                         <dbl>
 ##
                                                            ##
 ## 1
               5.1
                            3.5
                                    1.4
                                            0.2 setosa
                                                                           5.1
                                                            ##
      ... with abbreviated variable names
                                                            ##
                                                                           4.9
 ## #
        1: Petal.Length, 2: Petal.Width
                                                            ##
                                                                           4.7
                                                            ##
                                                                           4.6
                                                                           5
                                                            ##
 tbl iris[[1]]
                                                            ##
                                                                           5.4
 ##
      [1] 5.1 4.9 4.7 4.6 5.0 5.4 4.6 5.0 4.4 4.9 5.4
                                                                           4.6
                                                            ##
 ##
     [12] 4.8 4.8 4.3 5.8 5.7 5.4 5.1 5.7 5.1 5.4 5.1
                                                            ##
                                                                           5
 ##
     [23] 4.6 5.1 4.8 5.0 5.0 5.2 5.2 4.7 4.8 5.4 5.2
                                                            ##
                                                                           4.4
 ##
     [34] 5.5 4.9 5.0 5.5 4.9 4.4 5.1 5.0 4.5 4.4 5.0
                                                                           4.9
                                                            ## 10
 ##
     [45] 5.1 4.8 5.1 4.6 5.3 5.0 7.0 6.4 6.9 5.5 6.5
                                                            ## # ... with 140 more rows
7 ##
7 / 50
##
     [56] 5.7 6.3 4.9 6.6 5.2 5.0 5.9 6.0 6.1 5.6 6.7
                                                                                              ubăl
     [67] 5 6 5 8 6 2 5 6 5 9 6 1 6 3 6 1 6 4 6 6 6 8
```

More laziness: partial matching

Tibbles do not use partial matching when the \$ operator is used.

```
head( iris$Sp )
## [1] setosa setosa setosa setosa setosa setosa
## Levels: setosa versicolor virginica
head( iris$Species )
## [1] setosa setosa setosa setosa setosa setosa
## Levels: setosa versicolor virginica
head( tbl_iris$Sp )
## Warning: Unknown or uninitialised column: `Sp`.
## NULL.
head( tbl_iris$Species )
## [1] setosa setosa setosa setosa setosa setosa
## Levels: setosa versicolor virginica
```

Ubd

More laziness: stringsAsFactors

Tibbles also have always had stringsAsFactors = FALSE as default behavior.

```
(t = tibble(
                                            (d = data.frame(
 x = 1:3,
                                              x = 1:3.
 y = c("A", "B", "C"),
                                              v = c("A", "B", "C"),
 z = factor(c("X", "Y", "Z"))
                                              z = factor(c("X","Y","Z")).
))
                                              stringsAsFactors = TRUE
                                            ))
## # A tibble: 3 x 3
                                            ## x y z
        x y z
   <int> <chr> <fct>
##
## 1 1 A X
                                            ## 2 2 B Y
## 2 2 B Y
                                            ## 3 3 C Z
## 3 3 C
```

Tibbles and length coercion

Only vectors with length 1 will undergo length coercion. Everything else will throw an error.

```
data.frame(x = 1:4, y = 1)
                                              tibble(x = 1:4, y = 1)
                                              ## # A tibble: 4 x 2
## x y
## 1 1 1
                                                      x y
## 2 2 1
                                              ##
                                                   <int> <dbl>
## 3 3 1
                                              ## 1
 ## 4 4 1
                                              ## 2 2
                                              ## 3 3
                                              ## 4
data.frame(x = 1:4, y = 1:2)
## x y
                                              tibble(x = 1:4, y = 1:2)
## 1 1 1
                                              ## Error:
## 2 2 2
                                              ## ! Tibble columns must have compatible sizes.
## 3 3 1
## 4 4 2
                                              ## * Size 4: Existing data.
                                              ## * Size 2: Column `y`.
                                              ## i Only values of size one are recycled. \cup
10 / 50
```

Tidy data

Tibbles

Pipeline

Data manipulation

Reshaping data (Wide vs. Long

What is a pipe

In software engineering, a pipeline consists of a chain of processing elements (processes, threads, coroutines, functions, etc.), arranged so that the output of each element is the input of the next.

—Wikipedia

Magrittr's pipe is a new infix operator that allows us to link two functions together in a way that is readable from left to right.

The two code examples below are equivalent:

$$f(g(x=1, y=2), n=2)$$
 $g(x=1, y=2) \%\% f(n=2)$

Readability

Consider the following sequence of actions that describe the process of getting to campus in the morning:

I need to find my key, then unlock my car, then start my car, then drive to school, then park.

Expressed as a set of nested functions in R pseudocode this would look like:

```
park(drive(start_car(find("keys")), to = "campus"))
```

Writing it out using pipes give it a more natural (and easier to read) structure:

```
find("keys") %>%
    start_car() %>%
    drive(to = "campus") %>%
    park()
```

Approaches

All of the following are fine, it comes down to personal preference:

Nested:

```
h(g(f(x), y = 1), z = 1)
```

Piped:

```
f(x) %>%
g(y = 1) %>%
h(z = 1)
```

Intermediate:

```
res <- f(x)
res <- g(res, y = 1)
res <- h(res, z = 1)
```

What about other arguments?

Sometimes we want to send our results to an function argument other than first one or we want to use the previous result for multiple arguments. In these cases we can refer to the previous result using '.'.

```
data.frame(a = 1:3, b = 3:1) \% lm(a ~ b, data = .)
 ##
 ## Call:
 ## lm(formula = a \sim b, data = .)
##
 ## Coefficients:
   (Intercept)
 ##
 data.frame(a = 1:3, b = 3:1) %>% .[[1]]
 ## [1] 1 2 3
data.frame(a = 1:3, b = 3:1) %>% .[[length(.)]]
## [1] 3 2 1
14 / 50
```

```
Tidy data
```

Tibbles

Pipeline

```
Data manipulation
filter
slice
select
relocate
pull
arrange
Others
```

Reshaping data (Wide vs. Long)

A Grammar of Data Manipulation

dplyr is based on the concepts of functions as verbs that manipulate data frames.

Core single data frame functions / verbs:

- filter() / slice(): pick rows based on criteria
- select() / rename(): select columns by name
- pull(): grab a column as a vector
- arrange(): reorder rows
- mutate() / transmute(): create or modify columns
- distinct(): filter for unique rows
- summarise() / count(): reduce variables to values
- group_by() / ungroup(): modify other verbs to act on subsets
- relocate(): change column order
- ... (many more)



dplyr rules

- 1. First argument is always a data frame
- 2. Subsequent arguments say what to do with that data frame
- 3. Always return a data frame
- 4. Don't modify in place
- 5. Lazy evaluation magic

Example data

