Statistical Computing with R Masters in Data Science 503 (S8&9) Third Batch, SMS, TU, 2024

Shital Bhandary

Associate Professor

Statistics/Bio-statistics, Demography and Public Health Informatics

Patan Academy of Health Sciences, Lalitpur, Nepal

Faculty, Data Analysis and Decision Modeling, MBA, Pokhara University, Nepal

Faculty, FAIMER Fellowship in Health Professions Education, India/USA.

Review Preview (Unit 2, Part 2 & 3)

Data wrangling

Reading database in R

Data munching

• Big data in R

• Tidy data

Text Mining

 "dplyr" package and its use for data manipulation

Data wrangling (Course book Chapter 9-16) R for Data Science, https://r4ds.had.co.nz/index.html

 Data wrangling is the art of getting your data into R in a useful form for visualization and modelling.

• Data wrangling is very important: without it you can't work with your own data! There are three main parts to data wrangling:

- Import
- Tidy
- Transform

Import data in R:

We have already covered this in the previous classes

More here: https://r4ds.had.co.nz/data-import.html

 Reach this chapter well as there are some important import functions that are part of this course and may not have discussed so far

 We will discuss about reading "database" in the second part of this class

Tidy data in R

- Tidy data is a consistent way to organize your data in R.
- Getting your/our data into this format requires some upfront work, but that work pays off in the long term.
- Once you/we have tidy data and the tidy tools provided by packages in the tidyverse, you will spend much less time munging/cleaning data from one representation to another, allowing you to spend more time on the analytic questions at hand.
- Tidy data in tidyverse packages are stored as "tibble"

Let us see what is "tibble" first: https://r4ds.had.co.nz/tibbles.html

- The variant of the data frame used by "tidiverse" is called: **tibble**.
- Tibbles are data frames, but they tweak some older behaviours to make life a little easier.
- R is an old language, and some things that were useful 10 or 20 years ago now get in your way.
- It's difficult to change base R without breaking existing code, so most innovation occurs in packages.
- The tibble package provides opinionated data frames that make working in the tidyverse a little easier.
- It's particularly **useful for large datasets** because it only prints the first few rows.

Note:

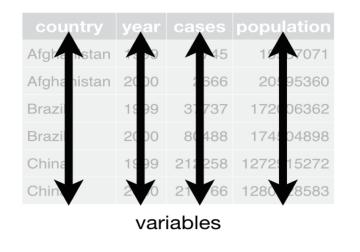
• All the functions of "tidyverse" package works fast with tibble so it will be wise to say that data frame/s should be converted to tibble before using functions of "tidyverse" package

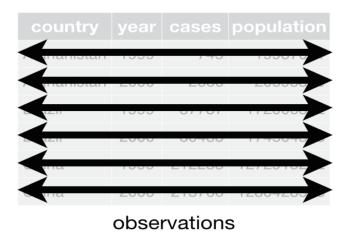
 However, most of the packages of the "tidyverse" super package works well with the data frame too!

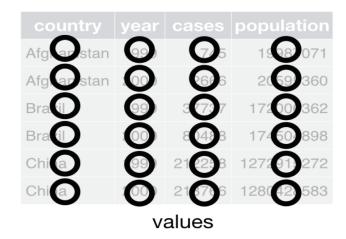
 There are two main differences in the usage of a data frame vs a tibble: printing, and subsetting. https://posit.co/blog/tibble-1-0-0/

There are three interrelated rules which make a dataset tidy:

- Each variable must have its own column.
- Each observation must have its own row.
- Each value must have its own cell.







