

# Session 2 - Data wrangling

## R training

---

María Reyes Retana

The World Bank | January 2025



# Government Analytics and R Training:

## Strengthening Public Sector Reporting and Data Analysis

January 13 – January 17, 2025






# Introduction

## Access to Training Materials

All the materials for this training are available at the following link:

<https://osf.io/r3fn5/>





Here's what you will find there:

-  **Data:** All datasets we'll use throughout the training sessions.
-  **Slides:** The presentation slides for each session.
-  **Solutions:** Will be added to the folder after each session.

# Introduction

## Course Structure

This training will cover the basics of coding in R. Below is the structure of the course:

-  **Day 1:** Introduction to R – Get familiar with the R environment and basic syntax.
-  **Day 2: Data Wrangling – Learn to clean and transform data effectively.**
-  **Day 3:** Descriptive Statistics – Explore methods to summarize and analyze data.
-  **Day 4:** Data Visualization – Create meaningful and impactful data visualizations.

# About this session

---

# About this session

Input



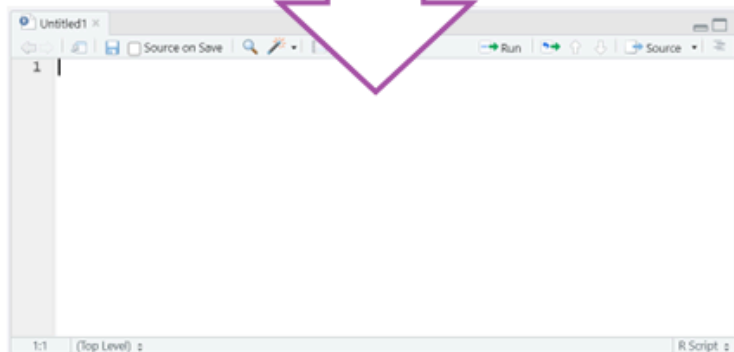
department\_staff\_list.csv

Data process

Instructions

1. Read data
2. Select columns
3. Obtain averages for all columns
4. Produce graph
5. Export graph in image

Output



# About this session

Input



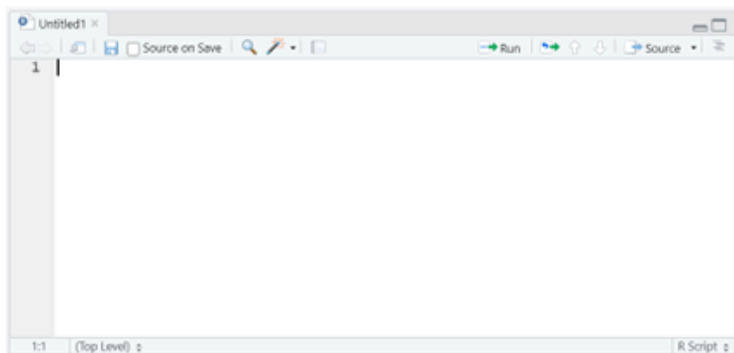
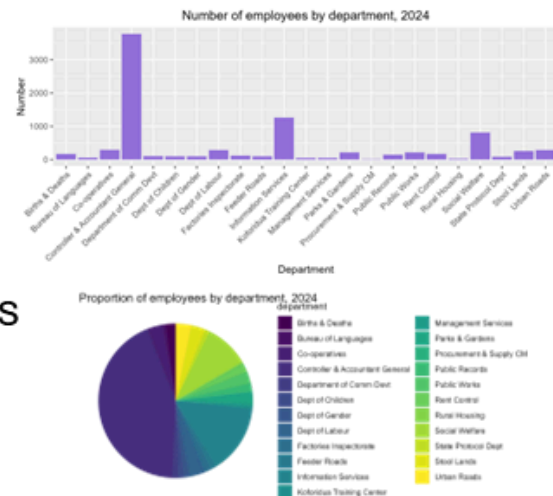
department\_staff\_list.csv

Data process

Instructions

1. Read data
2. Select columns
3. Obtain averages for all columns
4. Produce graph
5. Export graph in image

Output



This session: getting data ready to be used to create our outputs

# R Packages

---

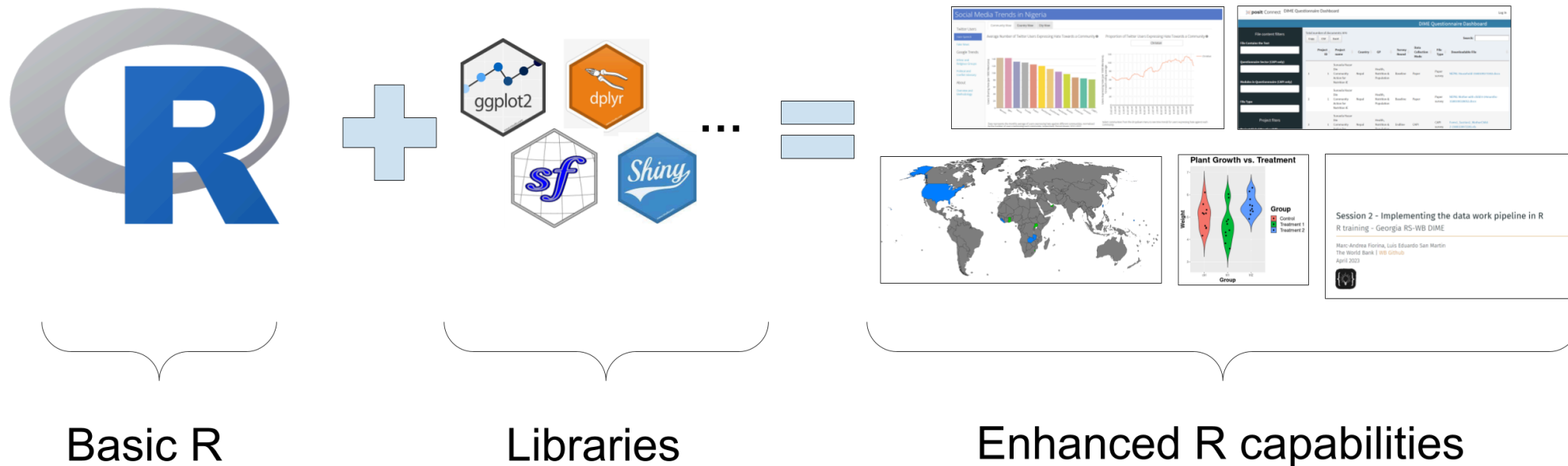


# R Packages

- Installing R in your computer gives you access to its basic functions
- Additionally, you can also install packages. Packages are a collection of R functions that allow you to do:
  - Operations that basic R functions don't do (example: work with geographic data)
  - Operations that basic R functions do, but easier (example: data wrangling)
- They contain code that other R users have prepared for the community.

# R Packages

In a nutshell:

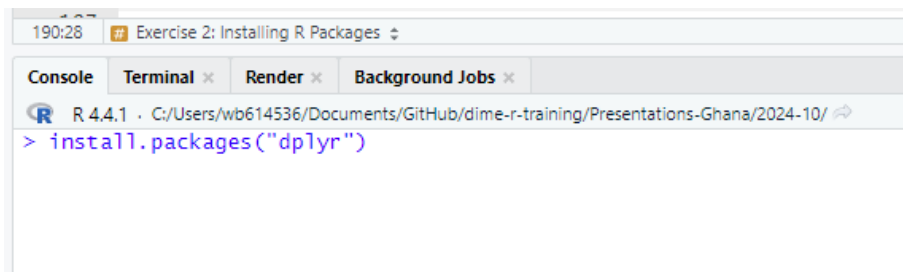


# R Packages

We'll use one package in today's session: `dplyr` and `tidyr` (really useful library for data cleaning and wrangling).