```
library(dplyr)
 library(nycflights13)
flights
## # A tibble: 336,776 x 19
 ##
        vear month
                      day dep_time sched_~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier flight
 ##
       <int> <int> <int>
                             <int>
                                       <int>
                                                <dbl>
                                                        <int>
                                                                 <int>
                                                                          <dbl> <chr>
                                                                                          <int>
        2013
                                517
                                         515
                                                    2
                                                          830
                                                                             11 UA
                                                                                           1545
 ##
                                                                   819
        2013
                                533
                                         529
                                                          850
                                                                   830
                                                                             20 UA
                                                                                           1714
 ##
                                                    4
 ##
        2013
                                542
                                         540
                                                           923
                                                                   850
                                                                             33 AA
                                                                                           1141
 ##
        2013
                                544
                                         545
                                                   -1
                                                          1004
                                                                  1022
                                                                            -18 B6
                                                                                            725
 ##
     5
        2013
                                554
                                         600
                                                   -6
                                                          812
                                                                   837
                                                                            -25 DL
                                                                                            461
 ##
        2013
                                554
                                         558
                                                   -4
                                                          740
                                                                   728
                                                                             12 UA
                                                                                           1696
 ##
        2013
                                555
                                         600
                                                   -5
                                                          913
                                                                   854
                                                                             19 B6
                                                                                            507
 ##
        2013
                                557
                                         600
                                                   -3
                                                          709
                                                                   723
                                                                            -14 EV
                                                                                           5708
     8
 ##
        2013
                                557
                                         600
                                                   -3
                                                          838
                                                                   846
                                                                             -8 B6
                                                                                             79
 ## 10
        2013
                                558
                                         600
                                                   -2
                                                           753
                                                                   745
                                                                              8 AA
                                                                                            301
 ## #
      ... with 336,766 more rows, 8 more variables: tailnum <chr>, origin <chr>,
## #
        dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
 ## #
        time_hour <dttm>, and abbreviated variable names 1: sched_dep_time,
 ## #
        2: dep delay. 3: arr time. 4: sched arr time. 5: arr delay
17 / 50
```

Ubö

filter(): March flights

```
flights %>% filter(month == 3)
## # A tibble: 28,834 x 19
                    day dep time sched ~1 dep d~2 arr t~3 sched~4 arr d~5 carrier flight
##
       vear month
                                                                       <dbl> <chr>
##
      <int> <int> <int>
                            <int>
                                     <int>
                                              <dbl>
                                                      <int>
                                                              <int>
                                                                                       <int>
       2013
                3
                                      2159
                                                125
                                                        318
                                                                  56
                                                                         142 B6
##
                                4
                                                                                          11
       2013
                3
                               50
                                      2358
                                                 52
                                                        526
                                                                438
                                                                          48 B6
                                                                                         707
##
                3
                                                152
                                                                         149 B6
##
       2013
                              117
                                      2245
                                                        223
                                                                2354
                                                                                         608
       2013
                3
                              454
                                       500
                                                                         -15 US
##
                                                 -6
                                                        633
                                                                 648
                                                                                        1117
       2013
                3
                              505
                                       515
                                                        746
                                                                810
                                                                         -24 UA
                                                                                         475
##
                                                -10
       2013
                3
                                        530
                                                 -9
                                                                827
                                                                         -14 UA
                                                                                        1714
##
                              521
                                                        813
       2013
                3
                              537
                                       540
                                                 -3
                                                        856
                                                                850
                                                                           6 AA
                                                                                        1141
##
       2013
                3
                                       545
                                                       1014
                                                                          -9 B6
                                                                                         725
##
    8
                              541
                                                 -4
                                                                1023
                3
                                                        639
##
       2013
                              549
                                       600
                                                -11
                                                                 703
                                                                         -24 US
                                                                                        2114
       2013
                3
                              550
                                       600
                                                                         -14 EV
## 10
                                                -10
                                                        747
                                                                801
                                                                                        4911
     ... with 28,824 more rows, 8 more variables: tailnum <chr>, origin <chr>,
## #
## #
       dest <chr>, air time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
## #
       time_hour <dttm>, and abbreviated variable names 1: sched_dep_time,
## #
       2: dep delay. 3: arr time. 4: sched arr time. 5: arr delay
```

Ubö

filter(): Flights in the first 7 days of March

```
flights %>% filter(month == 3, day <= 7)
## # A tibble: 6,530 x 19
                    day dep_time sched_~1 dep_d~2 arr_t~3 sched~4 arr d~5 carrier flight
##
##
      <int> <int> <int>
                            <int>
                                     <int>
                                              <dbl>
                                                      <int>
                                                              <int>
                                                                       <dbl> <chr>
                                                                                      <int>
       2013
                3
                                      2159
                                                125
                                                        318
                                                                 56
                                                                         142 B6
##
                                4
                                                                                          11
       2013
                3
                                      2358
                                                 52
                                                        526
                                                                          48 B6
                                                                                        707
##
                               50
                                                                438
                3
                                                                         149 B6
##
       2013
                              117
                                      2245
                                                152
                                                        223
                                                               2354
                                                                                        608
       2013
                3
                                       500
                                                                         -15 US
##
                              454
                                                 -6
                                                        633
                                                                648
                                                                                        1117
       2013
                3
                              505
                                       515
                                                        746
                                                                         -24 UA
                                                                                        475
##
                                                -10
                                                                810
       2013
                3
                                       530
                                                                827
                                                                         -14 UA
                                                                                        1714
##
                              521
                                                 -9
                                                        813
       2013
                3
                              537
                                       540
                                                 -3
                                                        856
                                                                850
                                                                           6 AA
                                                                                        1141
##
       2013
                3
                                                                          -9 B6
##
    8
                              541
                                       545
                                                 -4
                                                       1014
                                                               1023
                                                                                        725
                3
                                                        639
##
       2013
                              549
                                       600
                                                -11
                                                                703
                                                                         -24 US
                                                                                       2114
       2013
                3
## 10
                              550
                                       600
                                                -10
                                                        747
                                                                801
                                                                         -14 EV
                                                                                       4911
## #
     ... with 6,520 more rows, 8 more variables: tailnum <chr>, origin <chr>,
## #
       dest <chr>, air time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
## #
       time_hour <dttm>, and abbreviated variable names 1: sched_dep_time,
## #
       2: dep delay. 3: arr time. 4: sched arr time. 5: arr delay
```

ubä

filter(): Flights to LAX or JFK in March