Example: Which one is "tidy"? Why?

```
#Creating tibble:
table1 <- tibble(
country = c("Afghanistan", "Afghanistan", "Brazil", "Brazil", "China", "China"),
year = c(1999, 2000, 1999, 2000, 1999, 2000),
cases = c(745,2666,37737,80488,212258,213766),
population = c(19987071,20595360,172006362,
174504898, 1272915272,1280428583)
# Data frame to tibble: as_tibble(data_frame)
# Tibble to data frame: as.data.frame(tibble_data)
```

- table1
- A tibble: 6 x 4

•	year country	cases	population
•	<dbl> <chr></chr></dbl>	<dbl></dbl>	<dbl></dbl>
•	1 1999 Afghanistan	745	19987071
•	2 2000 Afghanistan"	2666	20595360
•	3 1999 Brazil	37737	172006362
•	4 2000 Brazil	80488	174504898
•	5 1999 China	212258	1272915272
•	6 2000 China	213766	1280428583

dbl = duble instead of number in "tibble"!

Example: Which one is "tidy"? Why?

• #> # ... with 6 more rows

```
• table2
                                                • table3
• #> # A tibble: 12 × 4
                                                • #> # A tibble: 6 × 3
• #> country
                 year
                                   count
                                                • #> country
                                                                                rate
                         type
                                                                    year
                                                • #> * <chr>
• #> <chr>
                  <int>
                         <chr>
                                    <int>
                                                                    <int>
                                                                               <chr>
• #> 1 Afghanistan 1999
                                     745
                                                • #> 1 Afghanistan 1999
                                                                            745/19987071
                          cases
                                                • #> 2 Afghanistan 2000
• #> 2 Afghanistan 1999 population 19987071
                                                                            2666/20595360
• #> 3 Afghanistan 2000
                                     2666
                                                • #> 3 Brazil
                                                                  1999
                                                                             37737/172006362
                          cases
• #> 4 Afghanistan 2000 population 20595360

    #> 4 Brazil

                                                                   2000
                                                                            80488/174504898
• #> 5 Brazil
                  1999
                                    37737
                                                                   1999
                                                                             212258/1272915272
                                                • #> 5 China
                          cases
                  1999 population 172006362
                                                • #> 6 China
                                                                   2000
                                                                             213766/1280428583

    #> 6 Brazil
```

Example: Which one is "tidy"? Why? # Spread across two tibbles

table4a # cases

- #> # A tibble: 3 × 3
- #> country `1999` `2000`
- #> * <chr> <int> <int>
- #> 1 Afghanistan 745 2666
- #> 2 Brazil 37737 80488
- #> 3 China 212258 213766

table4b # population

- #> # A tibble: 3 × 3
- #> country `1999` `2000`
- #> * <chr> <int> <int>
- #> 1 Afghanistan 19987071 20595360
- #> 2 Brazil 172006362 174504898
- #> 3 China 1272915272 1280428583

Why ensure that your data is tidy? There are two main advantages:

• There's a general advantage to picking one consistent way of storing data. If you have a consistent data structure, it's easier to learn the tools that work with it because they have an underlying uniformity.

• There's a specific advantage to placing variables in columns because it allows R's vectorized nature to shine.

• dplyr, ggplot2, and all the other packages in the tidyverse are designed to work with tidy data.

Tidy data: Pivoting – Longer to wider (To do standard statistical analysis)

- table2
- #> # A tibble: 12 × 4
- #> country year type count
- #> <chr> <int> <chr> <int>
- #> 1 Afghanistan 1999 cases 745
- #> 2 Afghanistan 1999 population 19987071
- #> 3 Afghanistan 2000 cases 2666
- #> 4 Afghanistan 2000 population 20595360
- #> 5 Brazil 1999 cases 37737
- #> 6 Brazil 1999 population 172006362
- #> # ... with 6 more rows

table2 %>%
 pivot_wider(names_from = type, values_from =
count)

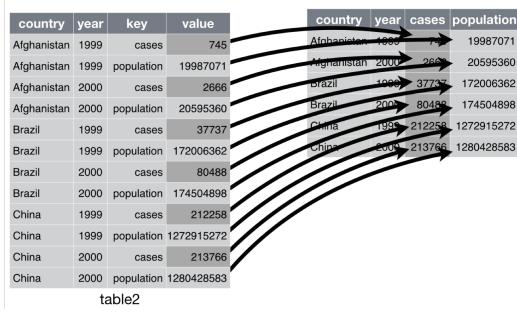


Figure 12.3: Pivoting table 2 into a wider, tidy form.

