

# Week 2: Tidying and Wrangling data using R

## Getting started

This week we will demonstrate various techniques for **tidying** and **wrangling** data in R. From the ‘Introduction to R Programming’ course we are familiar with a data frame in R: a rectangular spreadsheet-like representation of data in R where the rows correspond to observations and the columns correspond to variables describing each observation. In Week 1 of Data Analysis, we started exploring the data frame `flights` included in the `nycflights13` package by creating visualisations of the data contained within said data frame.

Here we will discover a type of data formatting called **tidy** data. You will see that having data stored in the **tidy** format is about more than what the colloquial definition of the term **tidy** might suggest of having your data “neatly organised” in a spreadsheet. Instead, we define the term **tidy** in a more rigorous fashion, outlining a set of rules by which data can be stored and the implications of these rules on analyses.

### Note

This session is based on Chapters 4 and 5 of the open-source book [An Introduction to Statistical and Data Science via R](#) which can be consulted at any point.

First, start by opening **RStudio** by going to **Desktop -> Maths-Stats -> RStudio**. Once RStudio has opened create a new R script by going to **File -> New File -> R Script**. Next go to **File -> Save As...** and save the script into your personal drive, either M: or K: (do not save it to the H: drive). We shall now load into R all of the libraries we will need for this session. This can be done by typing the following into your R script:

```
library(tidyverse)
library(nycflights13)
library(fivethirtyeight)
```

The **tidyverse** library is actually a collection of different R packages for transforming and visualising data. The final two libraries (**nycflights13** and **fivethirtyeight**) contain interesting data sets that we shall examine in this session. Notice that when loading the **tidyverse** package you get a message that tells you about conflicting functions of certain packages. This means that there is at least one or more functions with the same name loaded from different packages (and thus one the function will mask the other). You can use the function **tidyverse\_conflicts()** for getting a list of the conflicted packages:

```
tidyverse_conflicts()
```

```
-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()    masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

In here, we can see for example that the **filter** function from the **dplyr** package has a conflict with the **filter** function in base R **stats** library. A way of sorting that out is to load the **dplyr** library after base R so that R will only consider the version of the function that was last loaded. We can be more rigorous about this and load the **conflicted** library. This will prohibit us to use any functions that have some conflict with previously defined functions.

```
library(conflicted)
```

By doing this, we would need to be more specific about the source package from which the desired function should be loaded. There are two ways of doing this:

1. Using **::** after calling the package name every time we use the function from that package. E.g., **dplyr::filter(...)** will tell R to explicitly use the function **filter** from the **dplyr** library.
2. Using the **conflicts\_prefer("function","package")** function to explicitly declare which version of the function you want to use in the remaining R session (i.e. after **conflicts\_prefer()** is called, e.g., **conflict\_prefer("filter","dplyr")** .

#### Question

What do you think is the advantage of using the **conflicts\_prefer** as opposed to the first approach?

## What is tidy data?

What does it mean for your data to be **tidy**? Beyond just being organised, having **tidy** data means that your data follows a standardised format. This makes it easier for you and others to visualise your data, to wrangle/transform your data, and to model your data. We will follow Hadley Wickham's definition of **tidy data** here:

A data set is a collection of values, usually either numbers (if quantitative) or strings AKA text data (if qualitative). Values are organised in two ways. Every value belongs to a variable and an observation. A variable contains all values that measure the same underlying attribute (like height, temperature, duration) across units. An observation contains all values measured on the same unit (like a person, or a day, or a city) across attributes.

Tidy data is a standard way of mapping the meaning of a data set to its structure. A data set is messy or tidy depending on how rows, columns and tables are matched up with observations, variables and types. In tidy data:

1. Each variable forms a column.
2. Each observation forms a row.
3. Each type of observational unit forms a table.

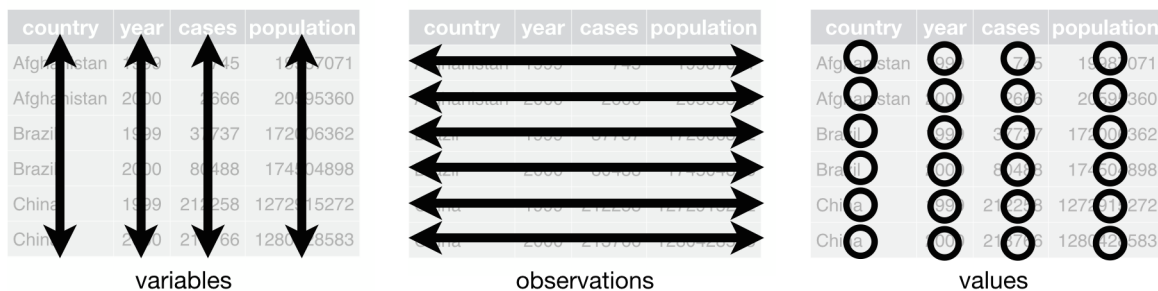


Figure 1: Tidy data graphic from <http://r4ds.had.co.nz/tidy-data.html>

For example, say the following table consists of stock prices:

Table 1: Stock Prices (Non-Tidy Format)

Date	Boeing Stock Price	Amazon Stock Price	Google Stock Price
2009-01-01	\$173.55	\$174.90	\$174.34
2009-01-02	\$172.61	\$171.42	\$170.04

Although the data are neatly organised in a spreadsheet-type format, they are not in tidy format since there are three variables corresponding to three unique pieces of information (Date, Stock Name, and Stock Price), but there are not three columns. In tidy data format each variable should be its own column, as shown below. Notice that both tables present the same information, but in different formats.

Table 2: Stock Prices (Tidy Format)

Date	Stock Name	Stock Price
2009-01-01	Boeing	\$173.55
2009-01-02	Boeing	\$172.61
2009-01-01	Amazon	\$174.90
2009-01-02	Amazon	\$171.42
2009-01-01	Google	\$174.34
2009-01-02	Google	\$170.04

However, consider the following table:

Table 3: Date, Boeing Price, Weather Data

Date	Boeing Price	Weather
2009-01-01	\$173.55	Sunny
2009-01-02	\$172.61	Overcast

In this case, even though the variable **Boeing Price** occurs again, the data *is* tidy since there are three variables corresponding to three unique pieces of information (Date, Boeing stock price, and the weather on that particular day).

The non-tidy data format in the original table is also known as **wide** format whereas the tidy data format in the second table is also known as **long/narrow** data format. In this course, we will work mostly with data sets that are already in the tidy format.

### Question

Consider the following data frame of average number of servings of beer, spirits, and wine consumption in three countries as reported in the FiveThirtyEight article [Dear Mona Followup: Where Do People Drink The Most Beer, Wine And Spirits?](#)

```
# A tibble: 3 x 4
  country      beer_servings spirit_servings wine_servings
  <chr>          <int>          <int>          <int>
1 Canada           240            122            100
2 South Korea       140             16             9
3 USA              249            158            84
```

This data frame is not in tidy format. What would it look like if it were?

I need a hint

Think of these data as being in a wide format. What variables in this data set could be placed in different columns?

See the solution

```
# A tibble: 9 x 3
  country    `beverages type` `number of servings`
  <chr>      <chr>                <int>
1 Canada    beer_servings           240
2 South Korea beer_servings           140
3 USA       beer_servings           249
4 Canada    spirit_servings          122
5 South Korea spirit_servings          16
6 USA       spirit_servings          158
7 Canada    wine_servings            100
8 South Korea wine_servings           9
9 USA       wine_servings            84
```

## Observational units

Recall the `nycflights13` package with data about all domestic flights departing from New York City in 2013 that we used in Week 1 to create visualisations. In particular, let's revisit the `flights` data frame:

```
dim(flights) # Returns the dimensions of a data frame (number obs. and variables)
```

```
[1] 336776    19
```

```
head(flights) # Returns the first 6 rows of the object
```

```
# A tibble: 6 x 19
  year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
  <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
1  2013     1     1     517           515           2     830           819
2  2013     1     1     533           529           4     850           830
```

```

3 2013      1      1      542      540      2      923      850
4 2013      1      1      544      545     -1     1004     1022
5 2013      1      1      554      600     -6      812      837
6 2013      1      1      554      558     -4      740      728
# i 11 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
#   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
#   hour <dbl>, minute <dbl>, time_hour <dtm>

```

```

glimpse(flights) # Lists the variables in an object with their first few values

```

```

Rows: 336,776

```

```

Columns: 19

```

```

$ year      <int> 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2~
$ month     <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1~
$ day       <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1~
$ dep_time  <int> 517, 533, 542, 544, 554, 554, 555, 557, 557, 558, 558, ~
$ sched_dep_time <int> 515, 529, 540, 545, 600, 558, 600, 600, 600, 600, 600, ~
$ dep_delay <dbl> 2, 4, 2, -1, -6, -4, -5, -3, -3, -2, -2, -2, -2, -2, -1~
$ arr_time  <int> 830, 850, 923, 1004, 812, 740, 913, 709, 838, 753, 849,~
$ sched_arr_time <int> 819, 830, 850, 1022, 837, 728, 854, 723, 846, 745, 851,~
$ arr_delay <dbl> 11, 20, 33, -18, -25, 12, 19, -14, -8, 8, -2, -3, 7, -1~
$ carrier   <chr> "UA", "UA", "AA", "B6", "DL", "UA", "B6", "EV", "B6", "~
$ flight    <int> 1545, 1714, 1141, 725, 461, 1696, 507, 5708, 79, 301, 4~
$ tailnum   <chr> "N14228", "N24211", "N619AA", "N804JB", "N668DN", "N394~
$ origin    <chr> "EWR", "LGA", "JFK", "JFK", "LGA", "EWR", "EWR", "LGA",~
$ dest      <chr> "IAH", "IAH", "MIA", "BQN", "ATL", "ORD", "FLL", "IAD",~
$ air_time  <dbl> 227, 227, 160, 183, 116, 150, 158, 53, 140, 138, 149, 1~
$ distance  <dbl> 1400, 1416, 1089, 1576, 762, 719, 1065, 229, 944, 733, ~
$ hour      <dbl> 5, 5, 5, 5, 6, 5, 6, 6, 6, 6, 6, 6, 6, 6, 6, 5, 6, 6, 6~
$ minute    <dbl> 15, 29, 40, 45, 0, 58, 0, 0, 0, 0, 0, 0, 0, 0, 0, 59, 0~
$ time_hour <dtm> 2013-01-01 05:00:00, 2013-01-01 05:00:00, 2013-01-01 0~

```

We see that `flights` has a rectangular shape with each row corresponding to a different flight and each column corresponding to a characteristic of that flight. This matches exactly with the first two properties of tidy data, namely:

1. Each variable forms a column.
2. Each observation forms a row.

But what about the third property?

3. Each type of observational unit forms a table.

The observational unit in the `flights` data set is an individual flight and we can see above that this data set consists of 336,776 flights with 19 variables. In other words, rows of this data set don't refer to a measurement on an airline or on an airport; they refer to characteristics/measurements on a given flight from New York City in 2013. This illustrates the third property of tidy data, i.e. each observational unit is fully described by a single data set.

Note that there is only one observational unit of interest in any analysis. For example, also included in the `nycflights13` package are data sets with different observational units:

- `airlines`
- `planes`
- `weather`
- `airports`

The organisation of this data follows the third **tidy** data property: observations corresponding to the same observational unit are saved in the same data frame.

#### Task

For each of the data sets listed above (other than `flights`), identify the observational unit and how many of these are described in each of the data sets.

- In the `airlines` data set the observational unit is
- (A) Type of plane
- (B) Flight number
- (C) airport code
- (D) IATA carrier codes and names

and there are \_\_\_\_ observational units. - In the `planes` data set the observational unit is

- (A) Flight
- (B) Manufacturer of the plane
- (C) Plane
- (D) Average cruising speed in mph

and there are \_\_\_\_\_ observational units. - In the **weather** data set the observational unit is the

- (A) Weather Station
- (B) Temperature
- (C) Relative humidity
- (D) Sea level pressure

and there are \_\_\_\_\_ observations on **average** across the three NYC airports. - In the **airports** data set the observational unit is the

- (A) Time zone
- (B) Airport
- (C) Altitude
- (D) Daylight savings time zone

and there are \_\_\_\_\_ observational units.

## Identification vs measurement variables

There is a subtle difference between the kinds of variables that you will encounter in data frames: **measurement variables** and **identification variables**. The **airports** data frame contains both these types of variables. Recall that in **airports** the observational unit is an airport, and thus each row corresponds to one particular airport. Let's pull them apart using the **glimpse** function:

```
glimpse(airports)
```

```
Rows: 1,458  
Columns: 8
```



```

$ faa    <chr> "04G", "06A", "06C", "06N", "09J", "0A9", "0G6", "0G7", "0P2", "~
$ name   <chr> "Lansdowne Airport", "Moton Field Municipal Airport", "Schaumbur~
$ lat    <dbl> 41.13047, 32.46057, 41.98934, 41.43191, 31.07447, 36.37122, 41.4~
$ lon    <dbl> -80.61958, -85.68003, -88.10124, -74.39156, -81.42778, -82.17342~
$ alt    <dbl> 1044, 264, 801, 523, 11, 1593, 730, 492, 1000, 108, 409, 875, 10~
$ tz     <dbl> -5, -6, -6, -5, -5, -5, -5, -5, -5, -8, -5, -6, -5, -5, -5, ~
$ dst    <chr> "A", "A", "A", "A", "A", "A", "A", "A", "U", "A", "A", "U", "A", ~
$ tzone  <chr> "America/New_York", "America/Chicago", "America/Chicago", "Ameri~

```

The variables `faa` and `name` are what we will call **identification variables**: variables that uniquely identify each observational unit. They are mainly used to provide a unique name to each observational unit, thereby allowing us to uniquely identify them. `faa` gives the unique code provided by the Federal Aviation Administration in the USA for that airport, while the `name` variable gives the longer more natural name of the airport. The remaining variables (`lat`, `lon`, `alt`, `tz`, `dst`, `tzone`) are often called **measurement** or **characteristic** variables: variables that describe properties of each observational unit, in other words each observation in each row. For example, `lat` and `lon` describe the latitude and longitude of each airport.

Furthermore, sometimes a single variable might not be enough to uniquely identify each observational unit: combinations of variables might be needed (see **Task** below). While it is not an absolute rule, for organisational purposes it is considered good practice to have your identification variables in the far left-most columns of your data frame.

#### Question

What properties of the observational unit do each of `lat`, `lon`, `alt`, `tz`, `dst`, and `tzone` describe for the `airports` data frame?

- (A) Carriers in each airport
- (B) Airport Flights
- (C) Airport appliances
- (D) Spatial location of the airport

#### Task

From the data sets listed above, find an example where combinations of variables are needed to uniquely identify each observational unit.

**Hint** think about the weather data set, can you identify each observational unit based

on the station id only?

## Importing spreadsheets into R

Up to this point, we have been using data stored inside of an R package. In the real world, your data will usually come from a spreadsheet file either on your computer or online. Spreadsheet data is often saved in one of two formats:

- A **Comma Separated Values** `.csv` file. You can think of a CSV file as a bare-bones spreadsheet where:
  - Each line in the file corresponds to one row of data/one observation.
  - Values for each line are separated with commas. In other words, the values of different variables are separated by commas.
  - The first line is often, but not always, a *header* row indicating the names of the columns/variables.
- An **Excel** `.xlsx` file. This format is based on Microsoft's proprietary Excel software. As opposed to bare-bones `.csv` files, `.xlsx` Excel files contain a lot of *metadata*, i.e. data about the data. Examples include the use of bold and italic fonts, colored cells, different column widths, and formula macros etc.

We'll cover two methods for importing data in R: one using the R console and the other using RStudio's graphical interface.

### Method 1: From the console

First, let's download a **Comma Separated Values** (CSV) file of ratings of the level of democracy in different countries spanning 1952 to 1992: [https://moderndive.com/data/dem\\_score.csv](https://moderndive.com/data/dem_score.csv). We use the `read_csv()` function from the `readr` package to read it off the web:

```
dem_score <- read_csv("https://moderndive.com/data/dem_score.csv")
```

```
# A tibble: 96 x 10
```

	country	`1952`	`1957`	`1962`	`1967`	`1972`	`1977`	`1982`	`1987`	`1992`
	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	Albania	-9	-9	-9	-9	-9	-9	-9	-9	5
2	Argentina	-9	-1	-1	-9	-9	-9	-8	8	7
3	Armenia	-9	-7	-7	-7	-7	-7	-7	-7	7
4	Australia	10	10	10	10	10	10	10	10	10

```

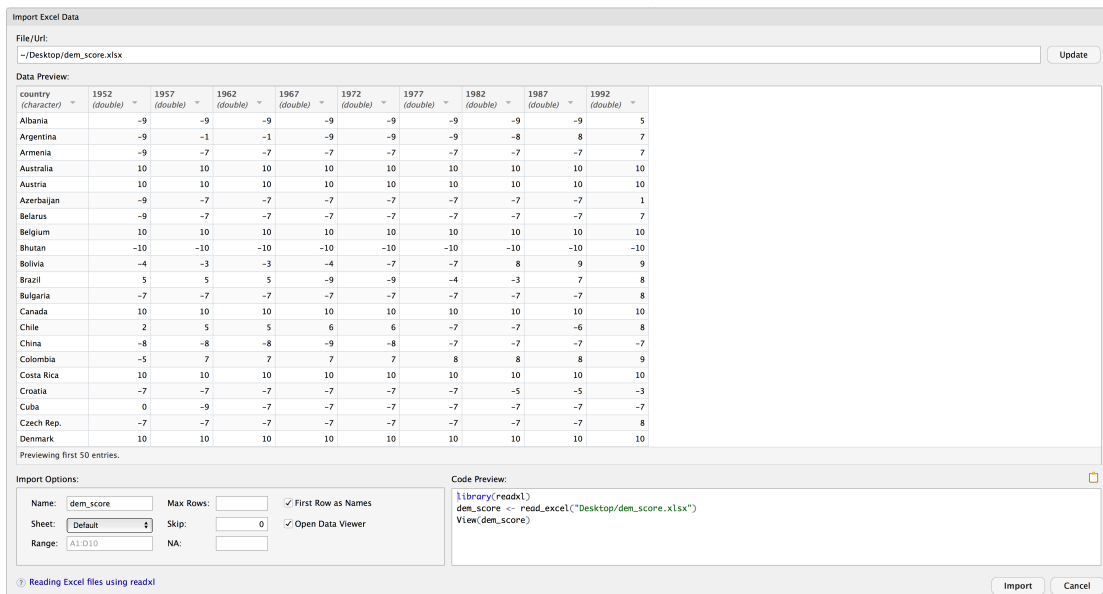
5 Austria      10      10      10      10      10      10      10      10      10
6 Azerbaijan  -9       -7      -7      -7      -7      -7      -7      -7      1
7 Belarus     -9       -7      -7      -7      -7      -7      -7      -7      7
8 Belgium     10       10      10      10      10      10      10      10      10
9 Bhutan      -10      -10     -10     -10     -10     -10     -10     -10     -10
10 Bolivia     -4       -3      -3      -4      -7      -7       8       9       9
# i 86 more rows

```

In this `dem_score` data frame, the minimum value of -10 corresponds to a highly autocratic nation whereas a value of 10 corresponds to a highly democratic nation.

## Method 2: Using RStudio's interface

Let's read in the same data saved in Excel format this time at [https://moderndive.com/data/dem\\_score.xlsx](https://moderndive.com/data/dem_score.xlsx), but using RStudio's graphical interface instead of via the R console. First download the Excel file, then go to the **Files -> Import Dataset -> From Excel...** and navigate to the directory where your downloaded `dem_score.xlsx` using **Browse...**. You should see something similar to the image below:



After clicking on the **Import** button on the bottom-right save this spreadsheet's data in a data frame called `dem_score` and display its contents in the spreadsheet viewer (`View()`). Furthermore you'll see the code that read in your data in the console; you can copy and paste this code to reload your data again later instead of repeating the above manual process.

### Caution

Note that if you use the `xlsx` package to import `.xlsx` files is important to have the latest version of java installed in your local PC. The `xlsx` package depends on the `rJava` package which requires the Java Runtime Environment 1.2 or above. Download and install the latest version of the Java Runtime Environment from [Oracle](#).

### Task

Read in the life expectancy data stored at [https://moderndive.com/data/le\\_mess.csv](https://moderndive.com/data/le_mess.csv), either using the R console or RStudio's interface.

## Converting to tidy data format

In this section, we will see how to convert a data set that is not in the **tidy** format i.e. **wide** format, to a data set that is in the **tidy** format i.e. **long/narrow** format. Let's use the `dem_score` data frame we loaded from a spreadsheet in the previous section but focus on only data corresponding to the country of Guatemala.