```
flights %>% filter(dest == "LAX" | dest == "JFK", month == 3)
## # A tibble: 1,178 x 19
       year month day dep_time sched_~1 dep_d~2 arr_t~3 sched~4 arr d~5 carrier flight
##
##
      <int> <int> <int>
                           <int>
                                    <int>
                                             <dbl>
                                                     <int>
                                                             <int>
                                                                      <dbl> <chr>
                                                                                     <int>
       2013
                3
                             607
                                       610
                                                -3
                                                       832
                                                                        -53 UA
                                                                                       797
##
                                                               925
       2013
                3
                             629
                                       632
                                                       844
                                                               952
                                                                        -68 UA
                                                                                      1702
##
                                                -3
                3
                                                                        -41 DL
##
       2013
                             657
                                       700
                                                -3
                                                       953
                                                              1034
                                                                                       763
       2013
                3
                                       715
                                                              1037
                                                                        -58 B6
##
                             714
                                                -1
                                                       939
                                                                                       671
       2013
                3
                             716
                                       710
                                                 6
                                                              1035
                                                                        -37 VX
##
                                                       958
                                                                                       399
       2013
                3
                             727
                                       730
                                                              1100
                                                                        -53 AA
                                                                                        33
##
                                                -3
                                                      1007
       2013
                3
                             836
                                       840
                                                      1111
                                                              1157
                                                                        -46 UA
                                                                                      1553
##
                                                -4
       2013
                3
                             857
                                                      1202
                                                              1221
                                                                        -19 DL
##
   8
                                       900
                                                -3
                                                                                       120
                3
                                                 3
                                                              1220
##
       2013
                             903
                                       900
                                                      1157
                                                                        -23 AA
       2013
                3
                                                33
                                                      1150
                                                              1151
## 10
                             904
                                       831
                                                                         -1 UA
                                                                                      1223
## #
     ... with 1,168 more rows, 8 more variables: tailnum <chr>, origin <chr>,
## #
       dest <chr>, air time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
## #
       time_hour <dttm>, and abbreviated variable names 1: sched_dep_time,
## #
       2: dep delay, 3: arr time, 4: sched arr time, 5: arr delay
```

ubā

slice(): First 10 flights

```
flights %>% slice(1:10)
## # A tibble: 10 x 19
                    day dep_time sched_~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier flight
##
##
      <int> <int> <int>
                            <int>
                                     <int>
                                              <dbl>
                                                      <int>
                                                               <int>
                                                                       <dbl> <chr>
                                                                                       <int>
       2013
                              517
                                        515
                                                        830
                                                                 819
                                                                          11 UA
                                                                                        1545
##
       2013
                              533
                                        529
                                                  4
                                                        850
                                                                 830
                                                                          20 UA
                                                                                        1714
##
       2013
                                        540
                                                  2
                                                                          33 AA
##
                              542
                                                        923
                                                                 850
                                                                                        1141
       2013
                                        545
                                                                          -18 B6
                                                                                         725
##
                              544
                                                 -1
                                                        1004
                                                                1022
       2013
                              554
                                        600
                                                        812
                                                                 837
                                                                          -25 DL
                                                                                         461
##
                                                 -6
       2013
                              554
                                        558
                                                        740
                                                                 728
                                                                          12 UA
                                                                                        1696
##
                                                 -4
       2013
                              555
                                        600
                                                 -5
                                                        913
                                                                          19 B6
                                                                                         507
##
                                                                 854
       2013
                              557
                                                                 723
                                                                         -14 EV
                                                                                        5708
##
    8
                                        600
                                                 -3
                                                        709
       2013
                              557
                                                                          -8 B6
##
                                        600
                                                 -3
                                                        838
                                                                 846
                                                                                          79
## 10
       2013
                              558
                                        600
                                                 -2
                                                         753
                                                                 745
                                                                           8 AA
                                                                                         301
     ... with 8 more variables: tailnum <chr>, origin <chr>, dest <chr>,
## #
## #
       air time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time hour <dttm>, and
## #
       abbreviated variable names 1: sched_dep_time, 2: dep_delay, 3: arr_time,
## #
       4: sched arr time. 5: arr delay
```

Uböl

slice(): Last 5 flights

```
flights \%% slice((n() - 4):n())
## # A tibble: 5 x 19
##
      year month day dep_time sched d~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier flight
##
     <int> <int> <int>
                          <int>
                                    <int>
                                            <dbl>
                                                    <int>
                                                            <int>
                                                                    <dbl> <chr>
                                                                                    <int>
                             NA
                                               NA
## 1
     2013
                    30
                                     1455
                                                       NA
                                                             1634
                                                                       NA 9E
                                                                                    3393
                    30
                             NA
                                     2200
                                               NA
                                                                                    3525
## 2
     2013
                                                       NA
                                                             2312
                                                                       NA 9E
## 3
     2013
                    30
                             NA
                                     1210
                                               NA
                                                             1330
                                                                       NA MQ
                                                                                    3461
                                                       NA
      2013
                    30
                             NA
                                     1159
                                               NΑ
                                                       NA
                                                                                    3572
## 4
                                                             1344
                                                                       NA MQ
## 5
      2013
               9
                    30
                             NA
                                      840
                                               NA
                                                       NA
                                                             1020
                                                                                    3531
                                                                       NA MQ
     ... with 8 more variables: tailnum <chr>, origin <chr>, dest <chr>,
## #
## #
       air time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time hour <dttm>, and
## #
       abbreviated variable names 1: sched_dep_time, 2: dep_delay, 3: arr_time,
## #
       4: sched arr time, 5: arr delay
```

slice(): Last 5 flights (cont.)