## Exercise 1: Installing R Packages

1. Install the R Packages by using `install.packages()`
  - `install.packages("dplyr")`
  - Note the quotes ( " ") in the packages names
  - **Introduce this code in the console**, not the script panel

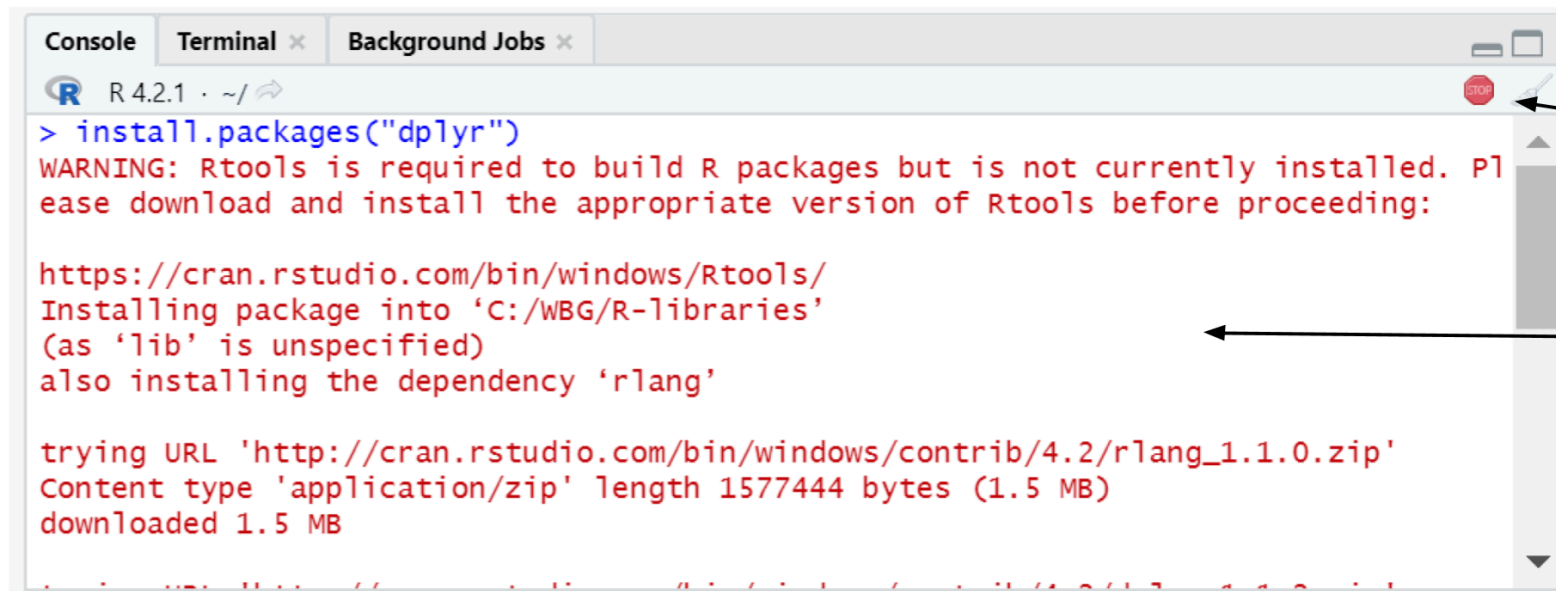


The screenshot shows an R console window with the following elements:

- Top bar: "190:28" and "Exercise 2: Installing R Packages"
- Tab bar: "Console", "Terminal", "Render", "Background Jobs"
- Console output: "R 4.4.1 · C:/Users/wb614536/Documents/GitHub/dime-r-training/Presentations-Ghana/2024-10/"
- Command prompt: "> install.packages('dplyr')"

# R packages

## Installing packages



```
R 4.2.1 · ~/
> install.packages("dplyr")
WARNING: Rtools is required to build R packages but is not currently installed. Please
download and install the appropriate version of Rtools before proceeding:

https://cran.rstudio.com/bin/windows/Rtools/
Installing package into 'C:/WBG/R-libraries'
(as 'lib' is unspecified)
also installing the dependency 'rlang'

trying URL 'http://cran.rstudio.com/bin/windows/contrib/4.2/rlang_1.1.0.zip'
Content type 'application/zip' length 1577444 bytes (1.5 MB)
downloaded 1.5 MB
```

The “STOP” sign means that the code is still running, just wait until it finishes

Note that this message is **not an error**

# R packages

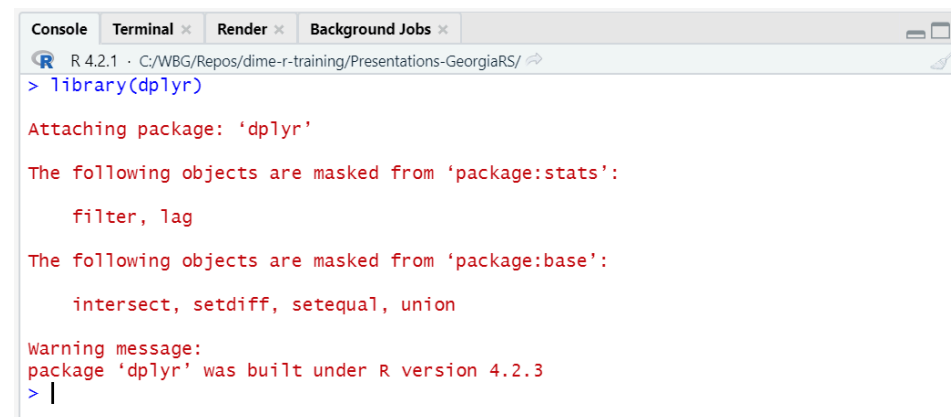
Now that `dplyr` is installed, we only need to load them to start using the functions they have.

## Exercise 2: Loading packages

1. Open a new script with **File** >> **New File** >> **R Script**
2. Load the packages with:

```
library(dplyr)
```

- Run this code from the new script you just opened
- Notice that we don't use quotes in the package names this time



```
Console Terminal x Render x Background Jobs x
R 4.2.1 · C:/WBG/Repos/dime-r-training/Presentations-GeorgiaRS/
> library(dplyr)

Attaching package: 'dplyr'

The following objects are masked from 'package:stats':

  filter, lag

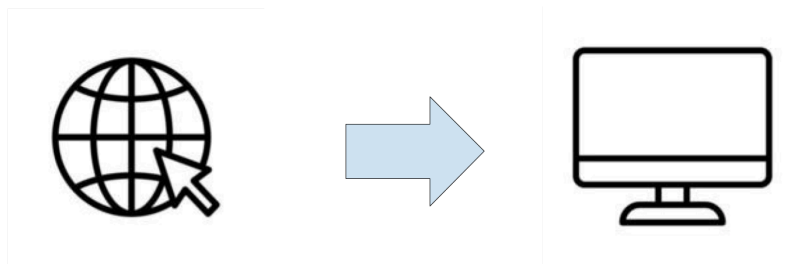
The following objects are masked from 'package:base':

  intersect, setdiff, setequal, union

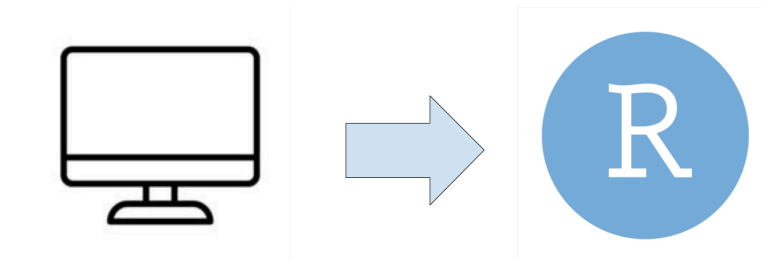
Warning message:
package 'dplyr' was built under R version 4.2.3
> |
```

# R packages

- Library installation:



- Library loading:



- You install R packages only once in your computer
- You load packages every time you open a new RStudio window (only load the packages you will use)