Tidy data: Pivoting – Wider to Longer (To do Variance components analysis)

- table4a
- #> # A tibble: 3 × 3
- #> country `1999` `2000`
- #> * <chr> <int> <int>
- #> 1 Afghanistan 745 2666
- #> 2 Brazil 37737 80488
- #> 3 China 212258 213766

```
table4a %>%
pivot_longer(c(`1999`, `2000`), names_to =
"year", values_to = "cases")
```

country	year	cases	country	1999	2000
Afghanistan	1999	745	Afghanistan	7/5	2666
Afghanistan	2000	2666	Brazil	37737	80488
Brazil	1999	37737	China	212258	213766
Brazil	2000	80488			
China	1999	212258			
China	2000	213766		table4	

Figure 12.2: Pivoting table4 into a longer, tidy form.

Tidy data: Separate

- table3
- #> # A tibble: 6 × 3
- #> country year rate
- #> * <chr> <int> <chr>
- #> 1 Afghanistan 1999 745/19987071
- #> 2 Afghanistan 2000 2666/20595360
- #> 3 Brazil 1999 37737/172006362
- #> 4 Brazil 2000 80488/174504898
- #> 5 China 1999 212258/1272915272
- #> 6 China 2000 213766/1280428583

table3 %>%
 separate(rate, into = c("cases", "population"))
OR
table3 %>%
 separate(rate, into = c("cases", "population"), sep = "/")

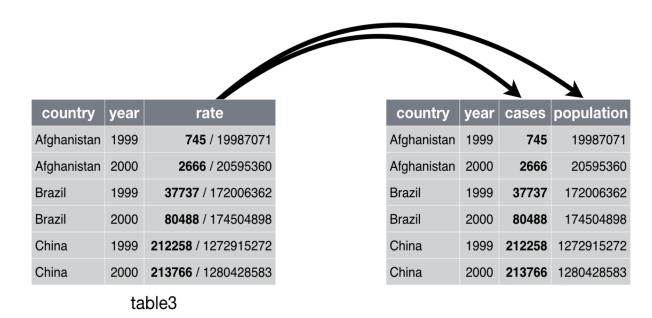


Figure 12.4: Separating table3 makes it tidy

Tidy data: Unite

• unite() is the inverse of separate(): it combines multiple columns into a single column.

 You'll need it much less frequently than **separate()**, but it's still a useful tool to have in your back pocket.

table5 %>% unite(new, century, year) OR table5 %>% unite(new, century, year, sep = "")

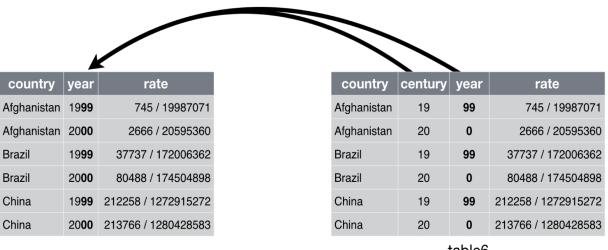


table6

Figure 12.5: Uniting table5 makes it tidy

Missing values

• Changing the representation of a dataset brings up an important subtlety of missing values.

Surprisingly, a value can be missing in one of two possible ways:

- **Explicitly**, i.e. flagged with NA.
- Implicitly, i.e. simply not present in the data.

Missing values: Example

#Create a tibble with missing values:

```
stocks <- tibble(
year = c(2015, 2015, 2015, 2015, 2016, 2016, 2016),
qtr = c( 1,  2,  3,  4,  2,  3,  4),
return = c(1.88, 0.59, 0.35, NA, 0.92, 0.17, 2.66)
)
```

Missing values: Example

There are two missing values in this dataset:

 The return for the fourth quarter of 2015 is explicitly missing, because the cell where its value should be instead contains NA.

• The return for the first quarter of 2016 is **implicitly missing**, because it simply does not appear in the dataset.