```
flights \%% slice tail(n = 5)
## # A tibble: 5 x 19
##
      year month day dep_time sched d~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier flight
##
     <int> <int> <int>
                          <int>
                                    <int>
                                            <dbl>
                                                    <int>
                                                             <int>
                                                                     <dbl> <chr>
                                                                                    <int>
                             NA
                                               NA
## 1
     2013
                    30
                                     1455
                                                       NA
                                                              1634
                                                                        NA 9E
                                                                                     3393
                    30
                             NA
                                     2200
                                               NA
                                                                                     3525
## 2
     2013
                                                       NA
                                                              2312
                                                                        NA 9E
## 3
                    30
                             NA
                                     1210
                                               NA
                                                              1330
                                                                        NA MQ
                                                                                     3461
     2013
                                                       NA
      2013
                    30
                             NA
                                     1159
                                               NA
                                                       NA
                                                                                     3572
## 4
                                                              1344
                                                                        NA MQ
## 5
      2013
               9
                    30
                             NA
                                      840
                                               NA
                                                       NA
                                                              1020
                                                                                     3531
                                                                        NA MQ
     ... with 8 more variables: tailnum <chr>, origin <chr>, dest <chr>,
## #
## #
       air time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time hour <dttm>, and
## #
       abbreviated variable names 1: sched_dep_time, 2: dep_delay, 3: arr_time,
## #
       4: sched arr time, 5: arr delay
```

Ubā

select(): Individual Columns

```
flights %>% select(year, month, day)
## # A tibble: 336,776 x 3
##
      year month day
##
      <int> <int> <int>
##
      2013
      2013
##
      2013
##
##
      2013
##
      2013
##
   6 2013
      2013
##
      2013
##
##
      2013
```

Ubö

10

2013

... with 336,766 more rows

select(): Exclude Columns

```
flights %>% select(-year, -month, -day) ## # A tibble: 336,776 x 16
      dep_time sched_dep~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier flight tailnum origin
##
##
         <int>
                       <int>
                               <dbl>
                                        <int>
                                                 <int>
                                                         <dbl> <chr>
                                                                          <int> <chr>
                                                                                         <chr>
##
            517
                         515
                                    2
                                          830
                                                   819
                                                             11 UA
                                                                           1545 N14228
                                                                                         EWR.
                                                             20 UA
##
            533
                         529
                                    4
                                          850
                                                   830
                                                                           1714 N24211
                                                                                         LGA
            542
                                    2
                                                             33 AA
##
                         540
                                          923
                                                   850
                                                                           1141 N619AA
                                                                                         JFK
                         545
                                                            -18 B6
                                                                            725 N804JB
                                                                                         JFK
##
            544
                                   -1
                                         1004
                                                  1022
            554
                         600
                                          812
                                                   837
                                                            -25 DL
                                                                            461 N668DN
                                                                                         LGA
##
                                   -6
            554
                         558
                                          740
                                                   728
                                                             12 UA
                                                                           1696 N39463
                                                                                         EWR
##
                                   -4
            555
                                                                            507 N516JB
##
                         600
                                   -5
                                          913
                                                   854
                                                             19 B6
                                                                                         F.WR.
            557
                                                   723
                                                            -14 EV
                                                                           5708 N829AS
##
    8
                         600
                                   -3
                                          709
                                                                                         LGA
            557
                                          838
                                                   846
                                                             -8 B6
##
                         600
                                   -3
                                                                             79 N593JB
                                                                                         JFK
            558
                                   -2
                                          753
                                                   745
                                                                            301 N3ALAA
## 10
                         600
                                                              8 44
                                                                                        T.GA
##
     ... with 336,766 more rows, 6 more variables: dest <chr>, air time <dbl>,
## #
       distance <dbl>, hour <dbl>, minute <dbl>, time hour <dttm>, and abbreviated
## #
       variable names 1: sched dep time, 2: dep delay, 3: arr time, 4: sched arr time,
       5: arr_delay
## #
```

Ubā

select(): Ranges

This will select all columns within the two variables specified.

```
flights %>% select(year:day)
## # A tibble: 336,776 x 3
##
      year month day
##
      <int> <int> <int>
##
      2013
##
      2013
##
      2013
##
      2013
##
      2013
##
      2013
##
      2013
##
      2013
##
      2013
## 10
      2013
## # ... with 336,766 more rows
```

ubä

select(): Exclusion Ranges

```
flights %>% select(-(year:day))
## #
     A tibble: 336,776 x 16
##
      dep_time sched_dep~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier flight tailnum origin
##
                              <dbl>
         <int>
                      <int>
                                       <int>
                                               <int>
                                                        <dbl> <chr>
                                                                        <int> <chr>
                                                                                      <chr>>
           517
                        515
                                   2
                                         830
                                                 819
                                                           11 UA
                                                                         1545 N14228
                                                                                      EWR
##
           533
                        529
                                   4
                                         850
                                                 830
                                                           20 UA
                                                                         1714 N24211
                                                                                      LGA
##
                                  2
##
           542
                        540
                                         923
                                                 850
                                                           33 AA
                                                                         1141 N619AA
                                                                                      JFK
##
           544
                        545
                                 -1
                                        1004
                                                1022
                                                          -18 B6
                                                                          725 N804JB
                                                                                      JFK
##
           554
                        600
                                 -6
                                         812
                                                 837
                                                          -25 DL
                                                                          461 N668DN
                                                                                      LGA
##
           554
                        558
                                 -4
                                         740
                                                 728
                                                           12 UA
                                                                         1696 N39463
                                                                                      EWR.
##
           555
                        600
                                 -5
                                         913
                                                 854
                                                           19 B6
                                                                          507 N516JB
                                                                                      EWR.
##
           557
                        600
                                 -3
                                         709
                                                 723
                                                          -14 EV
                                                                         5708 N829AS
                                                                                      LGA
##
           557
                        600
                                 -3
                                         838
                                                 846
                                                           -8 B6
                                                                           79 N593JB
                                                                                      JFK
## 10
           558
                        600
                                 -2
                                         753
                                                 745
                                                            AA 8
                                                                          301 N3ALAA
                                                                                      L.G.A
## #
     ... with 336.766 more rows, 6 more variables: dest <chr>, air time <dbl>,
## #
       distance <dbl>, hour <dbl>, minute <dbl>, time hour <dttm>, and abbreviated
## #
       variable names 1: sched_dep_time, 2: dep_delay, 3: arr_time, 4: sched_arr_time,
## #
       5: arr delay
```

Ubō

select(): Matching