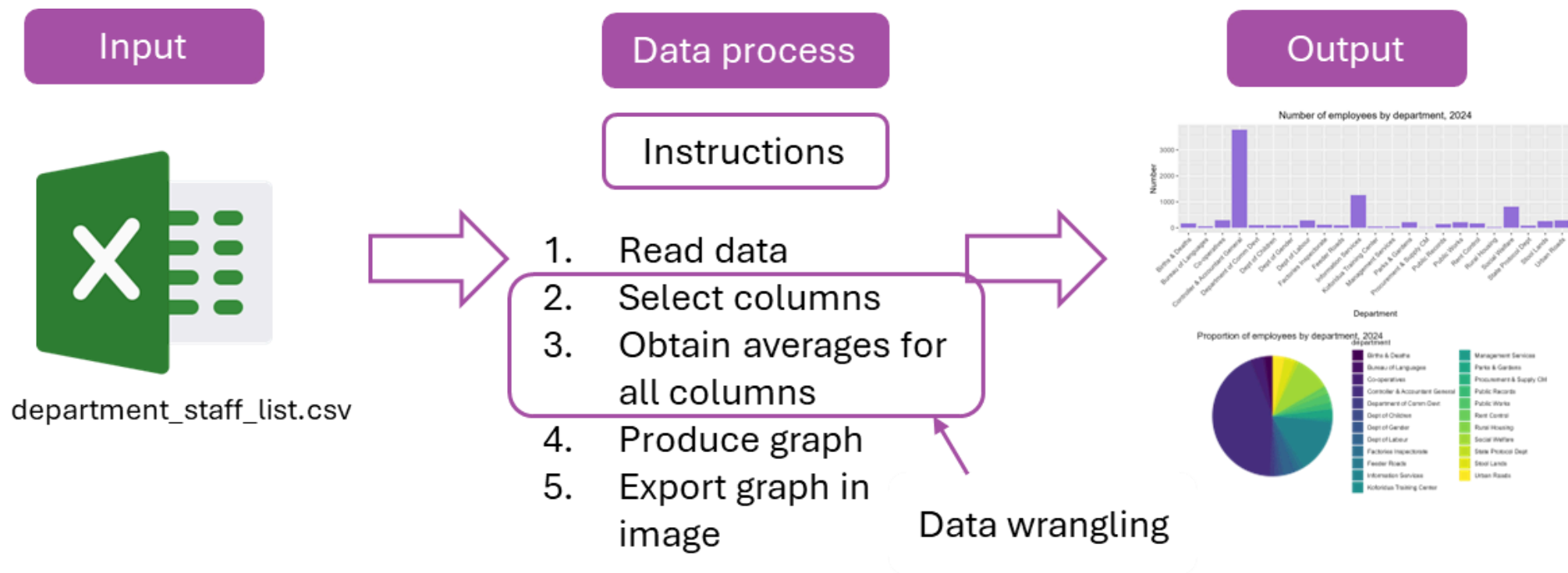
# Data wrangling

---

# Data wrangling

## Getting your data ready

- Data is rarely in a format where it can be converted in an output right away
- In statistical programming, the process of transforming data into a condition where it's ready to be converted into an output is called **data wrangling**







# Data wrangling

## Getting your data ready

- As we said before we'll use `dplyr` for data wrangling in this training
- You can also use basic R, but we recommend these packages because its functions are easier to use



# Data wrangling

## Exercise 3: Loading data

Note that this part of this is the same exercise we did in session 1, but it's okay to repeat it in order to start using a new RStudio session. **If you have RStudio open, start by closing the window and opening RStudio again.**

1. In your new RStudio window, go to **File** > **Import Dataset** > **From Text(base)** and select again the file **department\_staff\_list.csv**

- if you don't know where the file is, check in the **Downloads** folder
- if you need to download it again, it's here:

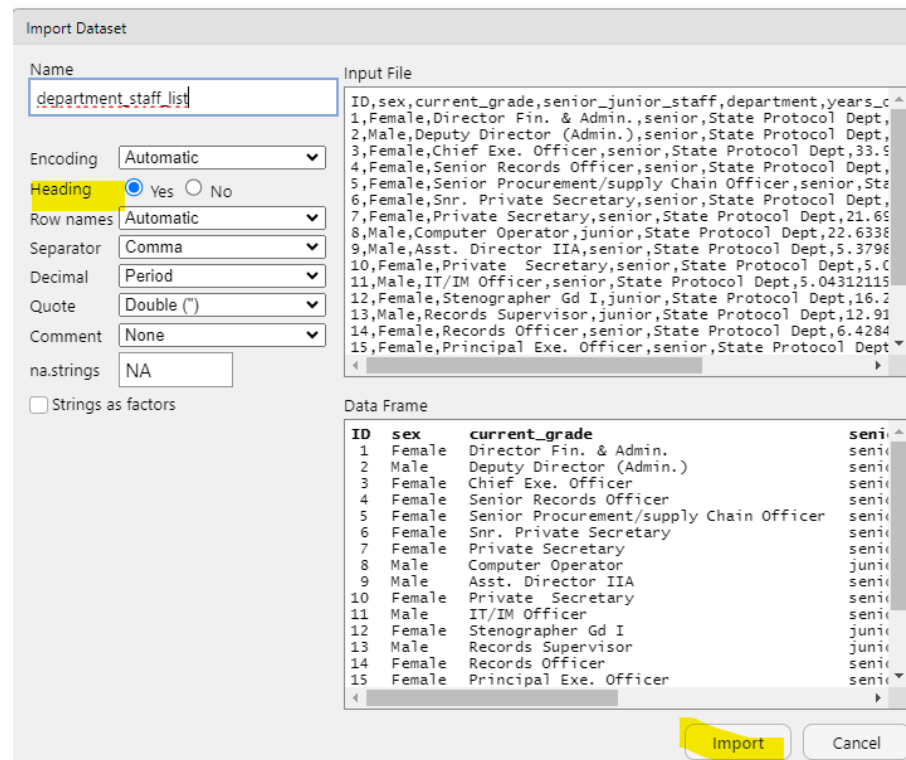
<https://osf.io/chdgj>

2. Make sure to select **Heading** > **Yes** in the next window

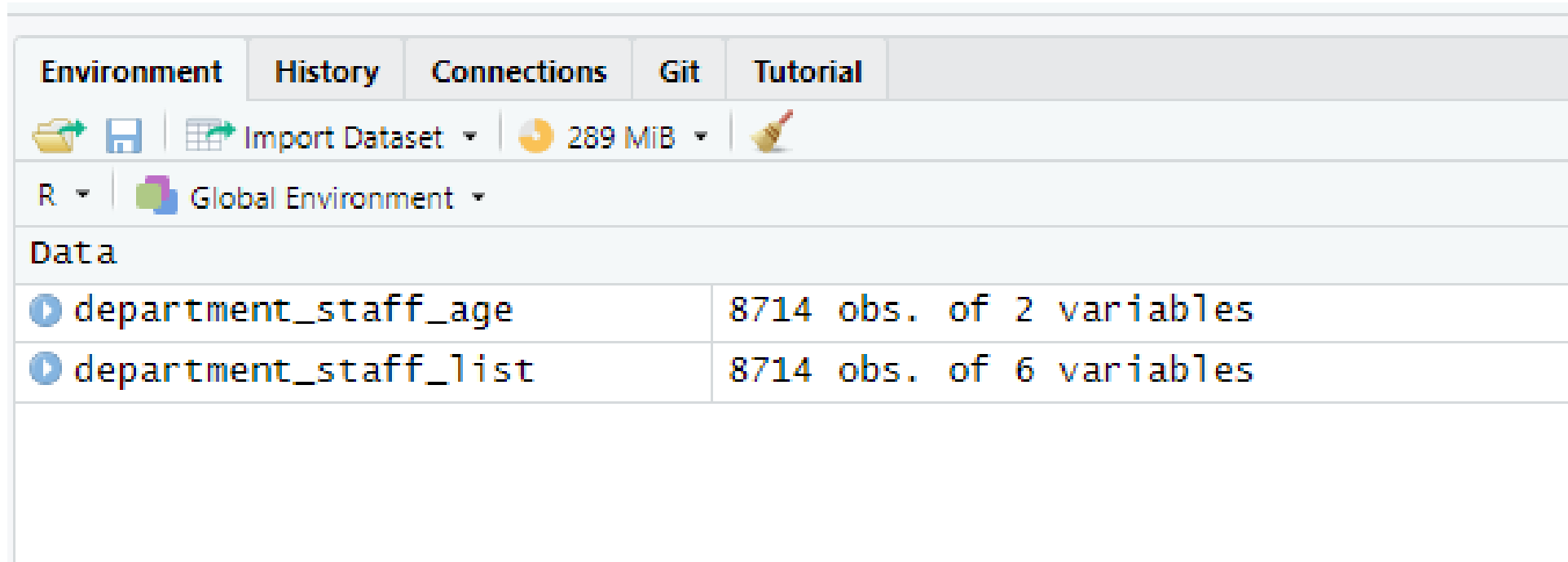
3. Select **Import**

4. Download this new file **department\_staff\_age.csv**:

<https://osf.io/qsmke> and repeat steps 1-3 with it



# Data wrangling



The screenshot shows the RStudio Environment pane. At the top, there are tabs for 'Environment', 'History', 'Connections', 'Git', and 'Tutorial'. Below the tabs is a toolbar with icons for file operations and a memory usage indicator showing '289 MiB'. The main area of the pane is titled 'Data' and contains a table of objects in the 'Global Environment'.