Missing values: Example

- stocks %>%
- pivot_wider(names_from = year, values_from = return)
- #> # A tibble: 4 × 3

• #>	qtr	`2015`	`2016`
• #>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
• #> 1	1	1.88	NA
• #> 2	2	0.59	0.92
• #> 3	3	0.35	0.17
• #> 4	4	NA	2.66

Missing values: What will happen now?

```
stocks %>%
      pivot wider(names from = year, values from = return) %>%
      pivot longer(
     cols = c(`2015`, `2016`),
     names to = "year",
     values to = "return",
     values_drop_na = TRUE
```

Missing values: We can use "complete" command!

- stocks %>%
- complete(year, qtr)
- #> # A tibble: 8 × 3
- #> year qtr return
- #> <dbl> <dbl>
- #> 1 2015 1 1.88
- #> 2 2015 2 0.59
- #> 3 2015 3 0.35
- #> 4 2015 4 NA
- #> 5 2016 1 NA
- #> 6 2016 2 0.92
- #> # ... with 2 more rows

Missing values: Another example (tibble by row or tribble!)

```
treatment <- tribble(</li>
     ~ person, ~ treatment, ~response,
     "Derrick Whitmore", 1, 7,
     NA, 2, 10,
     NA,
     "Katherine Burke", 1,
treatment
```

Missing values: fill() for another example

```
• treatment %>%
                        # "tidyr" package is required here!
• fill(person)
• #> # A tibble: 4 × 3
• #> person
                        treatment
                                           response
                         <dbl>
                                           <dbl>
• #> <chr>
• #> 1 Derrick Whitmore
• #> 2 Derrick Whitmore
                                            10
• #> 3 Derrick Whitmore
                                           9
• #> 4 Katherine Burke
```

Question/Queries so far?

Transform/manipulate data with "dplyr"

- To learn five key "dplyr" package functions that allow you to solve the vast majority of your data manipulation challenges:
 - Pick observations by their values (filter()).
 - Reorder the rows (arrange()).
 - Pick variables by their names (select()).
 - Create new variables with functions of existing variables (mutate()).
 - Collapse many values down to a single summary (summarise()).
- These can all be used in conjunction with **group_by()** which changes the scope of each function from operating on the entire dataset to operating on it group-by-group.

Data manipulation with "dplyr"

- These six functions provide the verbs for a language of data manipulation.
- All verbs work similarly:
 - The first argument is a data frame.
 - The subsequent arguments describe what to do with the data frame, using the variable names (without quotes).
 - The result is a new data frame.
- Together these properties make it easy to chain together multiple simple steps to achieve a complex result.

Let's use them with nycflighst13 data

- library(dplyr)
- library(nycflights13)
- flights
- #> # A tibble: 336,776 × 19
- #> year month day dep_time sched_dep...¹ dep_d...² arr_t...³ sched...⁴ arr_d...⁵ carrier

```
• #> <int> <int> <int> <int> <dbl> <int> <dbl> <int> <dbl> <
```

- #> 5 2013 1 1 554 600 -6 812 837 -25 DL
- #> 6 2013 1 1 554 558 -4 740 728 12 UA
- #> # ... with 336,770 more rows, 9 more variables

Filter: What will happen?

- filter(flights, month == 1, day == 1)
- #> # A tibble: 842 × 19
- #> year month day dep_time sched_dep...¹ dep_d...² arr_t...³ sched...⁴ arr d...⁵ carrier

```
<int> <dbl> <int> <dbl> <chr>
• #> <int> <int> <int>
                                 830
                                      819
• #> 1 2013
           1
                 517
                        515
                              2
                                            11 UA
• #> 2 2013 1 1
                 533
                        529
                                 850
                                      830
                                           20 UA
                              4
• #> 3 2013 1 1
                542
                        540 2
                                 923
                                      850
                                           33 AA
• #> 4 2013 1 1
                544
                        545
                                 1004 1022 -18 B6
                              -1
• #> 5 2013
           1 1 554
                        600
                              -6 812
                                      837
                                           -25 DL
           1
                                      728
• #> 6 2013
                 554
                        558
                              -4 740
                                            12 UA
```

• #> # ... with 836 more rows, 9 more variables

Are these better?