```
flights %>% select(contains("dep"),
                    contains("arr"))
## # A tibble: 336,776 x 7
##
      dep time sched dep time dep delay arr time sched arr time arr delay carrier
##
         <int>
                         <int>
                                    <dbl>
                                             <int>
                                                              <int>
                                                                         <dbl> <chr>
##
           517
                           515
                                        2
                                                830
                                                                819
                                                                            11 UA
##
           533
                           529
                                        4
                                                850
                                                                830
                                                                            20 UA
##
           542
                           540
                                                923
                                                                850
                                                                            33 AA
##
           544
                           545
                                       -1
                                               1004
                                                               1022
                                                                           -18 B6
##
           554
                           600
                                       -6
                                                812
                                                                837
                                                                           -25 DL
##
           554
                           558
                                       -4
                                                740
                                                                728
                                                                            12 UA
##
           555
                           600
                                       -5
                                                913
                                                                854
                                                                            19 B6
           557
                           600
                                       -3
                                                709
                                                                723
                                                                           -14 EV
##
##
           557
                           600
                                       -3
                                                838
                                                                846
                                                                            -8 B6
                                                753
## 10
           558
                           600
                                       -2
                                                                745
                                                                             8 AA
     ... with 336,766 more rows
```

Ubä

select(): Matching (cont.)

```
flights %>% select(starts_with("dep"),
                    starts with("arr"))
## # A tibble: 336.776 x 4
##
      dep_time dep_delay arr_time arr_delay
##
         <int>
                    <dbl>
                             <int>
                                        <dbl>
##
           517
                               830
                                           11
##
           533
                               850
                                           20
##
           542
                               923
                                           33
##
           544
                       -1
                               1004
                                          -18
##
           554
                       -6
                               812
                                          -25
##
           554
                       -4
                               740
                                           12
##
           555
                       -5
                               913
                                           19
##
           557
                       -3
                               709
                                          -14
##
           557
                       -3
                               838
                                           -8
## 10
           558
                       -2
                                753
## # ... with 336,766 more rows
```

Other helpers provide by tidyselect: starts_with, ends_with, everything, matches, num_range, one_of, everything, last_col.

Ubö

select() + where(): Get numeric columns

```
flights %>% select(where(is.numeric))
##
       tibble: 336,776 x 14
##
       year month day dep_time sched_~1 dep_d~2 arr_t~3 sched~4 arr_d~5 flight air t~6
##
      <int> <int> <int>
                            <int>
                                     <int>
                                              <dbl>
                                                      <int>
                                                               <int>
                                                                       <dbl>
                                                                              <int>
                                                                                       dbl>
       2013
                              517
                                        515
                                                  2
                                                        830
                                                                                1545
##
                                                                 819
                                                                          11
                                                                                         227
##
       2013
                              533
                                        529
                                                  4
                                                        850
                                                                 830
                                                                          20
                                                                                1714
                                                                                         227
                              542
##
       2013
                                        540
                                                        923
                                                                 850
                                                                          33
                                                                                1141
                                                                                         160
       2013
                              544
                                        545
                                                        1004
                                                                1022
                                                                                 725
##
                                                 -1
                                                                          -18
                                                                                         183
                                        600
                                                        812
##
       2013
                              554
                                                 -6
                                                                 837
                                                                         -25
                                                                                 461
                                                                                         116
       2013
                                        558
                                                        740
                                                                 728
                                                                          12
                                                                                1696
##
                              554
                                                 -4
                                                                                         150
##
       2013
                              555
                                        600
                                                 -5
                                                        913
                                                                 854
                                                                          19
                                                                                 507
                                                                                         158
                                                                                          53
##
       2013
                              557
                                        600
                                                 -3
                                                        709
                                                                 723
                                                                         -14
                                                                                5708
##
    9
       2013
                              557
                                        600
                                                 -3
                                                        838
                                                                 846
                                                                           -8
                                                                                  79
                                                                                         140
## 10
       2013
                              558
                                        600
                                                 -2
                                                         753
                                                                 745
                                                                           8
                                                                                 301
                                                                                         138
## #
     ... with 336,766 more rows, 3 more variables: distance <dbl>, hour <dbl>,
       minute <dbl>, and abbreviated variable names 1: sched dep time, 2: dep delay,
## #
## #
       3: arr time, 4: sched arr time, 5: arr delay, 6: air time
```

Ubd

relocate(): to the front

```
flights %>% relocate(carrier, origin, dest)
## # A tibble: 336,776 x 19
      carrier origin dest year month day dep_time sched_de~1 dep_d~2 arr_t~3 sched~4
##
                      <chr> <int> <int> <int>
##
      <chr>
              <chr>
                                                  <int>
                                                              <int>
                                                                      <dbl>
                                                                               <int>
                                                                                       <int>
    1 UA
              EWR
                      IAH
                             2013
                                                    517
                                                                515
                                                                          2
                                                                                 830
                                                                                         819
##
                                       1
    2 UA
              LGA
                             2013
                                                    533
                                                                529
                                                                                 850
                                                                                         830
##
                      IAH
                                                                          4
              JFK
                                                    542
                                                                540
##
    3 AA
                      MIA
                             2013
                                                                                 923
                                                                                         850
    4 B6
              JFK
                             2013
                                                    544
                                                                545
##
                      BON
                                                                          -1
                                                                                1004
                                                                                        1022
    5 DL
              LGA
                             2013
                                                    554
                                                                600
                                                                                 812
                                                                                         837
##
                      ATL
                                                                          -6
    6 UA
              EWR
                      ORD
                             2013
                                                    554
                                                                558
                                                                                 740
                                                                                         728
##
                                                                         -4
    7 B6
              EWR
                      FLL
                             2013
                                                    555
                                                                                 913
##
                                                                600
                                                                         -5
                                                                                         854
    8 EV
                                                    557
##
              LGA
                      IAD
                             2013
                                                                600
                                                                          -3
                                                                                 709
                                                                                         723
                                                    557
##
    9 B6
              JFK
                      MCO
                             2013
                                                                600
                                                                         -3
                                                                                 838
                                                                                         846
## 10 AA
              LGA
                                                    558
                                                                                 753
                      ORD
                             2013
                                                                600
                                                                          -2
                                                                                         745
     ... with 336,766 more rows, 8 more variables: arr_delay <dbl>, flight <int>,
       tailnum <chr>, air time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
## #
## #
       time_hour <dttm>, and abbreviated variable names 1: sched_dep_time,
## #
       2: dep delay. 3: arr time. 4: sched arr time
```

ubā

relocate(): to the end