Global Environment	
Data	
▶ department_staff_age	8714 obs. of 2 variables
▶ department_staff_list	8714 obs. of 6 variables

# Data wrangling

## Note: loading data with a function

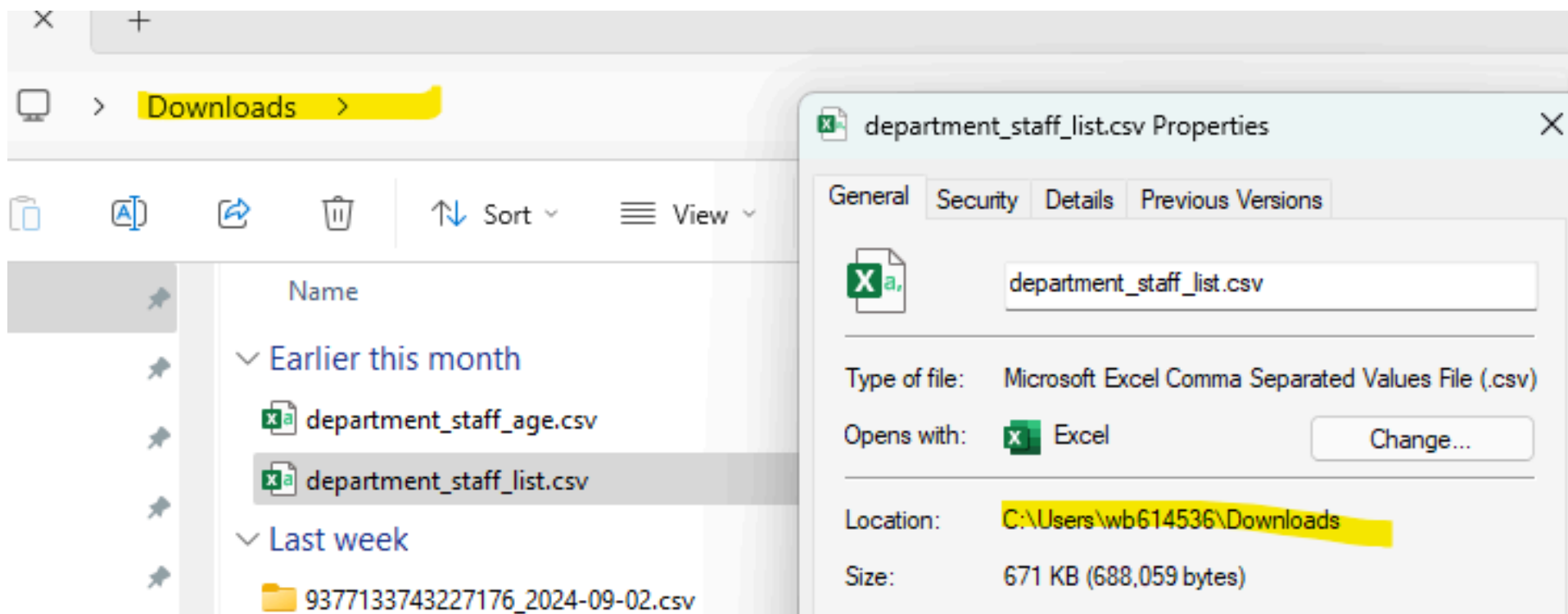
- You can also load csv files with the function `read.csv()` instead of using this point-and-click approach
- The **argument** of `read.csv()` is the path in your computer where your data is. For example

```
department_staff_list <- read.csv("C:/Users/wb614536/Downloads/department_staff_list.csv")
```

- As usual, you need to save the result of `read.csv()` into a dataframe object with the arrow operator (`<-`) for it to be stored in the environment

# Note on file paths

- A **file path** tells R where to find your file on your computer. It is like giving R the directions to your file.
- If you downloaded the data and haven't moved it, the data path for the department\_list dataset will be probably something like: `"C:/Users/wb614536/Downloads/department_staff_list.csv"`
- You can find the path of a file by right-clicking>properties>and seeing the location, or by clicking on the top bar (see below)



# Data wrangling

## Recap: knowing your data

- Dataframe `department_staff_list` is the same dataframe we used last session that contains the data from your department.

```
glimpse(department_staff_list)
```

```
## Rows: 8,713
## Columns: 6
## $ ID          <int> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15,...
## $ sex         <chr> "Female", "Male", "Female", "Female", "Female", "F...
## $ current_grade <chr> "Assist. Landscape Designer Grade II", "Co-Operati...
## $ senior_junior_staff <chr> "junior", "senior", "senior", "senior", "senior", ...
## $ department  <chr> "Parks & Gardens", "Co-operatives", "Controller & ...
## $ years_of_service <dbl> 24.87, 17.83, 25.65, 2.76, 12.56, 23.75, 4.59, 3.2...
```

# Data wrangling

## Recap: knowing your data

- Now we are introducing a second dataset that only has the ID number (one that I invented), and the age of that person.

```
glimpse(department_staff_age)
```

```
## Rows: 8,714
```

```
## Columns: 2
```

```
## $ ID   <int> 8603, 3695, 5563, 3529, 5032, 4391, 3069, 6048, 4680, 8551, 5107, ...
```

```
## $ age  <dbl> 18.00442, 18.00657, 18.00669, 18.00874, 18.01030, 18.01677, 18.018...
```



# Data wrangling

- We will only use this second dataframe in one of the next exercises, but we load it now because it's in general a good practice to have data loaded into the memory so it's ready to be used.
- For the next exercises, we will propose scenarios that could show up while doing your annual reports or in day-to-day operations.
- We will do everything using **functions** from the `dplyr` package, that is already in our environment.

# Filtering and sorting

---

# Filtering and sorting

`dplyr::filter()` KEEP ROWS THAT  
satisfy  
your **CONDITIONS**

keep rows from... this data... ONLY IF... type is "otter" AND site is "bay"  
`filter(df, type == "otter" & site == "bay")`



type	food	site
otter	urchin	bay
shark	seal	channel
otter	abalone	bay
otter	crab	wharf



The table shows the results of the filter operation. The first and third rows are kept (indicated by checkmarks), while the second and fourth rows are filtered out (indicated by X marks).

# Filtering and sorting

## Government analytics request

We will use the data from our dataframe `department_staff_list`

For this we can filter the Female rows from the sex category

exercises-session2.R* x		department_staff_list x	
Filter			
	ID	sex	current_grade
1	1	Female	Director Fin. & Admin.
2	2	Male	Deputy Director (Admin.)
3	3	Female	Chief Exe. Officer
4	4	Female	Senior Records Officer
5	5	Female	Senior Procurement/supply
6	6	Female	Snr. Private Secretary
7	7	Female	Private Secretary
8	8	Male	Computer Operator
9	9	Male	Asst. Director IIA
10	10	Female	Private Secretary
11	11	Male	IT/IM Officer
12	12	Female	Stenographer Gd I



Filter			
	ID	sex	current_grade
1	1	Female	Director Fin. & Admin.
2	3	Female	Chief Exe. Officer
3	4	Female	Senior Records Officer
4	5	Female	Senior Procurement/supply Chain Officer
5	6	Female	Snr. Private Secretary
6	7	Female	Private Secretary
7	10	Female	Private Secretary
8	12	Female	Stenographer Gd I
9	14	Female	Records Officer
10	15	Female	Principal Exe. Officer
11	16	Female	Steno Secretary
12	17	Female	senior Technical Asst.
13	22	Female	Principal Protocol Officer

# Filtering and sorting

## Government analytics request

Now let's say that we are also interested at a first glance of the females that recently joined the department. We would have to do the following

1. Keeping only the female employees
2. Sorting by years of service

✕ In Excel: We would filter and then use arrange by years of service.