```
guat_dem <- dem_score |>
  dplyr::filter(country == "Guatemala")
```

```
# A tibble: 1 x 10
  country `1952` `1957` `1962` `1967` `1972` `1977` `1982` `1987` `1992`
  <chr>    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 Guatemala      2     -6     -5      3      1     -3     -7      3      3
```

### Note

Here we have used the `filter` function from `dplyr` package to subset the data set. We will revisit this code for subsetting data later in the session.

Now let's produce a plot showing how the democracy scores have changed over the 40 years from 1952 to 1992 for Guatemala. Let's start by laying out how we would map our aesthetics to variables in the data frame:

- The data frame is `guat_dem` so we use `data = guat_dem`.

We would like to see how the democracy score has changed over the years in Guatemala. But we have a problem. We see that we have a variable named `country` but its only value is

Guatemala. We have other variables denoted by different year values. Unfortunately, we've run into a data set that is not in the appropriate format to apply the **Grammar of Graphics** in `ggplot2`. Remember that `ggplot2` is a package in the `tidyverse` and, thus, needs data to be in a tidy format. We'd like to finish off our mapping of aesthetics to variables by doing something like

- The aesthetic mapping is set by `aes(x = year, y = democracy_score)`,

but this is not possible with our wide-formatted data. We need to take the values of the current column names in `guat_dem` (aside from `country`) and convert them into a new variable that will act as a key called `year`. Then, we'd like to take the numbers on the inside of the table and turn them into a column that will act as values called `democracy_score`. Our resulting data frame will have three columns: `country`, `year`, and `democracy_score`.

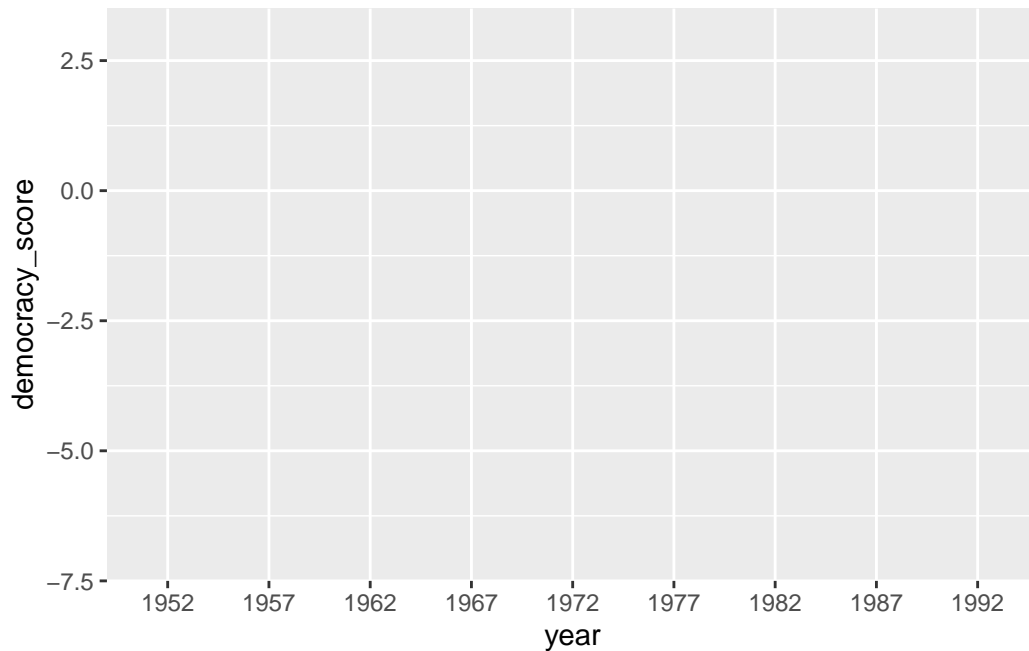
The `gather` function in the `tidyr` package can complete this task for us. The first argument to `gather`, just as with `ggplot2`, is the `data` argument where we specify which data frame we would like to tidy. The next two arguments to `gather` are `key` and `value`, which specify what we would like to call the new columns that convert our wide data into tidy/long format. Lastly, we include a specification for variables we would like to NOT include in the tidying process using a `-`.

```
guat_tidy <- gather(data = guat_dem,
                    key = year,
                    value = democracy_score,
                    - country)
```

```
# A tibble: 9 x 3
  country year democracy_score
  <chr>   <chr>           <dbl>
1 Guatemala 1952             2
2 Guatemala 1957            -6
3 Guatemala 1962            -5
4 Guatemala 1967             3
5 Guatemala 1972             1
6 Guatemala 1977            -3
7 Guatemala 1982            -7
8 Guatemala 1987             3
9 Guatemala 1992             3
```

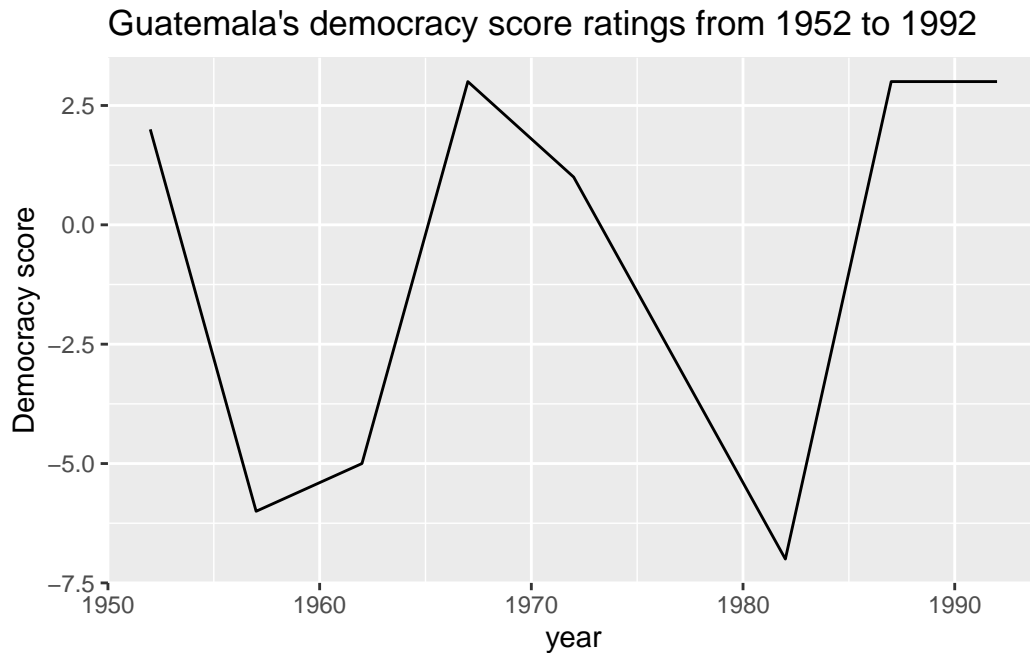
We can now create a plot showing how democracy score in Guatemala has changed from 1952 to 1992 using a linegraph and `ggplot2`.

```
ggplot(data = guat_tidy, mapping = aes(x = year, y = democracy_score)) +
  geom_line() +
  labs(x = "year")
```



Observe that the `year` variable in `guat_tidy` is stored as a character vector since we had to circumvent the naming rules in R by adding backticks around the different year columns in `guat_dem`. This is leading to `ggplot` not knowing exactly how to plot a line using a categorical variable. We can fix this by using the `parse_number` function in the `readr` package:

```
ggplot(data = guat_tidy, mapping = aes(x = parse_number(year), y = democracy_score)) +
  geom_line() +
  labs(x = "year", y = "Democracy score",
       title = "Guatemala's democracy score ratings from 1952 to 1992")
```



We'll see later how we could use the `mutate` function to change `year` to be a numeric variable during the tidying process (alternatively we could have added the argument `convert=T` in the `gather()` function to declare the `key` column values as an integers; see `?gather` for more details) . Notice now that the mappings of aesthetics to variables makes sense in the figure:

- The data frame is `guat_tidy` by setting `data = guat_tidy`;
- The `x` aesthetic is mapped to `year`;
- The `y` aesthetic is mapped to `democracy_score`; and
- The `geom_etry` chosen is `line`.

#### Task

Convert the `dem_score` data frame into a tidy data frame and assign the name of `dem_score_tidy` to the resulting long-formatted data frame.

Take hint

See the documentation for `gather()` (`?gather`). Try using the `convert= T` argument and comment on the output.

[Click here to see the solution](#)

```
dem_score_tidy <- gather(data = dem_score,
                        key = year,
                        convert= T,
                        value = democracy_score,
                        - country)

head(dem_score_tidy)
```

```
# A tibble: 6 x 3
  country    year democracy_score
  <chr>      <int>          <dbl>
1 Albania   1952            -9
2 Argentina 1952            -9
3 Armenia   1952            -9
4 Australia 1952             10
5 Austria   1952             10
6 Azerbaijan 1952            -9
```

### Task

Now try converting the life expectancy data set you created in a previous task into a tidy data frame.

## Introduction to data wrangling

We are now able to import data and perform basic operations on the data to get it into the **tidy** format. In this and subsequent sections we will use tools from the **dplyr** package to perform data **wrangling** which includes transforming, mapping and summarising variables.

### The pipe |>

Before we dig into data wrangling, let's first introduce the pipe operator (**|>**). Just as the **+** sign was used to add layers to a plot created using **ggplot**, the pipe operator allows us to chain together data wrangling functions. The pipe operator can be read as **then**. The **|>** operator allows us to go from one step in to the next easily so we can, for example:

- **filter** our data frame to only focus on a few rows **then**
- **group\_by** another variable to create groups **then**
- **summarize** this grouped data to calculate the mean for each level of the group.



The piping syntax will be our major focus throughout the rest of this course and you'll find that you'll quickly be addicted to the chaining with some practice.

## Data wrangling verbs

The **d** in **dplyr** stands for data frames, so the functions in **dplyr** are built for working with objects of the data frame type. For now, we focus on the most commonly used functions that help wrangle and summarise data. A description of these verbs follows, with each subsequent section devoted to an example of that verb, or a combination of a few verbs, in action.

1. **filter**: Pick rows based on conditions about their values
2. **summarize**: Compute summary measures known as “summary statistics” of variables
3. **group\_by**: Group rows of observations together
4. **mutate**: Create a new variable in the data frame by mutating existing ones
5. **arrange**: Arrange/sort the rows based on one or more variables
6. **join**: Join/merge two data frames by matching along a “key” variable. There are many different joins available. Here, we will focus on the **inner\_join** function.

All of the verbs are used similarly where you: take a data frame, pipe it using the `%>%` syntax into one of the verbs above followed by other arguments specifying which criteria you would like the verb to work with in parentheses.

## Filter observations using filter

### Subset Observations (Rows)



The **filter** function allows you to specify criteria about values of a variable in your data set and then chooses only those rows that match that criteria.

#### ! Important

Recall that the base R has already a **filter** function defined. So make sure to avoid any conflicts either by calling `dplyr::filter()` every time you use the function (specially if you have loaded the **conflicts** library) or alternatively run `theconflict_prefer()` function to let R know that it should use **dplyr**'s **filter** function as default.

Warning: package 'conflicted' was built under R version 4.2.3

```
conflict_prefer("filter", "dplyr")
```

[conflicted] Will prefer dplyr::filter over any other package.

We begin by focusing only on flights from New York City to Portland, Oregon. The `dest` code (or airport code) for Portland, Oregon is PDX. Run the following code and look at the resulting spreadsheet to ensure that only flights heading to Portland are chosen:

```
portland_flights <- flights |>
  filter(dest == "PDX")
# We do not display columns 6-11 so we can see the destination (dest) variable.
portland_flights[,-(6:12)]
```