```
flights %>%
  relocate(year, month, day, .after = last_col())
## # A tibble: 336,776 x 19
      dep time sched dep~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier flight tailnum origin
##
                              <dbl>
##
         <int>
                      <int>
                                      <int>
                                               <int>
                                                       <dbl> <chr>
                                                                       <int> <chr>
                                                                                      <chr>>
           517
                        515
                                  2
                                         830
                                                 819
                                                           11 UA
                                                                                      EWR
##
                                                                        1545 N14228
##
           533
                        529
                                  4
                                         850
                                                 830
                                                          20 UA
                                                                        1714 N24211
                                                                                      LGA
           542
                                  2
                                                 850
                                                          33 AA
                                                                        1141 N619AA
##
                        540
                                         923
                                                                                      JFK
##
           544
                        545
                                 -1
                                        1004
                                                1022
                                                         -18 B6
                                                                         725 N804JB
                                                                                      JFK
##
           554
                        600
                                 -6
                                         812
                                                 837
                                                         -25 DL
                                                                         461 N668DN
                                                                                      LGA
##
    6
           554
                        558
                                         740
                                                 728
                                                          12 UA
                                                                        1696 N39463
                                                                                      EWR.
                                 -4
           555
                        600
                                         913
                                                 854
                                                          19 B6
                                                                         507 N516JB
                                                                                      EWR.
##
                                 -5
           557
                        600
                                 -3
                                         709
                                                 723
                                                         -14 EV
                                                                        5708 N829AS
                                                                                      LGA
##
##
    9
           557
                        600
                                 -3
                                         838
                                                 846
                                                           -8 B6
                                                                          79 N593JB
                                                                                      JFK
## 10
           558
                        600
                                 -2
                                         753
                                                 745
                                                           8 AA
                                                                         301 N3ALAA
                                                                                      L.G.A
     ... with 336,766 more rows, 9 more variables: dest <chr>, air_time <dbl>,
## #
       distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>. year <int>.
## #
## #
       month <int>, day <int>, and abbreviated variable names 1: sched dep time.
## #
       2: dep_delay, 3: arr_time, 4: sched_arr_time, 5: arr_delay
```

Ubd

rename(): Change column names

```
flights %>% rename(tail_number = tailnum)
## # A tibble: 336,776 x 19
                     day dep_time sched_~1 dep_d~2 arr_t~3 sched~4 arr d~5 carrier flight
##
##
      <int> <int> <int>
                             <int>
                                       <int>
                                               <dbl>
                                                        <int>
                                                                 <int>
                                                                          <dbl> <chr>
                                                                                          <int>
##
       2013
                               517
                                         515
                                                          830
                                                                   819
                                                                             11 UA
                                                                                           1545
##
       2013
                               533
                                         529
                                                           850
                                                                   830
                                                                             20 UA
                                                                                           1714
##
       2013
                               542
                                         540
                                                           923
                                                                   850
                                                                             33 AA
                                                                                            1141
                                         545
                                                                            -18 B6
                                                                                            725
##
       2013
                               544
                                                   -1
                                                          1004
                                                                  1022
       2013
                               554
                                         600
                                                          812
                                                                   837
                                                                            -25 DL
##
                                                   -6
                                                                                            461
                                         558
                                                          740
                                                                   728
                                                                             12 UA
                                                                                            1696
##
       2013
                               554
                                                   -4
                                                                             19 B6
##
       2013
                               555
                                         600
                                                   -5
                                                          913
                                                                   854
                                                                                            507
                               557
                                                                            -14 EV
                                                                                           5708
##
       2013
                                         600
                                                   -3
                                                          709
                                                                   723
                               557
                                                                             -8 B6
##
       2013
                                         600
                                                   -3
                                                          838
                                                                   846
                                                                                              79
##
  10
       2013
                               558
                                         600
                                                   -2
                                                           753
                                                                              8 AA
                                                                   745
                                                                                             301
     ... with 336,766 more rows, 8 more variables: tail number <chr>, origin <chr>,
## #
## #
       dest <chr>, air time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
## #
       time hour <dttm>, and abbreviated variable names 1: sched dep time,
## #
       2: dep_delay, 3: arr_time, 4: sched_arr_time, 5: arr_delay
```

pull

```
names(flights)
## [1] "year"
                         "month"
                                          "day"
                                                           "dep time"
                                                           "sched arr time"
## [5] "sched dep time" "dep delay"
                                          "arr time"
                     "carrier"
                                                           "tailnum"
## [9] "arr delay"
                                          "flight"
## [13] "origin"
                       "dest"
                                          "air time"
                                                           "distance"
## [17] "hour"
                                          "time hour"
                        "minute"
flights %>% pull("year") %>% head()
## [1] 2013 2013 2013 2013 2013 2013
flights %>% pull(1) %>% head() # position from the left
## [1] 2013 2013 2013 2013 2013 2013
flights %>% pull(-1) %>% head() # position from the right
## [1] "2013-01-01 05:00:00 EST" "2013-01-01 05:00:00 EST" "2013-01-01 05:00:00 EST"
## [4] "2013-01-01 05:00:00 EST" "2013-01-01 06:00:00 EST" "2013-01-01 05:00:00 EST"
flights$year %>% head()
## [1] 2013 2013 2013 2013 2013 2013
34 / 50
```

arrange(): Sort data

```
flights %>% filter(month == 3,day == 2) %>% arrange(origin, dest)
    A tibble: 765 x 19
##
       year month day dep_time sched_~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier flight
##
      <int> <int> <int>
                            <int>
                                     <int>
                                             <dbl>
                                                      <int>
                                                              <int>
                                                                      <dbl> <chr>
                                                                                      <int>
       2013
                3
                      2
                             1336
                                      1329
                                                       1426
                                                               1432
                                                                         -6 EV
                                                                                       4263
##
       2013
                3
                      2
                              628
                                       629
                                                                        -12 DL
                                                                                        575
##
                                                 -1
                                                        837
                                                                849
       2013
                3
                              637
                                       640
                                                                        -12 EV
##
                                                 -3
                                                        903
                                                                915
                                                                                       4209
                3
                              743
                                       745
                                                -2
                                                        945
                                                               1010
                                                                        -25 DL
##
       2013
                                                                                        807
       2013
                3
                              857
                                       900
                                                -3
                                                       1117
                                                               1126
                                                                         -9 DL
                                                                                        485
##
##
       2013
                3
                      2
                             1027
                                      1030
                                                -3
                                                       1234
                                                               1247
                                                                        -13 DL
                                                                                       2343
                3
                       2
                                      1145
                                                       1332
                                                               1359
                                                                        -27 DL
##
       2013
                             1134
                                               -11
                                                                                        401
       2013
                3
                      2
                                      1415
                                                       1636
                                                               1630
##
                             1412
                                                 -3
                                                                          6 DL
                                                                                        935
##
       2013
                3
                      2
                             1633
                                      1636
                                                       1848
                                                               1908
                                                                        -20 EV
                                                                                       3273
                                                -3
       2013
                3
                      2
                             1655
                                      1700
                                                -5
                                                       1857
                                                                        -27 DL
## 10
                                                               1924
                                                                                       2042
## #
     ... with 755 more rows, 8 more variables: tailnum <chr>, origin <chr>, dest <chr>,
## #
       air time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time hour <dttm>, and
## #
       abbreviated variable names 1: sched_dep_time, 2: dep_delay, 3: arr_time,
       4: sched_arr_time, 5: arr_delay
## #
```

arrange() & desc(): Descending order