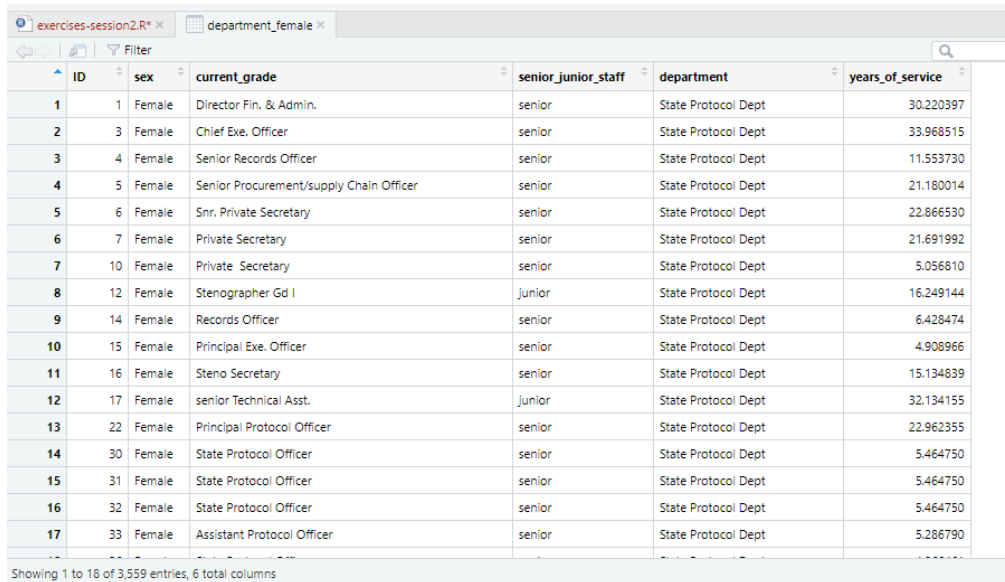
# Filtering and sorting

## 1. Keeping only the female employees

Use `filter()` for this:

Remember how we use functions in this case 1st argument is **data**, second argument is **filter** (sex takes the value female in this case)

```
temp1 <- filter(department_staff_list, sex == "Female")
```



	ID	sex	current_grade	senior_junior_staff	department	years_of_service
1	1	Female	Director Fin. & Admin.	senior	State Protocol Dept	30.220397
2	3	Female	Chief Exe. Officer	senior	State Protocol Dept	33.968515
3	4	Female	Senior Records Officer	senior	State Protocol Dept	11.553730
4	5	Female	Senior Procurement/supply Chain Officer	senior	State Protocol Dept	21.180014
5	6	Female	Snr. Private Secretary	senior	State Protocol Dept	22.866530
6	7	Female	Private Secretary	senior	State Protocol Dept	21.691992
7	10	Female	Private Secretary	senior	State Protocol Dept	5.056810
8	12	Female	Stenographer Gd I	junior	State Protocol Dept	16.249144
9	14	Female	Records Officer	senior	State Protocol Dept	6.428474
10	15	Female	Principal Exe. Officer	senior	State Protocol Dept	4.908966
11	16	Female	Steno Secretary	senior	State Protocol Dept	15.134839
12	17	Female	senior Technical Asst.	junior	State Protocol Dept	32.134155
13	22	Female	Principal Protocol Officer	senior	State Protocol Dept	22.962355
14	30	Female	State Protocol Officer	senior	State Protocol Dept	5.464750
15	31	Female	State Protocol Officer	senior	State Protocol Dept	5.464750
16	32	Female	State Protocol Officer	senior	State Protocol Dept	5.464750
17	33	Female	Assistant Protocol Officer	senior	State Protocol Dept	5.286790

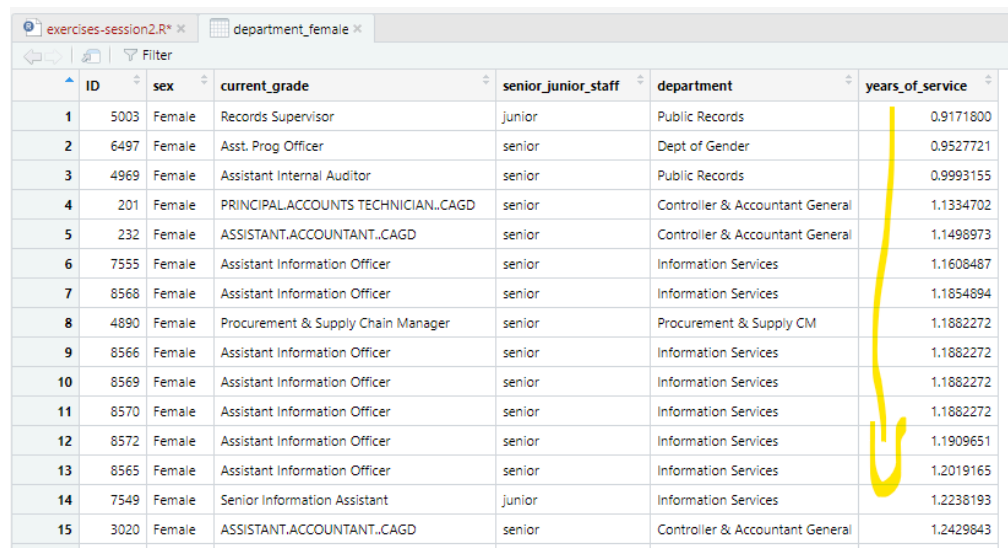
Showing 1 to 18 of 3,559 entries, 6 total columns

# Filtering and sorting

## 2. Sorting by years of service

Use the function `arrange()` to sort. Sortings are ascending by default in R (this will get us from less years of service to more years of service).

```
department_female <- arrange(temp1, years_of_service) # arrange by years of service
```



The screenshot shows an RStudio window with two tabs: 'exercises-session2.R\*' and 'department\_female'. The 'department\_female' tab is active, displaying a data table with 7 columns: ID, sex, current\_grade, senior\_junior\_staff, department, and years\_of\_service. The table is sorted by 'years\_of\_service' in ascending order. A yellow arrow points to the 'years\_of\_service' column.

	ID	sex	current_grade	senior_junior_staff	department	years_of_service
1	5003	Female	Records Supervisor	junior	Public Records	0.9171800
2	6497	Female	Asst. Prog Officer	senior	Dept of Gender	0.9527721
3	4969	Female	Assistant Internal Auditor	senior	Public Records	0.9993155
4	201	Female	PRINCIPAL.ACCOUNTS TECHNICIAN.,CAGD	senior	Controller & Accountant General	1.1334702
5	232	Female	ASSISTANT.ACCOUNTANT.,CAGD	senior	Controller & Accountant General	1.1498973
6	7555	Female	Assistant Information Officer	senior	Information Services	1.1608487
7	8568	Female	Assistant Information Officer	senior	Information Services	1.1854894
8	4890	Female	Procurement & Supply Chain Manager	senior	Procurement & Supply CM	1.1882272
9	8566	Female	Assistant Information Officer	senior	Information Services	1.1882272
10	8569	Female	Assistant Information Officer	senior	Information Services	1.1882272
11	8570	Female	Assistant Information Officer	senior	Information Services	1.1882272
12	8572	Female	Assistant Information Officer	senior	Information Services	1.1909651
13	8565	Female	Assistant Information Officer	senior	Information Services	1.2019165
14	7549	Female	Senior Information Assistant	junior	Information Services	1.2238193
15	3020	Female	ASSISTANT.ACCOUNTANT.,CAGD	senior	Controller & Accountant General	1.2429843

# Filtering and sorting

## Exercise 4: Now let's do it.. filter and sort your data

We can write the whole code for this in our exercise script. (You can copy and paste the code below)

1.- Filter by `Female`:

2.- Sort by `years_of_service`:

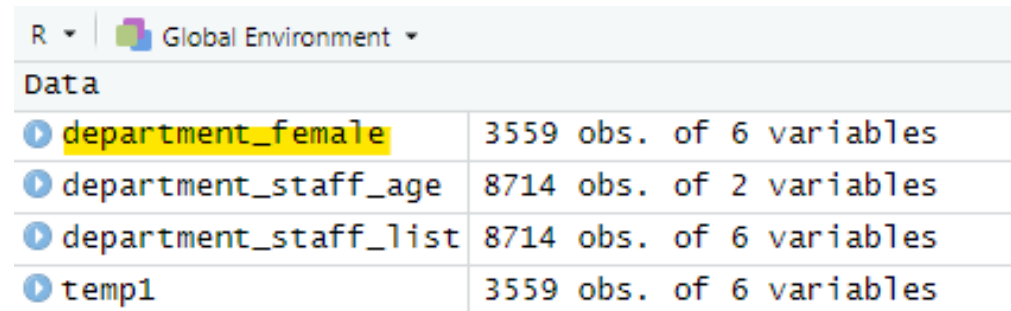
```
temp1 <- filter(department_staff_list, sex == "Female") # filter by female  
department_female <- arrange(temp1, years_of_service) # order by years of service
```



# Filtering and sorting

Some notes:

- `filter()` and `arrange()` are all functions from `dplyr`. Remember you have to always load `dplyr` first with `library(dplyr)` to be able to use them
- The resulting dataframe is `department_female` (and we also have `temp1`)



The screenshot shows the 'Global Environment' pane in R Studio. Under the 'Data' section, four objects are listed: 'department\_female', 'department\_staff\_age', 'department\_staff\_list', and 'temp1'. Each object has a blue play button icon to its left. The details for each object are: 'department\_female' (3559 obs. of 6 variables), 'department\_staff\_age' (8714 obs. of 2 variables), 'department\_staff\_list' (8714 obs. of 6 variables), and 'temp1' (3559 obs. of 6 variables).