- jan1 <- filter(flights, month == 1, day == 1)
- (jan1 <- filter(flights, month == 1, day == 1))

- dec25 <- filter(flights, month == 12, day == 25)
- (dec25 <- filter(flights, month == 12, day == 25))

- filter(flights, month = 1) #Why error?
- filter(flights, month == 1) #Works now? Why?

More with filter:

filter(flights, month == 11 | month == 12) #What?
filter(flights, month == 11 | 12) #Works?
nov dec <- filter(flights, month %in% c(11, 12)) #Works?

- De Morgan's Law:
- filter(flights, !(arr_delay > 120 | dep_delay > 120)) #Works?
- filter(flights, arr_delay <= 120, dep_delay <= 120) #Works?

Arrange: Example

- arrange(flights, year, month, day)
- #> # A tibble: 336,776 × 19
- #> year month day dep_time sched_dep...¹ dep_d...² arr_t...³ sched...⁴ arr_d...⁵ carrier

```
<int> <dbl> <int> <dbl> <chr>
• #> <int> <int> <int>
                                    819
• #> 1 2013
          1
             1
                 517
                       515
                             2 830
                                          11 UA
• #> 2 2013 1 1
                 533
                       529
                             4 850
                                    830
                                         20 UA
• #> 3 2013 1 1 542
                       540 2 923
                                    850 33 AA
• #> 4 2013 1 1 544
                       545
                            -1 1004
                                    1022 -18 B6
• #> 5 2013 1 1 554
                       600
                            -6 812
                                     837
                                         -25 DL
• #> 6 2013
             1
                 554
                       558
                                740
                                    728
                            -4
                                          12 UA
```

• #> # ... with 336,770 more rows, 9 more variables

What will happen now?

Arrange will sort the data in ascending order

arrange(flights, desc(dep_delay))

Use desc() to re-order by a column in descending order

Missing values are always sorted at the end

Select: Example

- # Select columns by name
- select(flights, year, month, day)
- #> # A tibble: 336,776 × 3
- #> year month day
- #> <int> <int>
- #> 1 2013 1 1
- #> 2 2013 1 1
- #> 3 2013 1 1
- #> 4 2013 1 1
- #> 5 2013 1 1
- #> 6 2013 1 1
- #> # ... with 336,770 more rows

- # Select all columns between year and day (inclusive)
- select(flights, year:day)
- #> # A tibble: 336,776 × 3
- #> year month day
- #> <int> <int>
- #> 1 2013 1 1
- #> 2 2013 1 1
- #> 3 2013 1 1
- #> 4 2013 1 1
- #> 5 2013 1 1
- #> 6 2013 1 1
- #> # ... with 336,770 more rows

Select: "except" example

- # Select all columns except those from year to day (inclusive)
- select(flights, -(year:day))
- #> # A tibble: 336,776 × 16
- #> dep_time sched...¹ dep_d...² arr_t...³ sched...⁴ arr_d...⁵ carrier flight tailnum origin
- #> <int> <int> <dbl> <int> <dbl> <chr> <int> <chr>
- #> 1 517 515 2 830 819 11 UA 1545 N14228 EWR
- #> 2 533 529 4 850 830 20 UA 1714 N24211 LGA
- #> 3 542 540 2 923 850 33 AA 1141 N619AA JFK
- #> 4 544 545 -1 1004 1022 -18 B6 725 N804JB JFK
- #> 5 554 600 -6 812 837 -25 DL 461 N668DN LGA
- #> 6 554 558 -4 740 728 12 UA 1696 N39463 EWR
- #> # ... with 336,770 more rows, 6 more variables

Select: More

- There are a number of helper functions you can use within select():
- starts_with("abc"): matches names that begin with "abc".
- ends_with("xyz"): matches names that end with "xyz".
- contains("ijk"): matches names that contain "ijk".

- matches("(.)\\1"): selects variables that match a **regular expression**.
- This one matches any variables that contain repeated characters.
- num_range("x", 1:3): matches x1, x2 and x3.
- See ?select for more details.