```
flights %>%
  filter(month == 3, day == 2) %>%
  arrange(desc(origin), dest) %>%
  select(origin, dest, tailnum)
## # A tibble: 765 x 3
##
      origin dest tailnum
##
      <chr>
             <chr> <chr>
##
    1 LGA
             ATL
                   N928AT
##
    2 LGA
             ATL
                   N623DL
##
    3 LGA
             ATL
                   N680DA
##
    4 LGA
             ATL
                   N996AT
##
    5 LGA
             ATL
                   N510MQ
```

##

##

##

##

6 LGA

7 LGA

8 LGA

9 LGA

10 LGA

ATL

ATL

ATL

ATL

ATL

... with 755 more rows

N663DN

N942DL

N511MQ

N910DE

N902DE

mutate(): Modify columns

```
flights %>%
  select(year:day) %>%
  mutate(date = paste(year, month, day, sep = "/"))
## # A tibble: 336,776 x 4
##
       year month day date
##
      <int> <int> <int> <chr>
##
       2013
                      1 2013/1/1
##
       2013
                      1 2013/1/1
                      1 2013/1/1
##
       2013
##
       2013
                      1 2013/1/1
##
       2013
                      1 2013/1/1
##
       2013
                      1 2013/1/1
       2013
                      1 2013/1/1
##
    8
       2013
                      1 2013/1/1
##
##
       2013
                      1 2013/1/1
## 10
       2013
                      1 2013/1/1
     ... with 336,766 more rows
```

Ubd

distinct(): Find unique rows

```
flights %>%
  select(origin, dest) %>%
 distinct() %>%
  arrange(origin, dest)
## # A tibble: 224 x 2
##
      origin dest
##
      <chr> <chr>
##
    1 EWR
             ALB
##
    2 EWR
             ANC
##
    3 EWR
             ATL
##
    4 EWR
             AUS
##
    5 EWR
             AVL
```

##

##

##

##

6 EWR

7 EWR.

8 EWR

9 EWR

10 EWR

BDL

BNA

BOS

BQN

BTV ## # ... with 214 more rows

summarise()

##

##

<int>

1 336776

<dbl>

-43

```
flights %>%
  summarize(n(), min(dep_delay), max(dep_delay))
## # A tibble: 1 x 3
   `n()` `min(dep_delay)` `max(dep_delay)`
##
##
      <int>
                    <dbl>
                                        <dbl>
                          NΑ
                                           NA
## 1 336776
flights %>%
  summarize(
   n = n()
   min dep delay = min(dep delay, na.rm = TRUE),
   max_dep_delay = max(dep_delay, na.rm = TRUE)
## # A tibble: 1 x 3
         n min_dep_delay max_dep_delay
```

39 / 50

<dbl>

1301

group_by()

```
flights %>% group_by(origin)
## # A tibble: 336,776 x 19
               origin [3]
##
     Groups:
##
       year month day dep_time sched_~1 dep_d~2 arr_t~3 sched~4 arr_d~5 carrier flight
##
      <int> <int> <int>
                            <int>
                                      <int>
                                              <dbl>
                                                       <int>
                                                                <int>
                                                                        <dbl> <chr>
                                                                                        <int>
       2013
                               517
                                        515
                                                         830
                                                                           11 UA
                                                                                         1545
##
                                                   2
                                                                  819
##
       2013
                               533
                                        529
                                                   4
                                                         850
                                                                  830
                                                                           20 UA
                                                                                         1714
       2013
                               542
                                        540
                                                         923
                                                                           33 AA
##
                                                                  850
                                                                                         1141
##
       2013
                               544
                                        545
                                                  -1
                                                        1004
                                                                 1022
                                                                          -18 B6
                                                                                          725
##
       2013
                               554
                                        600
                                                  -6
                                                         812
                                                                  837
                                                                          -25 DL
                                                                                          461
##
       2013
                               554
                                        558
                                                  -4
                                                         740
                                                                  728
                                                                           12 UA
                                                                                         1696
       2013
                               555
                                        600
                                                  -5
                                                         913
                                                                  854
                                                                           19 B6
                                                                                          507
##
##
       2013
                               557
                                        600
                                                  -3
                                                         709
                                                                  723
                                                                          -14 EV
                                                                                         5708
##
       2013
                               557
                                        600
                                                  -3
                                                         838
                                                                  846
                                                                           -8 B6
                                                                                           79
    9
## 10
       2013
                               558
                                        600
                                                  -2
                                                         753
                                                                  745
                                                                            8 AA
                                                                                          301
## #
     ... with 336,766 more rows, 8 more variables: tailnum <chr>, origin <chr>,
## #
       dest <chr>, air time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
## #
       time hour <dttm>, and abbreviated variable names 1: sched dep time,
## #
       2: dep delay, 3: arr time, 4: sched arr time, 5: arr delay
```

Ubä

summarise() with group_by()

1 EWR.

3 LGA

2 JFK 111279

120835

104662

-25

-43

-33

1126

1301

911

41 / 50

count()