Data	
▶ department_female	3559 obs. of 6 variables
▶ department_staff_age	8714 obs. of 2 variables
▶ department_staff_list	8714 obs. of 6 variables
▶ temp1	3559 obs. of 6 variables

# Filtering and sorting

- Filtering and sorting are two very common data wrangling operations in statistical programming
- Now we'll review a new data wrangling operation that is also quite common and useful: **mutate**, which basically means creating new variables

# Mutate

---

# Mutate

- `mutate` will take a statement like this:

```
mutate(variable_name = some_calculation)
```

- And attach variable\_name at the end of the dataset.



# Mutate

## Example

Using our data frame let's say that we want to include a variable that instead of years of service we want days of service. We would do this by doing the following:

```
example_mutate <- mutate(department_staff_list, days_of_service = years_of_service*365)
```

We will not use that variable for our government analytics examples, but this is a really useful data wrangling function.

# Questions?

# Merging dataframes

---

# Merging dataframes

Merging data is a common task in data analysis, especially when working with data from multiple departments.

Let's see how it would apply to our dataframes.

## Government analytics request

### **Scenario 2:**

*Let's imagine that for our annual report we are also interested in the age distribution from the employees* but these are in different datasets



# Merging dataframes

## Government analytics request

Use the data `department_staff_list` that you already know with the `department_staff_age`

We want to include the `age` variable to our dataframe

ID	sex	current_grade	senior_junior_staff	department	years_of_service
1	Female	Director Fin. & Admin.	senior	State Protocol Dept	30.220397
2	Male	Deputy Director (Admin.)	senior	State Protocol Dept	16.049281
3	Female	Chief Exe. Officer	senior	State Protocol Dept	33.968515
4	Female	Senior Records Officer	senior	State Protocol Dept	11.553730
5	Female	Senior Procurement/supply Chain Officer	senior	State Protocol Dept	21.180014
6	Female	Snr. Private Secretary	senior	State Protocol Dept	22.866530
7	Female	Private Secretary	senior	State Protocol Dept	21.691992
8	Male	Computer Operator	junior	State Protocol Dept	22.633812
9	Male	Asst. Director IIA	senior	State Protocol Dept	5.379877
10	Female	Private Secretary	senior	State Protocol Dept	5.056810
11	Male	IT/IM Officer	senior	State Protocol Dept	5.043121
12	Female	Stenographer Gd I	junior	State Protocol Dept	16.249144
13	Male	Records Supervisor	junior	State Protocol Dept	12.919918
14	Female	Records Officer	senior	State Protocol Dept	6.428474
15	Female	Principal Exe. Officer	senior	State Protocol Dept	4.908966
16	Female	Steno Secretary	senior	State Protocol Dept	15.134839



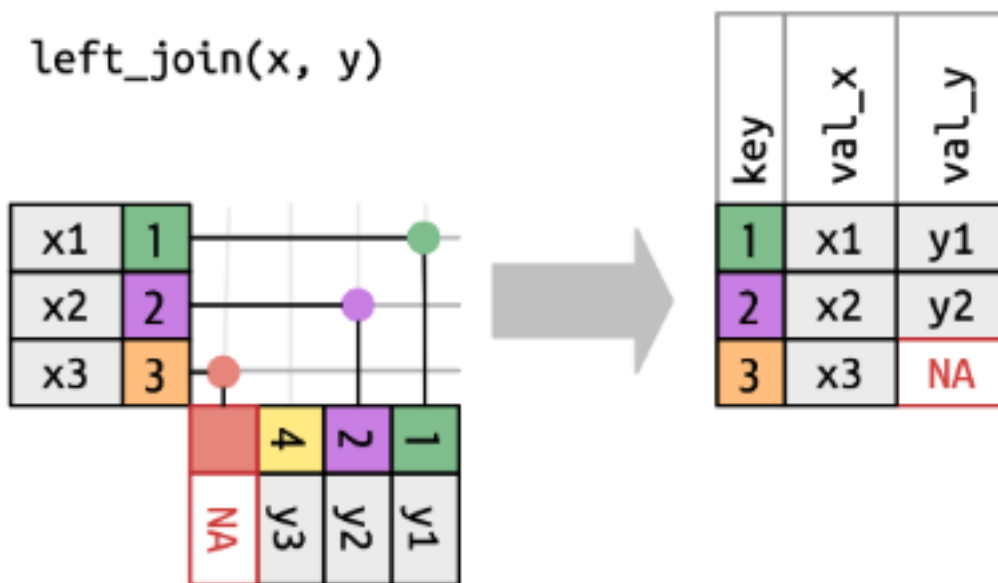
ID	age
1	60.15058
2	41.74949
3	59.43053
4	50.55989
5	51.28268
6	45.22656
7	41.90007
8	56.56400
9	29.09788
10	38.76523
11	32.14784
12	38.01506
13	34.16564
14	43.42505
15	44.99932
16	43.50171

# Merging dataframes

To do this we will use `left_join()` to merge the dataframes:

- The arguments of the function are 1. the "principal" dataset and 2. the dataframe we want to merge (paste) into that one

```
department_age <- left_join(department_staff_list, department_staff_age) # Our original data frame
```



# Merging dataframes

## Exercise 5: Now let's do it.

- You can copy and paste the code to your exercise script.

```
deparment_age <- left_join(deparment_staff_list, department_staff_age) # Our original data frame
```

The screenshot displays the RStudio environment. The main window shows a data frame with the following columns: **current\_grade**, **senior\_junior\_staff**, **department**, **years\_of\_service**, and **age**. The data includes rows for various roles such as Director Fin. & Admin., Deputy Director (Admin.), Chief Exe. Officer, Senior Records Officer, Senior Procurement/supply Chain Officer, Snr. Private Secretary, Private Secretary, Computer Operator, Asst. Director IIA, Private Secretary, IT/IM Officer, Stenographer Gd I, and Records Supervisor.

The Environment pane on the right lists the following objects:

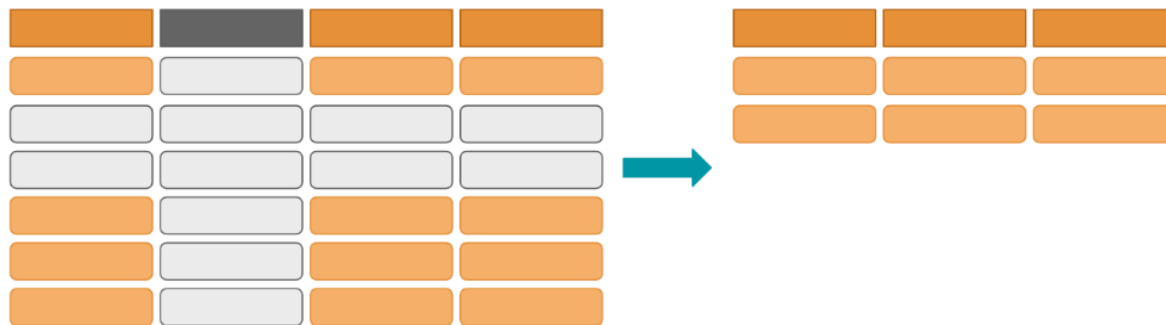
Object	Observations	Variables
deparment_age	8714 obs.	7 variables
department_female	3559 obs.	6 variables
department_staff_age	8714 obs.	2 variables
department_staff_list	8714 obs.	6 variables

# group\_by() and summarize()

We use this when we want to aggregate your data (by groups).

This is one of the most common operations. Sometimes we want to calculate statistics by groups

Customize with **group\_by()** and **summarize()**



# Grouping and summarizing

## Government analytics request

**Let's say that we are interested again in the number of employees by department**

We will use our `department_staff_list` dataset to do this.