More on regular expression are available here: https://cran.r-project.org/web/packages/stringr/vignettes/regular-expressions.html

Note:

- select() can be used to rename variables, but it's rarely useful because it drops all of the variables not explicitly mentioned.
- Instead, use rename(), which is a variant of select() that keeps all the variables that aren't explicitly mentioned
- rename(flights, tail_num = tailnum)

- Another option is to use select() in conjunction with the everything() helper.
- This is useful if you have a handful of variables you'd like to move to the start of the data frame.
- select(flights, time_hour, air_time, everything())

Mutate: Example

- Besides selecting sets of existing columns, it's often useful to add new columns that are functions of existing columns.
- That's the job of mutate().
- mutate() always adds new columns at the end of your dataset so we'll start by creating a narrower dataset so we can see the new variables.

```
#Addiing variables in flights sml:
flights sml <- select(flights,
       year:day,
       ends with("delay"),
       distance,
       air time
mutate(flights sml,
      gain = dep delay - arr delay,
   speed = distance / air_time * 60
```

Mutate: Example

- Besides selecting sets of existing columns, it's often useful to add new columns that are functions of existing columns.
- That's the job of mutate().
- mutate() always adds new columns at the end of your dataset so we'll start by creating a narrower dataset so we can see the new variables.

#Adding one more variable:

```
mutate(flights_sml,
  gain = dep_delay - arr_delay,
  hours = air_time / 60,
  gain_per_hour = gain / hours
)
```

#Note that you/we can refer to columns that you've just created

Transmute and other useful creation functions More@ https://r4ds.had.co.nz/transform.html

• If you only want to keep the new variables, use transmute()

```
transmute(flights,
gain = dep_delay - arr_delay,
hours = air_time / 60,
gain_per_hour = gain / hours
)
```

- Arithmetic operators: +, -, *, /, ^
- Modular arithmetic: %/% (integer division) and %% (remainder)
- Use: Compute hour and minute from dep_time with:
- transmute(flights,
 dep_time,
 hour = dep_time %/% 100,
 minute = dep_time %% 100)

Summarise: Works best for group summaries

- summarise(flights, delay = mean(dep_delay, na.rm = TRUE))
- #> # A tibble: 1 × 1
- #> delay
- #> <dbl>
- #> 1 12.6

- by_day <- group_by(flights, year, month, day)
- summarise(by_day, delay = mean(dep_delay, na.rm = TRUE))
- #> # A tibble: 365 × 4
- #> # Groups: year, month [12]
- #> year month day delay
- #> <int> <int> <dbl>
- #> 1 2013 1 1 11.5
- #> 2 2013 1 2 13.9
- #> 3 2013 1 3 11.0
- #> 4 2013 1 4 8.95
- #> 5 2013 1 5 5.73
- #> 6 2013 1 6 7.15
- #> # ... with 359 more rows

Multiple operations: pipes

```
#What will happen?
#What will happen?
delays <- flights %>%
                                     flights %>%
 group by(dest) %>%
                                      group by(year, month, day) %>%
 summarise(
                                       summarise(mean =
                                     mean(dep delay))
  count = n(),
  dist = mean(distance, na.rm =
TRUE),
                                     #And now?
  delay = mean(arr_delay, na.rm =
                                     flights %>%
TRUE)
                                      group by(year, month, day) %>%
 ) %>%
                                       summarise(mean =
                                     mean(dep delay, na.rm = TRUE))
 filter(count > 20, dest != "HNL")
```

How to remove cancelled flights? And, get summaries by groups!

```
filter(!is.na(dep_delay),
!is.na(arr_delay))

not_cancelled %>%
group_by(year, month, day) %>%
summarise(mean = mean(dep_delay))
```

not cancelled <- flights %>%

```
• #> # A tibble: 365 × 4
• #> # Groups: year, month [12]
• #> year month day mean
• #> <int> <int> <dbl>
• #> 1 2013 1 111.4
• #> 2 2013 1 2 13.7
• #> 3 2013 1 3 10.9
• #> 4 2013 1 4 8.97
• #> 5 2013 1
               5 5.73
• #> 6 2013 1 6 7.15
```

• #> # ... with 359 more rows

Counts: Example

 Whenever you do any aggregation, it's always a good idea to include either a count (n()), or a count of non-missing values (sum(!is.na(x))).