```
# A tibble: 1,354 x 12
   year month   day dep_time sched_dep_time origin dest  air_time distance
   <int> <int> <int>   <int>         <int> <chr>  <chr>    <dbl>    <dbl>
1  2013     1     1    1739           1740 JFK   PDX      341     2454
2  2013     1     1    1805           1757 EWR   PDX      336     2434
3  2013     1     1    2052           2029 JFK   PDX      331     2454
4  2013     1     2     804            805 EWR   PDX      310     2434
5  2013     1     2    1552           1550 JFK   PDX      305     2454
6  2013     1     2    1727           1720 EWR   PDX      351     2434
7  2013     1     2    1738           1740 JFK   PDX      322     2454
8  2013     1     2    2024           2029 JFK   PDX      325     2454
9  2013     1     3    1755           1745 JFK   PDX      325     2454
10 2013     1     3    1814           1727 EWR   PDX      320     2434
# i 1,344 more rows
# i 3 more variables: hour <dbl>, minute <dbl>, time_hour <dtm>
```

Note the following:

- The ordering of the commands:
  - Take the data frame `flights` **then**
  - `filter` the data frame so that only those where the `dest` equals PDX are included.
- The double equals sign `==` tests equality, and not a single equals sign `=`.

You can combine multiple criteria together using operators that make comparisons:

- `|` corresponds to **or**
- `&` corresponds to **and**

We can often skip the use of `&` and just separate our conditions with a comma. You'll see this in the example below.

In addition, you can use other mathematical checks (similar to `==`):

- `>` corresponds to **greater than**
- `<` corresponds to **less than**
- `>=` corresponds to **greater than or equal to**
- `<=` corresponds to **less than or equal to**
- `!=` corresponds to **not equal to**

To see many of these in action, let's select all flights that left JFK airport heading to Burlington, Vermont (BTV) or Seattle, Washington (SEA) in the months of October, November, or December. Run the following

```
btv_sea_flights_fall <- flights |>
  filter(origin == "JFK", (dest == "BTV" | dest == "SEA"), month >= 10)
btv_sea_flights_fall[,-(6:12)]
```

# A tibble: 815 x 12

	year	month	day	dep_time	sched_dep_time	origin	dest	air_time	distance
	<int>	<int>	<int>	<int>	<int>	<chr>	<chr>	<dbl>	<dbl>
1	2013	10	1	729	735	JFK	SEA	352	2422
2	2013	10	1	853	900	JFK	SEA	362	2422
3	2013	10	1	916	925	JFK	BTV	48	266
4	2013	10	1	1216	1221	JFK	BTV	49	266
5	2013	10	1	1452	1459	JFK	BTV	46	266
6	2013	10	1	1459	1500	JFK	SEA	348	2422
7	2013	10	1	1754	1800	JFK	SEA	338	2422
8	2013	10	1	1825	1830	JFK	SEA	366	2422
9	2013	10	1	1925	1930	JFK	SEA	332	2422
10	2013	10	1	2238	2245	JFK	BTV	48	266

# i 805 more rows

# i 3 more variables: hour <dbl>, minute <dbl>, time\_hour <dtm>

#### **i** Note

Even though colloquially speaking one might say “all flights leaving Burlington, Vermont *and* Seattle, Washington,” in terms of computer logical operations, we really mean “all flights leaving Burlington, Vermont *or* Seattle, Washington.” For a given row in the data, **dest** can be BTV, SEA, or something else, but not BTV **and** SEA at the same time.

Another example uses `!` to pick rows that *do not* match a condition. The `!` can be read as **not**. Here, we are selecting rows corresponding to flights that **did not** go to Burlington, VT or Seattle, WA.

```
not_BTV_SEA <- flights |>
  filter(!(dest == "BTV" | dest == "SEA"))
not_BTV_SEA[,-(6:12)]
```

```
# A tibble: 330,264 x 12
```

	year	month	day	dep_time	sched_dep_time	origin	dest	air_time	distance
	<int>	<int>	<int>	<int>	<int>	<chr>	<chr>	<dbl>	<dbl>
1	2013	1	1	517	515	EWB	IAH	227	1400
2	2013	1	1	533	529	LGA	IAH	227	1416
3	2013	1	1	542	540	JFK	MIA	160	1089
4	2013	1	1	544	545	JFK	BQN	183	1576
5	2013	1	1	554	600	LGA	ATL	116	762
6	2013	1	1	554	558	EWB	ORD	150	719
7	2013	1	1	555	600	EWB	FLL	158	1065
8	2013	1	1	557	600	LGA	IAD	53	229
9	2013	1	1	557	600	JFK	MCO	140	944
10	2013	1	1	558	600	LGA	ORD	138	733

```
# i 330,254 more rows
```

```
# i 3 more variables: hour <dbl>, minute <dbl>, time_hour <dtm>
```

```
# We do not display columns 6-11 so we can see the "origin" and "dest" variables.
```

As a final note we point out that `filter` should often be the first verb you'll apply to your data. This narrows down the data to just the observations you are interested in.

#### Task

What is another way of using the **not** operator `!` to filter only the rows that are not going to Burlington, VT nor Seattle, WA in the `flights` data frame?

Take a hint

Try using the `%in%` operator

[Click here to see the solution](#)

```
flights |>
  filter(!dest %in% c("BTV","SEA")) |>
  head()
```

```
# A tibble: 6 x 19
```

```

  year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
<int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
1  2013     1     1     517           515           2     830           819
2  2013     1     1     533           529           4     850           830
3  2013     1     1     542           540           2     923           850
4  2013     1     1     544           545          -1    1004          1022
5  2013     1     1     554           600          -6     812           837
6  2013     1     1     554           558          -4     740           728
# i 11 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
#   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
#   hour <dbl>, minute <dbl>, time_hour <dtm>

```

## Summarise variables using summarize

The next common task is to be able to summarise data: take a large number of values and summarise them with a single value. While this may seem like a very abstract idea, something as simple as the sum, the smallest value, and the largest values are all summaries of a large number of values.



We can calculate the standard deviation and mean of the temperature variable `temp` in the `weather` data frame of `nycflights13` in one step using the `summarize` (or equivalently using the UK spelling `summarise`) function in `dplyr`

```

summary_temp <- weather |>
  summarize(mean = mean(temp), std_dev = sd(temp))

```

mean	std_dev
NA	NA

We have created a small data frame here called `summary_temp` that includes both the mean (`mean`) and standard deviation (`std_dev`) of the `temp` variable in `weather`. Notice, the data frame `weather` went from many rows to a single row of just the summary values in the data frame `summary_temp`.

But why are the values returned `NA`? This stands for **not available or not applicable** and is how R encodes **missing values**; if in a data frame for a particular row and column no value exists, `NA` is stored instead. Furthermore, by default any time you try to summarise a number

of values (using `mean()` and `sd()` for example) that has one or more missing values, then `NA` is returned.

Values can be missing for many reasons. Perhaps the data was collected but someone forgot to enter it? Perhaps the data was not collected at all because it was too difficult? Perhaps there was an erroneous value that someone entered that was changed to read as missing? You'll often encounter issues with missing values.

You can summarise all non-missing values by setting the `na.rm` argument to `TRUE` (`rm` is short for remove). This will remove any `NA` missing values and only return the summary value for all non-missing values. So the code below computes the mean and standard deviation of all non-missing values. Notice how the `na.rm=TRUE` are set as arguments to the `mean` and `sd` functions, and not to the `summarize` function.

```
summary_temp <- weather |>
  summarize(mean = mean(temp, na.rm = TRUE), std_dev = sd(temp, na.rm = TRUE))
```

mean	std_dev
55.26039	17.78785

Another very useful function that allows you to summarise multiple columns is the `summarise_at()` function. Here, we can supply a vector of variables we want to summarise and a list of functions that we want to apply to each of the variables. See the following example where we compute the minimum and maximum values of temperature and relative humidity (notice that we also specified the argument `na.rm=T` to remove missing values - this arguments gets passed on to all the functions in the list):

```
weather |>
  summarise_at(.vars = c("temp","humid"),
    .funs = list(min = min, max = max),
    na.rm = TRUE)
```

```
# A tibble: 1 x 4
  temp_min humid_min temp_max humid_max
  <dbl>      <dbl>    <dbl>    <dbl>
1    10.9     12.7    100.     100
```

It is **not** good practice to include `na.rm = TRUE` in your summary commands by default; you should attempt to run code first without this argument as this will alert you to the presence of missing data. Only after you have identified where missing values occur and have thought about the potential issues of these should you consider using `na.rm = TRUE`. In the upcoming

Tasks we will consider the possible ramifications of blindly sweeping rows with missing values under the rug.

What other summary functions can we use inside the `summarize` verb? Any function in R that takes a vector of values and returns just one. Here are just a few:

- `mean`: the mean (or average)
- `sd`: the standard deviation, which is a measure of spread
- `min` and `max`: the minimum and maximum values, respectively
- `IQR`: the interquartile range
- `sum`: the sum
- `n`: a count of the number of rows/observations in each group. This particular summary function will make more sense when `group_by` is used in the next section.

#### Question

Say a doctor is studying the effect of smoking on lung cancer for a large number of patients who have records measured at five year intervals. She notices that a large number of patients have missing data points because the patient has died, so she chooses to ignore these patients in her analysis. What is wrong with this doctor's approach?

- (A) Introduces a selection bias since patient who died due to lung cancer are excluded from the analysis, leading to an underestimation of the true impact of smoking on lung cancer risk
- (B) There is no problem, smaller datasets with fewer missing values may require less computational resources, leading to faster processing times.
- (C) Removing patients with missing data reduces the sample size. Hence, conclusions may not be as easily generalizable to the broader population, as the excluded patients may represent a different subset with unique characteristics.
- (D) Removing missing values can result in a dataset with fewer errors and inconsistencies, which can lead to more accurate analyses.

#### Question

Modify `summary_temp` from above to also use the `n` summary function: `summarize(count = n())`. What does the returned value correspond to?

- (A) Number of weather stations

- (B) Number of columns in the data
- (C) Sample size

### Question

Why does the code below not work?

```
summary_temp <- weather |>
  summarize(mean = mean(temp, na.rm = TRUE)) |>
  summarize(std_dev = sd(temp, na.rm = TRUE))
```

```
Error in `summarize()` :
i In argument: `std_dev = sd(temp, na.rm = TRUE)`.
Caused by error in `is.data.frame()` :
! object 'temp' not found
```

Take hint

Run the code line by line instead of all at once, and then look at the data. In other words, run `summary_temp <- weather |> summarize(mean = mean(temp, na.rm = TRUE))` first.

Answer

The first line of code computes the temperature mean for the weather data set. Then, the output gets passed on to `weather |> summarize(std_dev= sd(temp, na.rm = TRUE))`. However, the temperature value is no longer present in the first result, hence the error: `! object 'temp' not found`

## Group rows using grouping structures



It is often more useful to summarise a variable based on the groupings of another variable. Let's say we are interested in the mean and standard deviation of temperatures but *grouped by month*. To be more specific: we want the mean and standard deviation of temperatures

1. split by month.



2. sliced by month.
3. aggregated by month.
4. collapsed over month.

Run the following code:

```
summary_monthly_temp <- weather |>
  summarize(mean = mean(temp, na.rm = TRUE),
            std_dev = sd(temp, na.rm = TRUE),
            .by = month)
```

This code is identical to the previous code that created `summary_temp`, with an extra `.by = month` added. This kind per-operation grouping allow us to do the grouping within the operation where the summarisation takes place without changing the structure of the data .

#### **i** Note

Previous versions of `dplyr` relied on the specification of a `group_by` function within the pipeline to do the grouping. For example, in the next line of code the `weather` data set is initially grouped by `month` and then passed as a new grouped data frame into `summarize`. Yielding to the same data frame that shows the mean and standard deviation of temperature for each month in New York City:

```
summary_monthly_temp <- weather |>
  group_by(month) |>
  summarize(mean = mean(temp, na.rm = TRUE),
            std_dev = sd(temp, na.rm = TRUE))
```

While `group_by` doesn't change the data frame, it sets *meta-data* (data about the data), specifically the group structure of the data. If we wanted to remove this group structure meta-data, we could add the `.groups = "drop"` option.