1. The first step is to define the group, in this case `department`

```
employees_by_department <- group_by(department_staff_list, department)
```

If we **only** do this, this won't do anything, to complete the function we need to use this with `summarise()`

# Grouping and summarizing

## Exercise 6: Government analytics request

summarize works in a similar way to mutate:

```
variable_name = some_calculation
```

In this case the `some_calculation` will be to count the number of employees

```
temp1 <- group_by(department_staff_list, department)
employees_by_department <- summarise(temp1, number = n())
```

# Grouping and summarizing

## Government analytics request

This will create the following dataframe/table:

```
## # A tibble: 23 × 2
##   department          number
##   <chr>              <int>
## 1 Births & Deaths      172
## 2 Bureau of Languages    58
## 3 Co-operatives        293
## 4 Controller & Accountant General 3776
## 5 Department of Comm Devt   109
## 6 Dept of Children       101
## 7 Dept of Gender         100
## 8 Dept of Labour         290
## 9 Factories Inspectorate  121
## 10 Feeder Roads          102
## # i 13 more rows
```

# More wrangling operations

These were two examples we chose to show different possible data wrangling operations. A summary of these and other common operations are:

Operation	Function in <code>dplyr</code>
Subset columns	<code>select()</code>
Subset rows (based on condition)	<code>filter()</code>
Create new columns	<code>mutate()</code>
Create new columns based on condition	<code>mutate()</code> and <code>case_when()</code>
Create new rows	<code>add_row()</code>
Merge dataframes	<code>inner_join()</code> , <code>left_join()</code> , <code>right_join()</code> , <code>full_join()</code>
Append dataframes	<code>bind_rows()</code>
Deduplicate	<code>distinct()</code>
Collapse and create summary indicators	<code>group_by()</code> , <code>summarize()</code>
Pass a result as the first argument for the next function	<code>%&gt;%</code> (operator, not function ( <b>tomorrow</b> ))

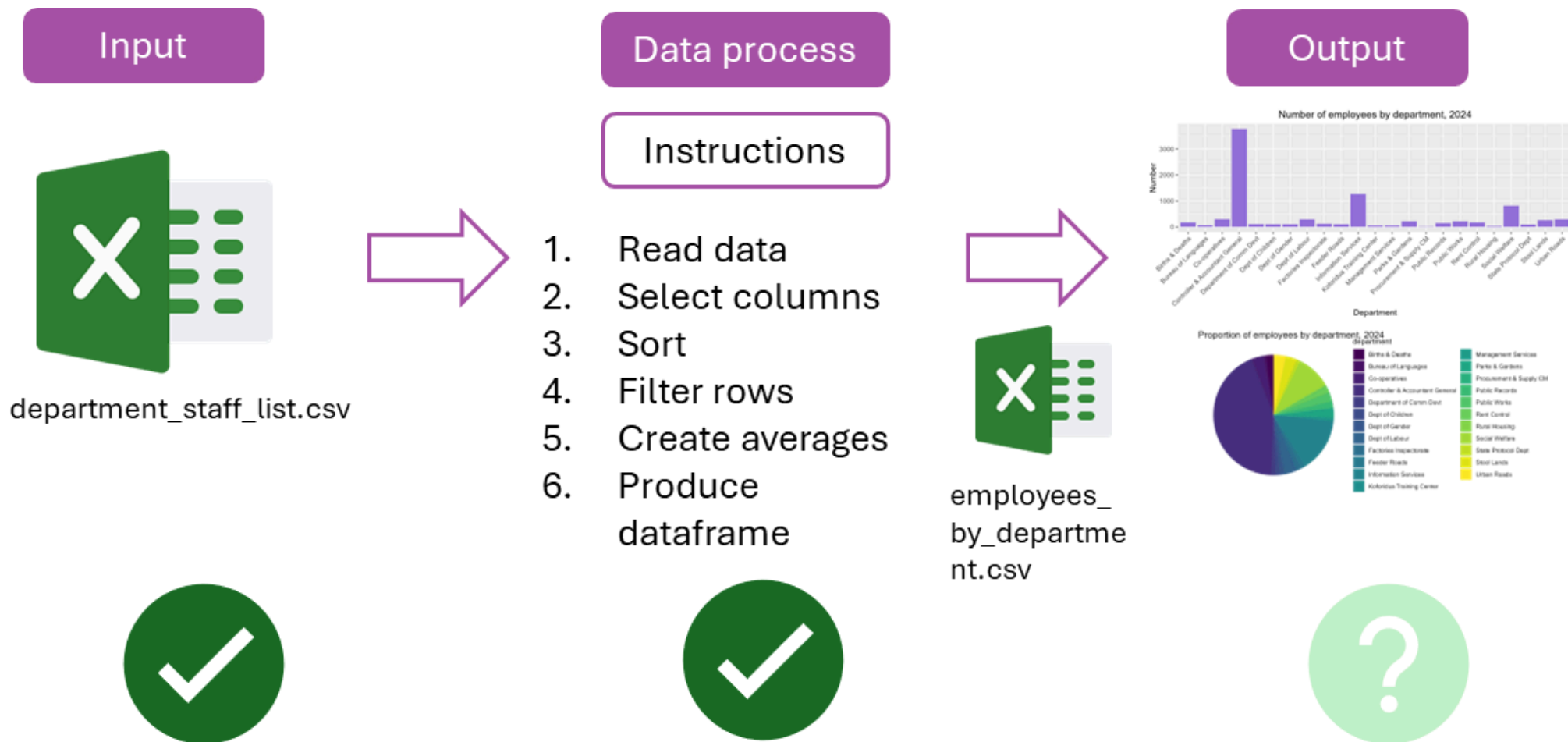


# Exporting outputs

---

# Exporting outputs

- Until now, we've seen full examples of part 1 and 2 of the Government analytics pipeline
- What about exporting outputs?



# Exporting outputs

## Exporting dataframes

- We can export it as a csv file with the function `write.csv()`
- `write.csv()` creates a csv file with the dataframe
- It takes two basic arguments:
  1. The name of the object you want to export
  2. A file path to export the object to

# Exporting outputs

## Exercise 7 (if time allows): Export `employees_by_department` and `department_staff_final`

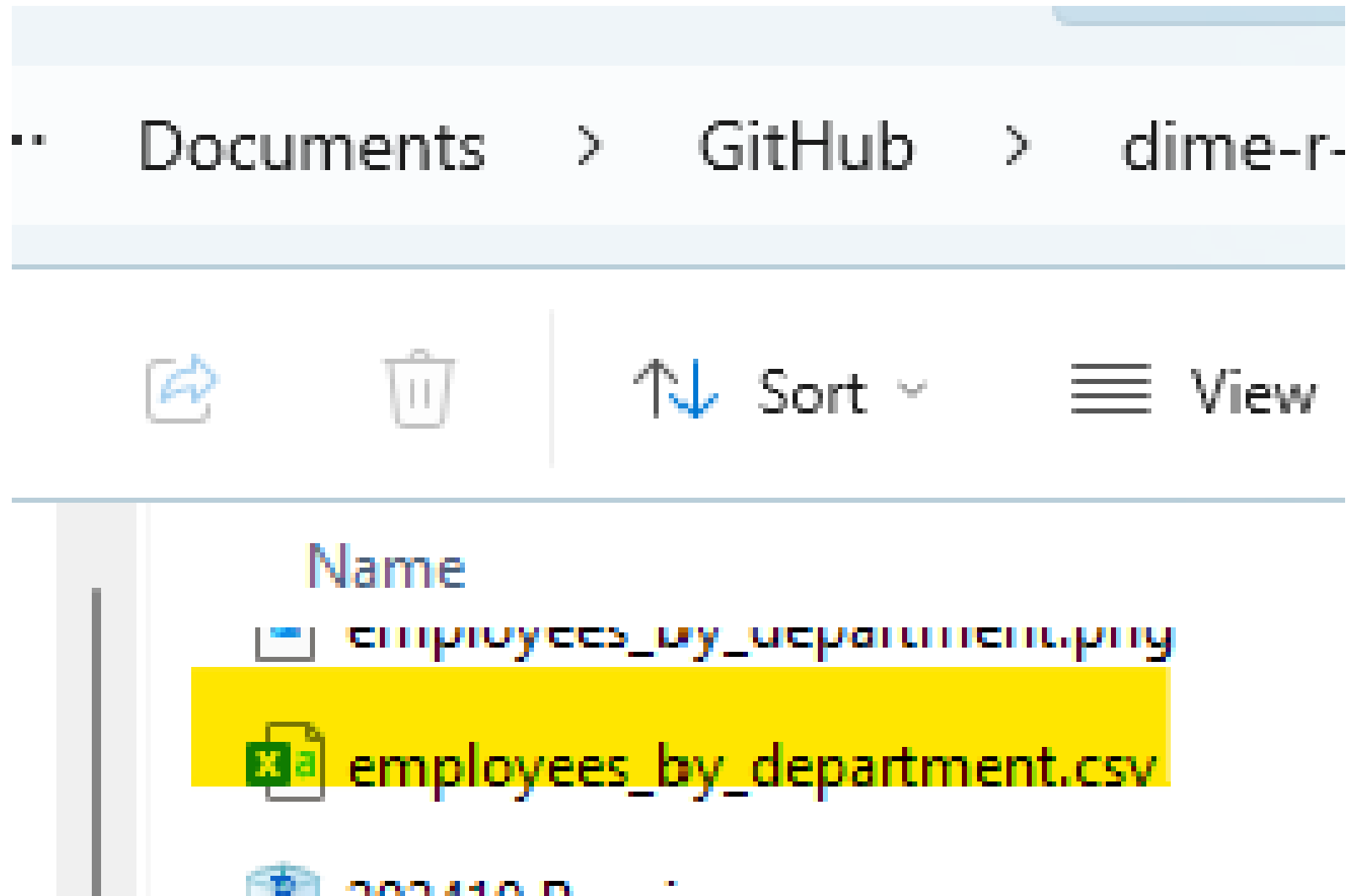
1. Use this code to export the results of the last two exercises:

```
write.csv(employees_by_department, "employees_by_department", row.names = FALSE)  
write.csv(department_age, "department_staff_final", row.names = FALSE)
```

Note: These files are already in our data folder online, so if you don't have time to do it don't worry.