```
flights %>%
                                                  flights %>%
  group by(origin, carrier) %>%
                                                    count(origin, carrier)
  summarise(n = n(), groups = "drop")
                                                  ## # A tibble: 35 x 3
  # or "drop_last" or "keep"
                                                  ##
                                                        origin carrier
## # A tibble: 35 x 3
                                                        <chr> <chr>
                                                  ##
                                                                        <int>
##
                                                      1 EWR
                                                               9E
                                                                         1268
      origin carrier
                        n
                                                  ##
##
      <chr>
             <chr>
                      <int>
                                                  ##
                                                      2 EWR
                                                               AA
                                                                         3487
    1 EWR
             9E
                       1268
                                                      3 EWR
                                                               AS
                                                                         714
##
                                                  ##
##
    2 EWR
             AA
                       3487
                                                      4 EWR
                                                               В6
                                                                         6557
                                                  ##
    3 EWR
             AS
                        714
                                                      5 EWR
                                                                         4342
##
                                                  ##
                                                               DL
    4 EWR
##
             B6
                       6557
                                                  ##
                                                      6 EWR
                                                               EV
                                                                        43939
##
    5 EWR
                       4342
                                                      7 EWR
                                                                         2276
             DL
                                                  ##
                                                               MQ
    6 EWR
##
             EV
                      43939
                                                  ##
                                                      8 EWR
                                                               00
                                                                            6
##
    7 EWR
             MQ
                       2276
                                                      9 EWR
                                                               UA
                                                                        46087
                                                  ##
##
    8 EWR
             00
                          6
                                                  ## 10 EWR
                                                               US
                                                                         4405
##
    9 EWR.
             IJΑ
                      46087
                                                  ## # ... with 25 more rows
## 10 EWR
             US
                       4405
## # ... with 25 more rows
```

OPQ

mutate() with group_by()

```
flights %>% group_by(origin) %>%
  mutate(n = n()) \%
  select(origin, n)
## # A tibble: 336,776 x 2
## #
    Groups:
             origin [3]
##
     origin
              n
      <chr>
              <int>
##
    1 EWR
             120835
##
    2 LGA
             104662
##
##
    3 JFK
             111279
##
    4 JFK
             111279
##
    5 LGA
             104662
##
    6 EWR
             120835
    7 EWR
             120835
##
##
   8 LGA
             104662
##
   9 JFK
             111279
## 10 LGA
             104662
## # ... with 336.766 more rows
```

Tidy data

Tibbles

Pipeline

Data manipulation

Reshaping data (Wide vs. Long)

$\textbf{Wide} \rightarrow \textbf{Long}$

country	1999	2000		country	year	cases
Α	0.7K	2K	\rightarrow	Α	1999	0.7K
В	37K	80K		В	1999	37K
С	212K	213K		С	1999	212K
				Α	2000	2K
				В	2000	80K
				С	2000	213K

pivot_longer (previously gather)

From data import cheatsheet.

44 / 50 UOO

Syntax

```
(d <- tibble::tribble(</pre>
                                               pivot_longer(d, cols = "1999": "2000",
  ~country, ~"1999", ~"2000",
                                                             names_to = "year",
        "A", "0.7K", "2K",
                                                             values to = "cases")
        "B", "37K", "80K",
                                               ## # A tibble: 6 x 3
        "C". "212K". "213K"
                                               ##
                                                    country year cases
))
                                               ##
                                                    <chr>
                                                             <chr> <chr>
## # A tibble: 3 x 3
                                               ## 1 A
                                                             1999 0.7K
##
     country `1999` `2000`
                                               ## 2 A
                                                             2000
                                                                   2K
##
     <chr>
             <chr> <chr>
                                               ## 3 B
                                                             1999
                                                                  37K
## 1 A
             0.7K
                    2K
                                               ## 4 B
                                                             2000
                                                                   80K
## 2 B
             37K
                    80K
                                               ## 5 C
                                                             1999
                                                                  212K
## 3 C
             212K
                    213K
                                               ## 6 C
                                                             2000
                                                                   213K
```



$\textbf{Long} \, \rightarrow \, \textbf{Wide}$

country	year	type	count	
Α	1999	cases	0.7K	_
Α	1999	рор	19M	-
Α	2000	cases	2K	
Α	2000	рор	20M	
В	1999	cases	37K	
В	1999	рор	172M	
В	2000	cases	80K	
В	2000	рор	174M	
С	1999	cases	212K	
С	1999	рор	1T	
С	2000	cases	213K	
С	2000	рор	1T	

country	year	cases	рор
Α	1999	0.7K	19M
Α	2000	2K	20M
В	1999	37K	172M
В	2000	80K	174M
С	1999	212K	1T
С	2000	213K	1T

pivot_wider (previously spread)

Syntax

```
d <- tibble tribble
  ~country, ~year, ~type, ~count,
      "A". 1999. "cases", "0.7K",
      "A", 1999, "pop", "19M",
      "A".
            2000, "cases", "2K",
      "A", 2000, "pop", "20M",
      "B". 1999, "cases", "37K",
      "B".
            1999, "pop", "172M",
      "B".
            2000, "cases", " 80K",
      "B".
            2000, "pop", "174M",
      "C".
            1999, "cases", "212K",
      "C".
            1999, "pop", "1T",
      "C".
            2000. "cases". "213K".
            2000, "pop", "1T"
      "C".
```

```
pivot_wider(d, id_cols = country:year,
            names_from = type,
            values_from = count)
## # A tibble: 6 \times 4
##
     country year cases
                          pop
             <dbl> <chr> <chr>
##
     <chr>
## 1 A
              1999 "0.7K" 19M
## 2 A
              2000 "2K"
                          20M
## 3 B
              1999 "37K" 172M
## 4 B
              2000 " 80K" 174M
## 5 C
              1999 "212K" 1T
## 6 C
              2000 "213K" 1T
```

Ubö

separate()

country	year	rate		country	year	cases	рор
Α	1999	0.7K/19M		Α	1999	0.7K	19M
Α	2000	2K/20M	\rightarrow	Α	2000	2K	20M
В	1999	37K/172M		В	1999	37K	172
В	2000	80K/174M		В	2000	80K	174

```
separate(d, rate, sep = "/", into = c("cases", "pop"))
```

/ 50

unite()

country	century	year		country	year
Α	19	99		Α	19 <mark>99</mark>
Α	20	00	\rightarrow	Α	2000
В	19	99		В	19 <mark>99</mark>
В	20	00		В	2000

```
unite(d, century, year, col = "year", sep = "")
```

9 / 50

Merging two tibbles

Check out the explanation here: https://github.com/gadenbuie/tidyexplain

- left_join(x, y)
- right_join(x, y)
- full_join(x, y)
- inner_join(x, y)
- union(x, y) and union_all(x, y)
- intersect(x, y)
- setdiff(x, y)