```
summary_monthly_temp <- weather |>
  group_by(month) |>
  summarize(mean = mean(temp, na.rm = TRUE),
            std_dev = sd(temp, na.rm = TRUE),
            .groups = "drop")
```

The advantage of using the `.by` argument is that the grouping occurs within the `summarize` function and thus the resulting data frame is no longer grouped.

### Question

The `drop_na()` function can be used in the pipeline to remove missing observations from a data set. Try running the following code to compute the mean and standard deviation of the temperature in the `weather` data set and comment on the output. Why is this different from one we had before?

```
summary_monthly_temp <- weather |>
  drop_na() |>
  summarize(mean = mean(temp),
            std_dev = sd(temp),
            .by = month)
```

### Answer

The `drop_na()` function remove all missing observation from the data set while specifying `na.rm = T` in each summarizing function only removes the missing values for the specific variable to which the function is applied.

We now revisit the `n` counting summary function we introduced in the previous section. For example, suppose we would like to get a sense for how many flights departed each of the three airports in New York City:

```
by_origin <- flights |>
  summarize(count = n(),
            .by = origin)
```

origin	count
EWB	120835
JFK	111279
LGA	104662

We see that Newark (EWB) had the most flights departing in 2013 followed by JFK and lastly by LaGuardia (LGA). Note, there is a subtle but important difference between `sum` and `n`. While `sum` simply adds up a large set of numbers, the latter counts the number of times each of many different values occur.

## Grouping by more than one variable

You are not limited to grouping by one variable. Say you wanted to know the number of flights leaving each of the three New York City airports *for each month*, we can also group by a second variable `month`:

```
by_origin_monthly <- flights |>
  summarize(count = n(),
            .by = c(origin, month))
by_origin_monthly
```

```
# A tibble: 36 x 3
  origin month count
  <chr>   <int> <int>
1 EWR      1  9893
2 LGA      1  7950
3 JFK      1  9161
4 EWR     10 10104
5 JFK     10  9143
6 LGA     10  9642
7 JFK     11  8710
8 EWR     11  9707
9 LGA     11  8851
10 JFK     12  9146
# i 26 more rows
```

We see there are 36 rows for `by_origin_monthly` because there are 12 months times 3 airports (EWR, JFK, and LGA). Let's now pose a question.

1. First, what if we reverse the order of the grouping, i.e. `.by = c(month, origin)`?

```
by_monthly_origin <- flights |>
  summarize(count = n(),
            .by = c(month, origin))
by_monthly_origin
```

```
# A tibble: 36 x 3
  month origin count
  <int> <chr>   <int>
1      1 EWR    9893
2      1 LGA    7950
3      1 JFK    9161
4     10 EWR   10104
5     10 JFK    9143
6     10 LGA    9642
7     11 JFK    8710
8     11 EWR    9707
```

```

 9      11 LGA      8851
10      12 JFK      9146
# i 26 more rows

```

In `by_monthly_origin` the `month` column is now first and the rows are sorted by `month` instead of `origin`. If you compare the values of `count` in `by_origin_monthly` and `by_monthly_origin` using the `View` function, you'll see that the values are actually the same, just presented in a different order.

### Question

Recall from Week 1 when we looked at plots of temperatures by months in NYC. What does the standard deviation column in the `summary_monthly_temp` data frame tell us about temperatures in New York City throughout the year?

- (A) Temperature are lower in the winter
- (B) Temperature variability increases during winter
- (C) Spring is the season with more outliers

### Task

Write code to produce the mean and standard deviation temperature for each day in 2013 for NYC

Take a hint

See the documentation for `plot()` (`?plot`)

[Click here to see the solution](#)

```

weather |>
  summarize(mean = mean(temp, na.rm = TRUE),
            std_dev = sd(temp, na.rm = TRUE),
            .by = day)

# A tibble: 31 x 3
   day mean std_dev
<int> <dbl> <dbl>
1     1  57.6   17.4
2     2  55.7   20.2
3     3  53.8   18.9
4     4  54.0   18.8

```

```
5      5  55.6    16.2
6      6  55.7    15.6
7      7  55.6    17.4
8      8  55.0    17.6
9      9  56.6    17.4
10     10  56.9    17.8
# i 21 more rows
```

### Task

How could we identify how many flights left each of the three airports for each **carrier**?

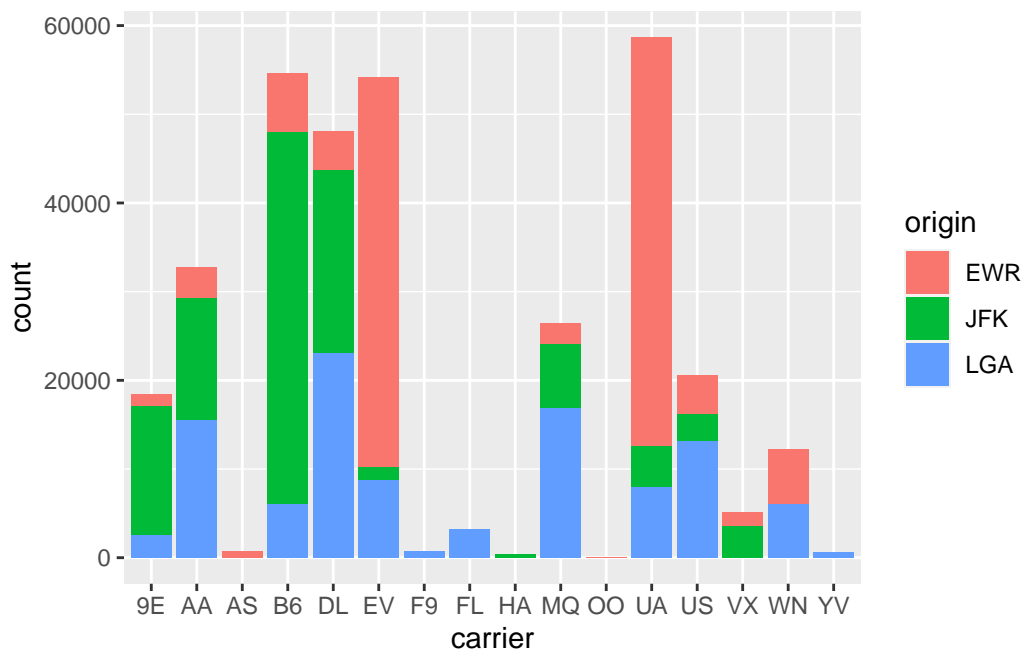
Can you create a bar plot showing these results?

Take a hint

You can count how many flights left each of the three airports by summarising the data using the `n()` function while grouping by the origin and carrier. Then, you can pass the resulting data frame to `ggplot` using the pipeline command `|>` and use a `geom_col` layer as we saw in the previous week.

[Click here to see the solution](#)

```
flights |>
  summarise(count = n(),
            .by = c(origin, carrier)) |>
  ggplot(aes(x = carrier, y = count, fill = origin)) + geom_col()
```



## Create new variables/change old variables using mutate

### Make New Variables



When looking at the `flights` data set, there are some clear additional variables that could be calculated based on the values of variables already in the data set. Passengers are often frustrated when their flights depart late, but change their mood a bit if pilots can make up some time during the flight to get them to their destination close to when they expected to land. This is commonly referred to as “gain” and we will create this variable using the `mutate` function. Note that we will be overwriting the `flights` data frame with one including the additional variable `gain` here, or put differently, the `mutate` command outputs a new data frame which then gets saved over the original `flights` data frame.

```
flights <- flights |>
  mutate(gain = dep_delay - arr_delay)
```

Let's take a look at `dep_delay`, `arr_delay`, and the resulting `gain` variables in our new `flights` data frame:

```
# A tibble: 336,776 x 3
  dep_delay arr_delay gain
  <dbl>      <dbl> <dbl>
1         2         11  -9
2         4         20 -16
3         2         33 -31
4        -1        -18  17
5        -6        -25  19
6        -4         12 -16
7        -5         19 -24
8        -3        -14  11
9        -3         -8   5
10       -2          8 -10
# i 336,766 more rows
```

The flight in the first row departed 2 minutes late but arrived 11 minutes late, so its “gained time in the air” is actually a loss of 9 minutes, hence its `gain` is `-9`. Contrast this to the flight in the fourth row which departed a minute early (`dep_delay` of `-1`) but arrived 18 minutes early (`arr_delay` of `-18`), so its “gained time in the air” is 17 minutes, hence its `gain` is `+17`.

Why did we overwrite `flights` instead of assigning the resulting data frame to a new object, like `flights_with_gain`? As a rough rule of thumb, as long as you are not losing information that you might need later, it's acceptable practice to overwrite data frames. However, if you overwrite existing variables and/or change the observational units, recovering the original information might prove difficult. In this case, it might make sense to create a new data object.

Let's look at summary measures of this `gain` variable and plot it in the form of a histogram:

```
gain_summary <- flights |>
  summarize(
    min = min(gain, na.rm = TRUE),
    q1 = quantile(gain, 0.25, na.rm = TRUE),
    median = quantile(gain, 0.5, na.rm = TRUE),
    q3 = quantile(gain, 0.75, na.rm = TRUE),
    max = max(gain, na.rm = TRUE),
    mean = mean(gain, na.rm = TRUE),
    sd = sd(gain, na.rm = TRUE),
    missing = sum(is.na(gain))
```

min	q1	median	q3	max	mean	sd	missing
-196	-3	7	17	109	5.659779	18.04365	9430

We have recreated the `summary` function we saw in Week 1 here using the `summarize` function in `dplyr`. Lets make a histogram for the new created `gain` variable.