 That way you can check that you're not drawing conclusions based on very small amounts of data.

```
# What happens now?
delays <- not cancelled %>%
     group by(tailnum) %>%
     summarise(
     delay = mean(arr delay)
hist(delays$delay)
```

What happens now?

```
delays <- not cancelled %>%
                                    # Plots
      group by(tailnum) %>%
                                    # Can you interpret them?
      summarise(
      delay = mean(arr delay,
                                    hist(delays$n)
na.rm = TRUE),
      n = n()
                                    hist(delays$delay)
                                    plot(delays$n, delays$delay)
```

Useful summary functions: https://r4ds.had.co.nz/transform.html

```
# When do the first and last flights
leave each day?
not cancelled %>%
 group by(year, month, day) %>%
 summarise(
  first = min(dep time),
  last = max(dep time)
```

• # Why is distance to some destinations more variable than to others? not cancelled %>% group by(dest) %>% summarise(distance sd = sd(distance)) %>% arrange(desc(distance_sd))

Useful summary functions: https://r4ds.had.co.nz/transform.html

```
# Which destinations have the
most carriers?
not_cancelled %>%
  group_by(dest) %>%
  summarise(carriers =
n_distinct(carrier)) %>%
  arrange(desc(carriers))
```

 # How many flights left before 5am? (these usually indicate delayed flights from the previous day)
 not_cancelled %>%
 group_by(year, month, day) %>%
 summarise(n_early = sum(dep_time < 500))

Useful summary functions: https://r4ds.had.co.nz/transform.html

```
# What proportion of flights are delayed by more than an hour?
```

```
#Find all groups bigger than a threshold:
```

```
not_cancelled %>% popular_dests <- flights %>% group_by(year, month, day) %>% group_by(dest) %>% summarise(hour_prop = filter(n() > 365) mean(arr_delay > 60)) popular_dests
```

Popular destination: head and tail (Are these results VALID?)

- head(popular_dests\$dest)
- [1] "IAH" "IAH" "MIA" "BQN" "ATL" "ORD"

- IAH = Texas
- MIA = Miami
- BQN = Puerto Rico
- ATL = Atalanta
- ORD = Chichago

- tail(popular_dests\$dest)
- [1] "BNA" "DCA" "SYR" "BNA" "CLE" "RDU"
- BNA = Nashville
- DCA = Washigton (Reagan Nat.)
- SYR = New York (Syracuse)
- CLE = Cleveland
- RDU = North Carolina

Bonus: dplyr "slice" function with examples https://dplyr.tidyverse.org/reference/slice.html

```
#What will happen?
```

flights %>% slice(1L)

flights %>% slice(n())

flights %>% slice(5:n())

slice(flights,-(1:4))

- flights %>% slice_sample(n=5)
- flights %>% slice_sample(n=5, replace = TRUE)
- set seed(123)
- train_data <- flights %>% slice_sample(prop=0.8)
- train_data
- test_data <- flights %>% slice_sample(prop=0.2)
- test_data

Question/Queries so far?

Readings for the Next class

• R and Relational database: Chapter 13, R for Data Science, First Edition https://r4ds.had.co.nz/relational-data.html

• R for Data Science, 2nd Edition, Chapter 22: Databases https://r4ds.hadley.nz/databases

• R and Big Data: https://rviews.rstudio.com/2019/07/17/3-big-data-strategies-for-r/

Reading continued ... Big Data in R:

https://www.columbia.edu/~sjm2186/EPIC R/EPIC R BigData.pdf

- Option I: Make the data smaller
 - Issue SQL query directly from R to database
- Option II: Get a bigger computer
- Option III: Use data.table rather than data.frame
- Option IV: Buffer the data set on disk
- Option V: Split it up

Question/Queries?

Thank you!

@shitalbhandary