```
ggplot(data = flights, mapping = aes(x = gain)) +
  geom_histogram(color = "white", bins = 20)
```

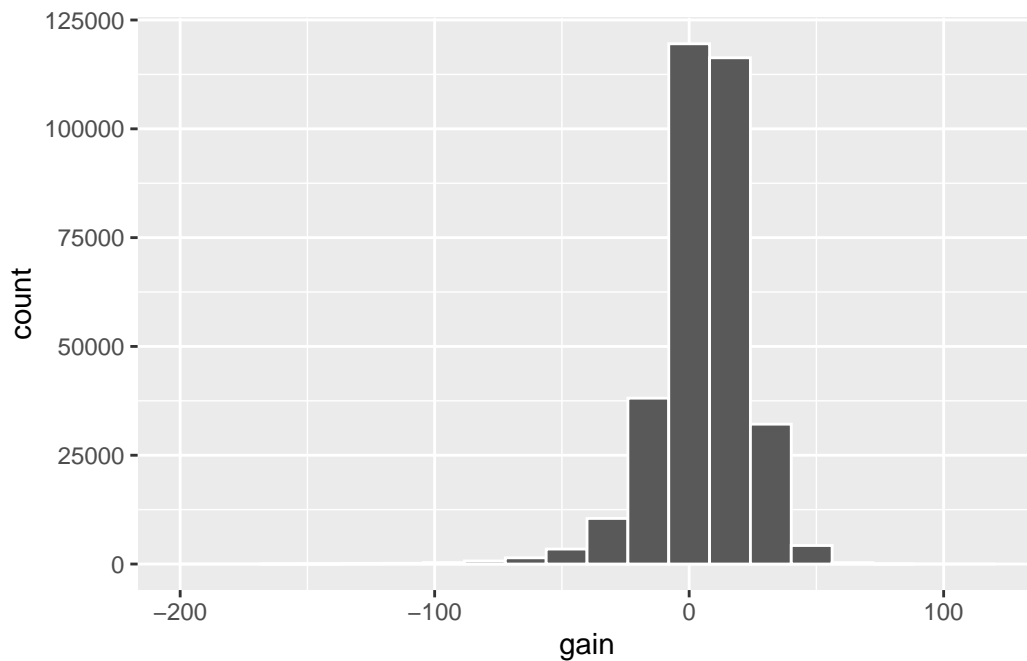


Figure 2: Histogram of gain variable.

We can also create multiple columns at once and even refer to columns that were just created in a new column.

```
flights <- flights |>
  mutate(
    gain = dep_delay - arr_delay,
    hours = air_time / 60,
    gain_per_hour = gain / hours
  )
```



### Question

What do positive values of the `gain` variable in `flights` correspond to?

- (A) Departure delays are greater than arrivals delays
- (B) Departure delays are lower than arrivals delays
- (C) Departures and arrivals delays are the same

What about negative values?

- (A) Departure delays are greater than arrivals delays
- (B) Departure delays are lower than arrivals delays
- (C) Departures and arrivals delays are the same

And what about a zero value?

- (A) Departure delays are greater than arrivals delays
- (B) Departure delays are lower than arrivals delays
- (C) Departures and arrivals delays are the same

### Question

Could we create the `dep_delay` and `arr_delay` columns by simply subtracting `dep_time` from `sched_dep_time` and similarly for arrivals? Try the code out and explain any differences between the result and what actually appears in `flights`.

```
flights |>
  mutate(dep_delay = sched_dep_time - dep_time ,
         arr_delay = sched_arr_time - arr_time)
```

Take a hint

See the description of the variables `arr_time`, `dep_time`, `sched_dep_time` and

`sched_arr_time` in the `flights` data set? `flights`

Answer

The differences are due to departure and arrival times have a HHMM or HMM format. E.g., if we compute the difference between a flight scheduled to arrive by 923 and its actual arrival time at 850, the result would be a difference of 73, while in reality there was only a 33 min difference if we consider the correct time format!

## Reorder the data frame using `arrange`

One of the most common things people working with data would like to do is sort the data frames by a specific variable in a column. Have you ever been asked to calculate a median by hand? This requires you to put the data in order from smallest to highest in value. The `dplyr` package has a function called `arrange` that we will use to sort/reorder our data according to the values of the specified variable. This is often used after we have grouped and summarized the data as we will see.

Let's suppose we were interested in determining the most frequent destination airports from New York City in 2013:

```
freq_dest <- flights |>
  summarize(num_flights = n(),
            .by = dest)
```

```
# A tibble: 105 x 2
  dest   num_flights
  <chr>         <int>
1 IAH             7198
2 MIA            11728
3 BQN              896
4 ATL            17215
5 ORD            17283
6 FLL            12055
7 IAD              5700
8 MCO            14082
9 PBI              6554
10 TPA            7466
# i 95 more rows
```

You'll see that by default the values of `dest` are displayed in alphabetical order here. We are interested in finding those airports that appear most:

```
freq_dest |>
  arrange(num_flights)
```

```
# A tibble: 105 x 2
  dest  num_flights
  <chr>      <int>
1 LEX          1
2 LGA          1
3 ANC          8
4 SBN         10
5 HDN         15
6 MTJ         15
7 EYW         17
8 PSP         19
9 JAC         25
10 BZN        36
# i 95 more rows
```

This is actually giving us the opposite of what we are looking for. It tells us the least frequent destination airports first. To switch the ordering to be descending instead of ascending we use the `desc` (descending) function:

```
freq_dest |>
  arrange(desc(num_flights))
```

```
# A tibble: 105 x 2
  dest  num_flights
  <chr>      <int>
1 ORD     17283
2 ATL     17215
3 LAX     16174
4 BOS     15508
5 MCO     14082
6 CLT     14064
7 SFO     13331
8 FLL     12055
9 MIA     11728
10 DCA      9705
# i 95 more rows
```

## Joining data frames

Another common task is joining (merging) two different data sets. For example, in the `flights` data, the variable `carrier` lists the carrier code for the different flights. While `UA` and `AA` might be somewhat easy to guess for some (United and American Airlines), what are `VX`, `HA`, and `B6`? This information is provided in a separate data frame `airlines`.

```
airlines
```

```
# A tibble: 16 x 2
  carrier name
  <chr>    <chr>
1 9E      Endeavor Air Inc.
2 AA      American Airlines Inc.
3 AS      Alaska Airlines Inc.
4 B6      JetBlue Airways
5 DL      Delta Air Lines Inc.
6 EV      ExpressJet Airlines Inc.
7 F9      Frontier Airlines Inc.
8 FL      AirTran Airways Corporation
9 HA      Hawaiian Airlines Inc.
10 MQ     Envoy Air
11 OO     SkyWest Airlines Inc.
12 UA     United Air Lines Inc.
13 US     US Airways Inc.
14 VX     Virgin America
15 WN     Southwest Airlines Co.
16 YV     Mesa Airlines Inc.
```

We see that in `airports`, `carrier` is the carrier code while `name` is the full name of the airline. Using this table, we can see that `VX`, `HA`, and `B6` correspond to Virgin America, Hawaiian Airlines Inc., and JetBlue Airways, respectively. However, will we have to continually look up the carrier's name for each flight in the `airlines` data set? No! Instead of having to do this manually, we can have R automatically do the “looking up” for us.

Note that the values in the variable `carrier` in `flights` match the values in the variable `carrier` in `airlines`. In this case, we can use the variable `carrier` as a *key variable* to join/merge/match the two data frames by. Key variables are almost always identification variables that uniquely identify the observational units as we saw back in the **Identification vs Measurement Variable** section. This ensures that rows in both data frames are appropriately matched during the join. This diagram helps us understand how the different data sets are linked by various key variables:

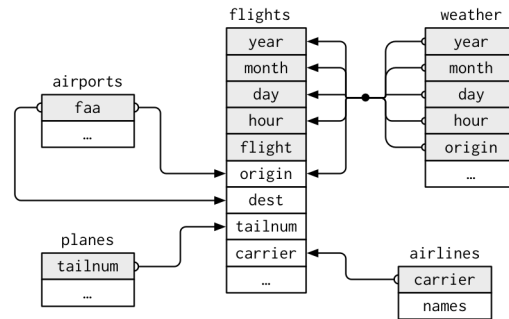


Figure 3: Data relationships in nycflights13 from R for Data Science, Hadley and Garrett (2016).

## Joining by “key” variables

In both `flights` and `airlines`, the key variable we want to join/merge/match the two data frames with has the same name in both data sets: `carriers`. We make use of the `inner_join` function to join by the variable `carrier`.

```
flights_joined <- flights |>
  inner_join(airlines,
            by = join_by(carrier))
```

```
flights
```

```
# A tibble: 336,776 x 22
```

	year	month	day	dep_time	sched_dep_time	dep_delay	arr_time	sched_arr_time
	<int>	<int>	<int>	<int>	<int>	<dbl>	<int>	<int>
1	2013	1	1	517	515	2	830	819
2	2013	1	1	533	529	4	850	830
3	2013	1	1	542	540	2	923	850
4	2013	1	1	544	545	-1	1004	1022
5	2013	1	1	554	600	-6	812	837
6	2013	1	1	554	558	-4	740	728
7	2013	1	1	555	600	-5	913	854
8	2013	1	1	557	600	-3	709	723
9	2013	1	1	557	600	-3	838	846
10	2013	1	1	558	600	-2	753	745

```
# i 336,766 more rows
```

```
# i 14 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
#   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
#   hour <dbl>, minute <dbl>, time_hour <dtm>, gain <dbl>, hours <dbl>,
```

```
# gain_per_hour <dbl>
```

```
flights_joined
```

```
# A tibble: 336,776 x 23
```

```
   year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
   <int> <int> <int>   <int>         <int>         <dbl>   <int>         <int>
1  2013     1     1     517           515           2     830           819
2  2013     1     1     533           529           4     850           830
3  2013     1     1     542           540           2     923           850
4  2013     1     1     544           545          -1    1004          1022
5  2013     1     1     554           600          -6     812           837
6  2013     1     1     554           558          -4     740           728
7  2013     1     1     555           600          -5     913           854
8  2013     1     1     557           600          -3     709           723
9  2013     1     1     557           600          -3     838           846
10 2013     1     1     558           600          -2     753           745
# i 336,766 more rows
# i 15 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
#   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
#   hour <dbl>, minute <dbl>, time_hour <dtm>, gain <dbl>, hours <dbl>,
#   gain_per_hour <dbl>, name <chr>
```

We observe that the `flights` and `flights_joined` are identical except that `flights_joined` has an additional variable `name` whose values were drawn from `airlines`.

A visual representation of the `inner_join` is given below:

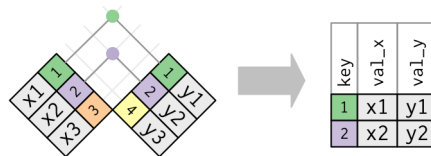


Figure 4: Diagram of inner join from R for Data Science.

There are more complex joins available, but the `inner_join` will solve nearly all of the problems you will face here.

## Joining by “key” variables with different names

Say instead, you are interested in all the destinations of flights from NYC in 2013 and ask yourself:

- “What cities are these airports in?”
- “Is ORD Orlando?”
- “Where is FLL?”

The `airports` data frame contains airport codes:

```
airports
```

```
# A tibble: 1,458 x 8
```

	faa	name	lat	lon	alt	tz	dst	tzone
	<chr>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<chr>	<chr>
1	04G	Lansdowne Airport	41.1	-80.6	1044	-5	A	America/~
2	06A	Moton Field Municipal Airport	32.5	-85.7	264	-6	A	America/~
3	06C	Schaumburg Regional	42.0	-88.1	801	-6	A	America/~
4	06N	Randall Airport	41.4	-74.4	523	-5	A	America/~
5	09J	Jekyll Island Airport	31.1	-81.4	11	-5	A	America/~
6	0A9	Elizabethton Municipal Airport	36.4	-82.2	1593	-5	A	America/~
7	0G6	Williams County Airport	41.5	-84.5	730	-5	A	America/~
8	0G7	Finger Lakes Regional Airport	42.9	-76.8	492	-5	A	America/~
9	0P2	Shoestring Aviation Airfield	39.8	-76.6	1000	-5	U	America/~
10	0S9	Jefferson County Intl	48.1	-123.	108	-8	A	America/~

```
# i 1,448 more rows
```

However, looking at both the `airports` and `flights` and the visual representation of the relations between the data frames in the figure above, we see that in:

- `airports` the airport code is in the variable `faa`
- `flights` the airport code is in the variable `origin`

So to join these two data sets, our `inner_join` operation involves a logical operator `==` argument that accounts for the different names.

```
flights |>
  inner_join(airports,
            by = join_by(dest == faa))
```

We can read the code out loud as:

```
""Take the flights data frame and inner join it to the airports data frame by the entries
  where the variable* 'dest' *is equal to* 'faa'""
```

Let's construct the sequence of commands that computes the number of flights from NYC to each destination, but also includes information about each destination airport:

```

named_dests <- flights|>
  summarize(num_flights = n(),
            .by = dest) |>
  arrange(desc(num_flights)) |>
  inner_join(airports, by = join_by(dest == faa)) %>%
  rename(airport_name = name)

# A tibble: 101 x 9
  dest   num_flights airport_name      lat   lon   alt   tz dst tzone
  <chr>     <int> <chr>          <dbl> <dbl> <dbl> <dbl> <chr> <chr>
1 ORD       17283 Chicago Ohare Intl    42.0 -87.9   668   -6 A Amer~
2 ATL       17215 Hartsfield Jackson At~  33.6 -84.4  1026   -5 A Amer~
3 LAX       16174 Los Angeles Intl     33.9 -118.   126   -8 A Amer~
4 BOS       15508 General Edward Lawren~  42.4 -71.0    19   -5 A Amer~
5 MCO       14082 Orlando Intl         28.4 -81.3    96   -5 A Amer~
6 CLT       14064 Charlotte Douglas Intl   35.2 -80.9   748   -5 A Amer~
7 SFO       13331 San Francisco Intl     37.6 -122.    13   -8 A Amer~
8 FLL       12055 Fort Lauderdale Holly~  26.1 -80.2     9   -5 A Amer~
9 MIA       11728 Miami Intl           25.8 -80.3     8   -5 A Amer~
10 DCA        9705 Ronald Reagan Washing~  38.9 -77.0    15   -5 A Amer~
# i 91 more rows

```

In case you didn't know, `ORD` is the airport code of Chicago O'Hare airport and `FLL` is the main airport in Fort Lauderdale, Florida, which we can now see in our `named_dests` data frame.

## Joining by multiple “key” variables

Say instead we are in a situation where we need to join by multiple variables. For example, in the first figure in this section we see that in order to join the `flights` and `weather` data frames, we need more than one key variable: `year`, `month`, `day`, `hour`, and `origin`. This is because the combination of these 5 variables act to uniquely identify each observational unit in the `weather` data frame: hourly weather recordings at each of the 3 NYC airports.

We achieve this by specifying a vector of key variables to join by.

```

flights_weather_joined <- flights |>
  inner_join(weather,
            by = join_by(year,month,day,hour,origin))

```

```

# A tibble: 335,220 x 32

```



	year	month	day	dep_time	sched_dep_time	dep_delay	arr_time	sched_arr_time
	<int>	<int>	<int>	<int>	<int>	<dbl>	<int>	<int>
1	2013	1	1	517	515	2	830	819
2	2013	1	1	533	529	4	850	830
3	2013	1	1	542	540	2	923	850
4	2013	1	1	544	545	-1	1004	1022
5	2013	1	1	554	600	-6	812	837
6	2013	1	1	554	558	-4	740	728
7	2013	1	1	555	600	-5	913	854
8	2013	1	1	557	600	-3	709	723
9	2013	1	1	557	600	-3	838	846
10	2013	1	1	558	600	-2	753	745

```
# i 335,210 more rows
# i 24 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
#   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
#   hour <dbl>, minute <dbl>, time_hour.x <dtm>, gain <dbl>, hours <dbl>,
#   gain_per_hour <dbl>, temp <dbl>, dewp <dbl>, humid <dbl>, wind_dir <dbl>,
#   wind_speed <dbl>, wind_gust <dbl>, precip <dbl>, pressure <dbl>,
#   visib <dbl>, time_hour.y <dtm>
```

### Question

Looking at the first figure in this section, when joining `flights` and `weather` (or, in other words, matching the hourly weather values with each flight), why do we need to join by all of `year`, `month`, `day`, `hour`, and `origin`, and not just `hour`?

Answer

`year`, `month`, `day`, `hour`, `origin` are the key variables that allow us to uniquely identify the observational units.

## Other verbs

### Select variables using `select`

We've seen that the `flights` data frame in the `nycflights13` package contains many different variables. The `names` function gives a listing of all the columns in a data frame; in our case you would run `names(flights)`. You can also identify these variables by running the `glimpse` function in the `dplyr` package:

```
glimpse(flights)
```

## Subset Variables (Columns)



Figure 5: Select diagram from Data Wrangling with dplyr and tidyr cheatsheet.

```

Rows: 336,776
Columns: 22
$ year      <int> 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2013, 2~
$ month     <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1~
$ day       <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1~
$ dep_time  <int> 517, 533, 542, 544, 554, 554, 555, 557, 557, 558, 558, ~
$ sched_dep_time <int> 515, 529, 540, 545, 600, 558, 600, 600, 600, 600, 600, ~
$ dep_delay <dbl> 2, 4, 2, -1, -6, -4, -5, -3, -3, -2, -2, -2, -2, -2, -1~
$ arr_time  <int> 830, 850, 923, 1004, 812, 740, 913, 709, 838, 753, 849,~
$ sched_arr_time <int> 819, 830, 850, 1022, 837, 728, 854, 723, 846, 745, 851,~
$ arr_delay <dbl> 11, 20, 33, -18, -25, 12, 19, -14, -8, 8, -2, -3, 7, -1~
$ carrier   <chr> "UA", "UA", "AA", "B6", "DL", "UA", "B6", "EV", "B6", "~
$ flight    <int> 1545, 1714, 1141, 725, 461, 1696, 507, 5708, 79, 301, 4~
$ tailnum   <chr> "N14228", "N24211", "N619AA", "N804JB", "N668DN", "N394~
$ origin    <chr> "EWR", "LGA", "JFK", "JFK", "LGA", "EWR", "EWR", "LGA",~
$ dest      <chr> "IAH", "IAH", "MIA", "BQN", "ATL", "ORD", "FLL", "IAD",~
$ air_time  <dbl> 227, 227, 160, 183, 116, 150, 158, 53, 140, 138, 149, 1~
$ distance  <dbl> 1400, 1416, 1089, 1576, 762, 719, 1065, 229, 944, 733, ~
$ hour      <dbl> 5, 5, 5, 5, 6, 5, 6, 6, 6, 6, 6, 6, 6, 6, 6, 5, 6, 6, 6~
$ minute    <dbl> 15, 29, 40, 45, 0, 58, 0, 0, 0, 0, 0, 0, 0, 0, 0, 59, 0~
$ time_hour <dtm> 2013-01-01 05:00:00, 2013-01-01 05:00:00, 2013-01-01 0~
$ gain      <dbl> -9, -16, -31, 17, 19, -16, -24, 11, 5, -10, 0, 1, -9, 1~
$ hours     <dbl> 3.7833333, 3.7833333, 2.6666667, 3.0500000, 1.9333333, ~
$ gain_per_hour <dbl> -2.3788546, -4.2290749, -11.6250000, 5.5737705, 9.82758~

```

However, say you only want to consider two of these variables, say `carrier` and `flight`. You can `select` these:

```

flights |>
  select(carrier, flight)

```

```
# A tibble: 336,776 x 2
```

```

      carrier flight
      <chr>    <int>
1 UA          1545
2 UA          1714
3 AA          1141
4 B6           725
5 DL           461
6 UA          1696
7 B6           507
8 EV          5708
9 B6            79
10 AA          301
# i 336,766 more rows

```

This function makes navigating data sets with a very large number of variables easier for humans by restricting consideration to only those of interest, like `carrier` and `flight` above. So for example, this might make viewing the data set using the `View` spreadsheet viewer more digestible. However, as far as the computer is concerned it does not care how many additional variables are in the data set in question, so long as `carrier` and `flight` are included.

Another example involves the variable `year`. If you remember the original description of the `flights` data frame (or by running `?flights`), you will remember that this data corresponds to flights in 2013 departing New York City. The `year` variable isn't really a variable here in that it doesn't vary, the `flights` data set actually comes from a larger data set that covers many years. We may want to remove the `year` variable from our data set since it won't be helpful for analysis in this case. We can deselect `year` by using the `-` sign:

```

flights_no_year <- flights |>
  select(-year)

```

Or we could specify a ranges of columns:

```

flight_arr_times <- flights |>
  select(month:dep_time, arr_time:sched_arr_time)

```

```

# A tibble: 336,776 x 5
  month   day dep_time arr_time sched_arr_time
  <int> <int>   <int>   <int>         <int>
1     1     1     517     830           819
2     1     1     533     850           830
3     1     1     542     923           850
4     1     1     544    1004          1022

```

```

5      1      1      554      812      837
6      1      1      554      740      728
7      1      1      555      913      854
8      1      1      557      709      723
9      1      1      557      838      846
10     1      1      558      753      745
# i 336,766 more rows

```

The `select` function can also be used to reorder columns in combination with the `everything` helper function. Let's suppose we would like the `hour`, `minute`, and `time_hour` variables, which appear at the end of the `flights` data set, to actually appear immediately after the `day` variable:

```

flights_reorder <- flights |>
  select(month:day, hour:time_hour, everything())

```

```

[1] "month"      "day"      "hour"      "minute"
[5] "time_hour"  "year"     "dep_time"  "sched_dep_time"
[9] "dep_delay"  "arr_time" "sched_arr_time" "arr_delay"
[13] "carrier"    "flight"   "tailnum"   "origin"
[17] "dest"       "air_time" "distance"  "gain"
[21] "hours"      "gain_per_hour"

```

in this case `everything()` picks up all remaining variables. Lastly, the helper functions `starts_with`, `ends_with`, and `contains` can be used to choose **variables / column names** that match those conditions:

```

flights_begin_a <- flights |>
  select(starts_with("a"))

```

```

# A tibble: 336,776 x 3
  arr_time arr_delay air_time
  <int>     <dbl>    <dbl>
1     830         11     227
2     850         20     227
3     923         33     160
4    1004        -18     183
5     812        -25     116
6     740         12     150
7     913         19     158

```

```

      8      709      -14      53
      9      838       -8     140
     10      753       8     138
# i 336,766 more rows

```

```

flights_delays <- flights |>
  select(ends_with("delay"))

```

```
# A tibble: 336,776 x 2
```

```

  dep_delay arr_delay
    <dbl>     <dbl>
1         2         11
2         4         20
3         2         33
4        -1        -18
5        -6        -25
6        -4         12
7        -5         19
8        -3        -14
9        -3         -8
10       -2          8

```

```
# i 336,766 more rows
```

```

flights_time <- flights |>
  select(contains("time"))

```

```
# A tibble: 336,776 x 6
```

```

  dep_time sched_dep_time arr_time sched_arr_time air_time time_hour
    <int>         <int>     <int>         <int>     <dbl> <dtm>
1     517         515      830           819     227 2013-01-01 05:00:00
2     533         529      850           830     227 2013-01-01 05:00:00
3     542         540      923           850     160 2013-01-01 05:00:00
4     544         545     1004          1022     183 2013-01-01 05:00:00
5     554         600      812           837     116 2013-01-01 06:00:00
6     554         558      740           728     150 2013-01-01 05:00:00
7     555         600      913           854     158 2013-01-01 06:00:00
8     557         600      709           723      53 2013-01-01 06:00:00
9     557         600      838           846     140 2013-01-01 06:00:00
10    558         600      753           745     138 2013-01-01 06:00:00

```

```
# i 336,766 more rows
```

## Rename variables using rename

Another useful function is `rename`, which as you may suspect renames one column to another name. Suppose we wanted `dep_time` and `arr_time` to be `departure_time` and `arrival_time` instead in the `flights_time` data frame:

```
flights_time <- flights |>
  select(contains("time")) |>
  rename(departure_time = dep_time,
         arrival_time = arr_time)
```

```
[1] "departure_time" "sched_dep_time" "arrival_time"   "sched_arr_time"
[5] "air_time"       "time_hour"
```

Note that in this case we used a single `=` sign with `rename`. eg. `departure_time = dep_time`. This is because we are not testing for equality like we would using `==`, but instead we want to assign a new variable `departure_time` to have the same values as `dep_time` and then delete the variable `dep_time`.

## Find the top number of values using slice

We can also use the `slice_max()` function which automatically tells us the most frequent `num_flights`. We specify the top 10 airports here:

```
named_dests |>
  slice_max(num_flights, n = 10)
```

# A tibble: 10 x 9

	dest	num_flights	airport_name	lat	lon	alt	tz	dst	tzone
	<chr>	<int>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<chr>	<chr>
1	ORD	17283	Chicago Ohare Intl	42.0	-87.9	668	-6	A	Amer~
2	ATL	17215	Hartsfield Jackson At~	33.6	-84.4	1026	-5	A	Amer~
3	LAX	16174	Los Angeles Intl	33.9	-118.	126	-8	A	Amer~
4	BOS	15508	General Edward Lawren~	42.4	-71.0	19	-5	A	Amer~
5	MCO	14082	Orlando Intl	28.4	-81.3	96	-5	A	Amer~
6	CLT	14064	Charlotte Douglas Intl	35.2	-80.9	748	-5	A	Amer~
7	SFO	13331	San Francisco Intl	37.6	-122.	13	-8	A	Amer~
8	FLL	12055	Fort Lauderdale Holly~	26.1	-80.2	9	-5	A	Amer~
9	MIA	11728	Miami Intl	25.8	-80.3	8	-5	A	Amer~
10	DCA	9705	Ronald Reagan Washing~	38.9	-77.0	15	-5	A	Amer~

We can find the most frequent flights in a single pipeline as follows:

```
ten_freq_dests <- flights |>
  summarize(num_flights = n(),
            .by = dest) |>
  slice_max(num_flights, n = 10)
```

```
# A tibble: 10 x 2
  dest   num_flights
  <chr>     <int>
1 ORD       17283
2 ATL       17215
3 LAX       16174
4 BOS       15508
5 MCO       14082
6 CLT       14064
7 SFO       13331
8 FLL       12055
9 MIA       11728
10 DCA        9705
```

#### Task

How could one use `starts_with`, `ends_with`, and `contains` to select columns from the `flights` data frame? Provide three different examples in total: one for `starts_with`, one for `ends_with`, and one for `contains`.

[Click here to see the solution](#)

```
# Select arrival time and arrival delay columns
flights |>
  select(starts_with("arr"))
```

```
# A tibble: 336,776 x 2
  arr_time arr_delay
  <int>     <dbl>
1     830         11
2     850         20
3     923         33
4    1004        -18
5     812        -25
6     740         12
7     913         19
```

```

8      709      -14
9      838      -8
10     753       8
# i 336,766 more rows

# Select departure and arrival delay columns
flights |>
  select(ends_with("delay"))

# A tibble: 336,776 x 2
  dep_delay arr_delay
    <dbl>      <dbl>
1         2         11
2         4         20
3         2         33
4        -1        -18
5        -6        -25
6        -4         12
7        -5         19
8        -3        -14
9        -3         -8
10       -2          8
# i 336,766 more rows

# Select departure times, schedule departure and departure delay columns
flights |>
  select(contains("dep"))

# A tibble: 336,776 x 3
  dep_time sched_dep_time dep_delay
    <int>         <int>      <dbl>
1     517           515          2
2     533           529          4
3     542           540          2
4     544           545         -1
5     554           600         -6
6     554           558         -4
7     555           600         -5
8     557           600         -3
9     557           600         -3
10    558           600         -2
# i 336,766 more rows

```



## Task

Create a new data frame that shows the top 5 airports with the largest average arrival delays from NYC in 2013.

Take a hint

Compute the mean arrival delay from each destination. You can then join the resulting data set with the `airports` data which contains the airports names and search for the top 5 entries.

[Click here to see the solution](#)

```
flights|>
  summarize(mean_arr_delay = mean(arr_delay,na.rm=T),
            .by = dest) |>
  inner_join(airports, by = join_by(dest == faa)) |>
  rename(airport_name = name) |>
  slice_max(mean_arr_delay,n=5)
```

## Vectorised if-else thru case\_when

`case_when` serves as a method to streamline multiple if-else statements by vectorizing them. It allows us to assess a condition expression and make decisions accordingly. For instance, consider a scenario where we need to categorize weather conditions according to the meteorological data contained in the `weather` data set.

Let suppose that we want to categorize the temperature variable into three categories:

- **low** for temperatures  $< 39.9$
- **medium** for temperature values  $\geq 39.9$  and  $\leq 70$
- **high** for temperature values  $> 70$

We can achieve this with the following code:

```
weather |>
  mutate(
    temp_cat = case_when(
      is.na(temp) ~ NA,
      temp < 39.9 ~ "low",
      between(temp,39.9,70)~ "medium",
      .default = "large"
    )
  )
```

```
) |>
relocate(temp,temp_cat)
```

```
# A tibble: 26,115 x 16
   temp temp_cat origin  year month  day  hour  dewp humid wind_dir wind_speed
  <dbl> <chr>   <chr> <int> <int> <int> <int> <dbl> <dbl>   <dbl>   <dbl>
1  39.0 low      EWR    2013     1     1     1  26.1  59.4     270    10.4
2  39.0 low      EWR    2013     1     1     2  27.0  61.6     250     8.06
3  39.0 low      EWR    2013     1     1     3  28.0  64.4     240    11.5
4  39.9 medium    EWR    2013     1     1     4  28.0  62.2     250    12.7
5  39.0 low      EWR    2013     1     1     5  28.0  64.4     260    12.7
6  37.9 low      EWR    2013     1     1     6  28.0  67.2     240    11.5
7  39.0 low      EWR    2013     1     1     7  28.0  64.4     240    15.0
8  39.9 medium    EWR    2013     1     1     8  28.0  62.2     250    10.4
9  39.9 medium    EWR    2013     1     1     9  28.0  62.2     260    15.0
10 41    medium    EWR    2013     1     1    10  28.0  59.6     260    13.8
# i 26,105 more rows
# i 5 more variables: wind_gust <dbl>, precip <dbl>, pressure <dbl>,
#   visib <dbl>, time_hour <dtm>
```

Here we use the `mutate` command to create new variable named `temp_cat`. The `case_when` will then set to NA those values in the original `temp` variable that are missing. Then if the values of `temp` are `< 30.9` it will assign them the label of `low`. If they lie between 39.9 and 70 it will assign them the label of `medium` and finally set to `large` any of the values that do not meet any of the aforementioned conditions. We can also use the function `relocate` to change the columns position so that the `temp` and `temp_cat` appears first on the data frame.

#### Task

Create a new variable called `extreme_weather` that takes the value of `extreme` if the wind speed exceeds 64 mph and the temperature is less than 40 °F and **not** `extreme` otherwise. Then, relocate this new variable along with the variables used to create it at the first columns of the data frame, and sort them out based on `wind_speed`.

Take a hint

Use the conditional operators `|` and `&` to add multiple conditions.

[Click here to see the solution](#)

```

weather |>
  mutate(
    extreme_weather = case_when(
      is.na(temp)|is.na(wind_speed) ~ NA,
      temp < 40 & wind_speed > 64 ~ "extreme",
      .default = "not extreme"
    )
  ) |>
  relocate(extreme_weather,temp,wind_speed) |>
  arrange(desc(wind_speed))

```

```
# A tibble: 26,115 x 16
```

	extreme_weather	temp	wind_speed	origin	year	month	day	hour	dewp	humid
	<chr>	<dbl>	<dbl>	<chr>	<int>	<int>	<int>	<int>	<dbl>	<dbl>
1	extreme	39.0	1048.	EWR	2013	2	12	3	27.0	61.6
2	not extreme	57.2	42.6	EWR	2013	1	31	6	53.6	87.7
3	not extreme	53.6	42.6	JFK	2013	1	31	4	53.1	100
4	not extreme	60.8	40.3	EWR	2013	1	31	4	59	93.8
5	not extreme	59	40.3	LGA	2013	1	31	4	55.4	93.7
6	not extreme	46.0	39.1	EWR	2013	1	31	8	30.0	53.3
7	not extreme	41	38.0	JFK	2013	3	6	14	28.9	61.9
8	not extreme	53.1	36.8	JFK	2013	1	31	3	52.0	100
9	not extreme	51.8	36.8	JFK	2013	1	31	7	46.4	81.7
10	not extreme	28.0	36.8	JFK	2013	11	24	10	-0.04	29.2

```
# i 26,105 more rows
```

```
# i 6 more variables: wind_dir <dbl>, wind_gust <dbl>, precip <dbl>,
```

```
# pressure <dbl>, visib <dbl>, time_hour <dtm>
```

## Summary

The table below lists a selection of the data wrangling verbs and summarises what they do. Using these verbs and the pipe `|>` operator, you'll be able to write easily legible code to perform almost all the data wrangling necessary for the rest of this course.

Table 4: Summary of data wrangling verbs

Verb	Operation
<code>filter()</code>	Pick out a subset of rows
<code>summarize()</code>	Summarise many values to one using a summary statistic function like <code>mean()</code> , <code>median()</code> , etc.
<code>mutate()</code>	Create new variables by mutating existing ones
<code>arrange()</code>	Arrange rows of a data variable in ascending (default) or descending order
<code>inner_join()</code>	Join/merge two data frames, matching rows by a key variable
<code>select()</code>	Pick out a subset of columns to make data frames easier to view