# Exporting outputs

Now `employees_by_department.csv` (probably in your `Documents` folder).



# Exporting outputs

## Some notes on file paths

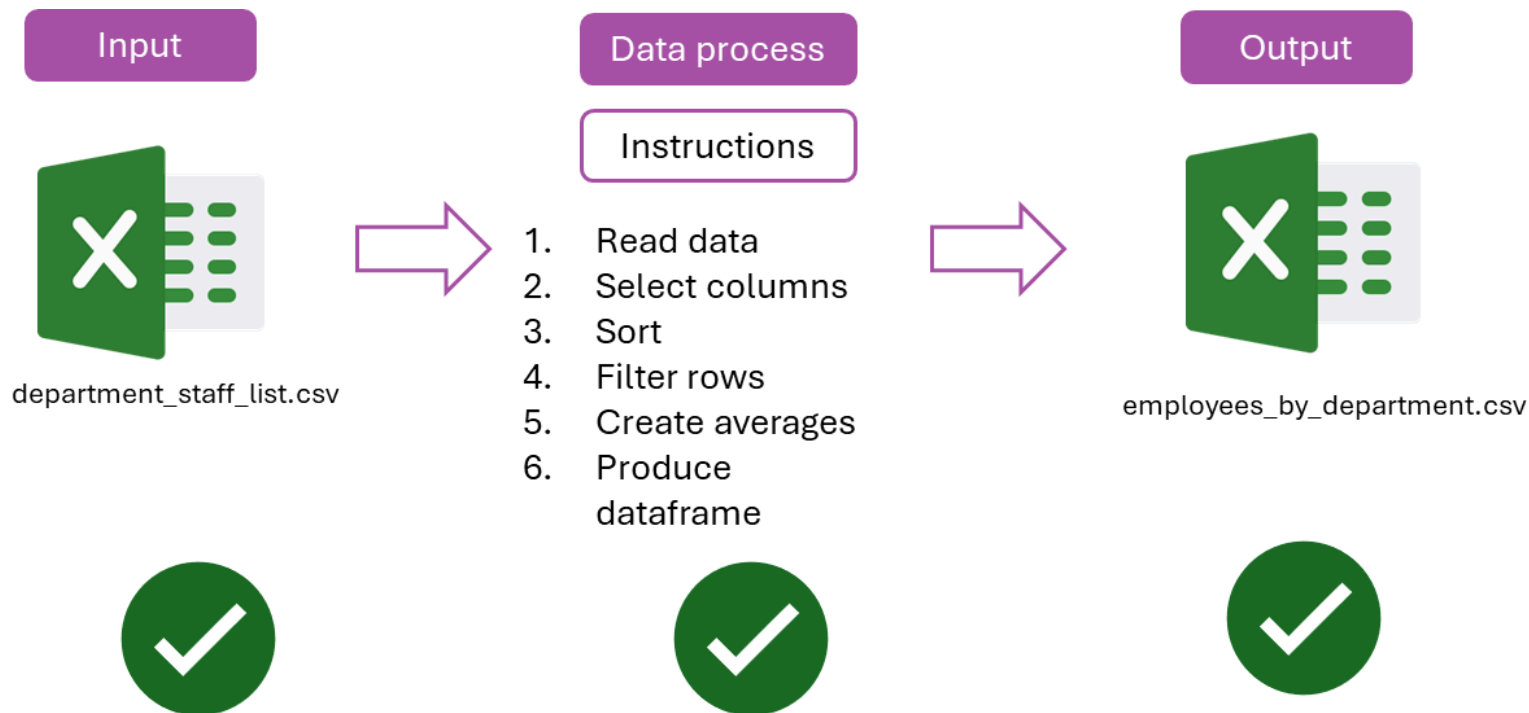
- The second argument of `write.csv()` specifies the file path we export the dataframe to

```
write.csv(employees_by_department, "employees_by_department", row.names = FALSE)
```

- You can include any path in your computer and R will write the file in that location
  - For example: `"C:/Users/wb614536/OneDrive - WBG/Desktop"` exports the file to the desktop of my computer (this will not work in other computers)
  - Note that file paths in R use forward slashes (`/`). Back slashes (`\`) **do not work in R**

# Exporting outputs

Our data pipeline has been fully implemented at this point. Great!



# Wrapping up

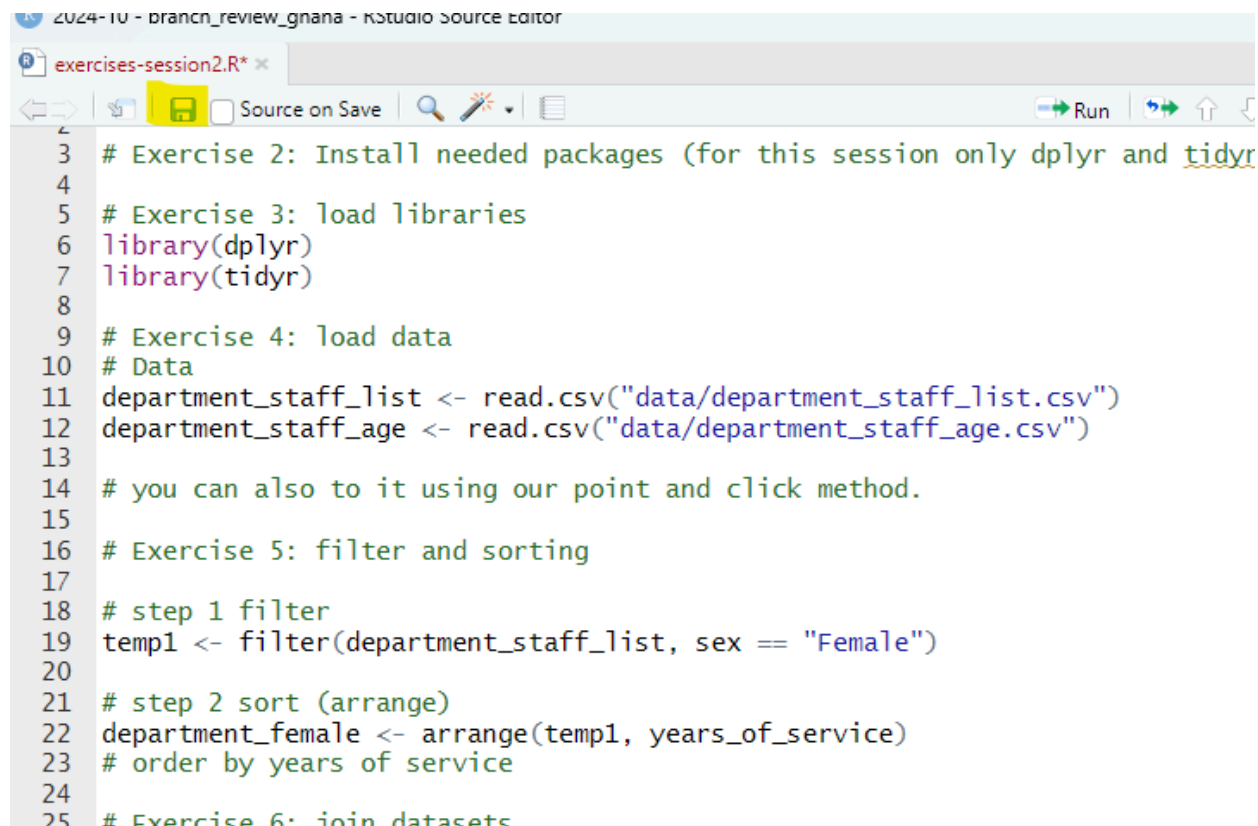
---



# Wrapping up

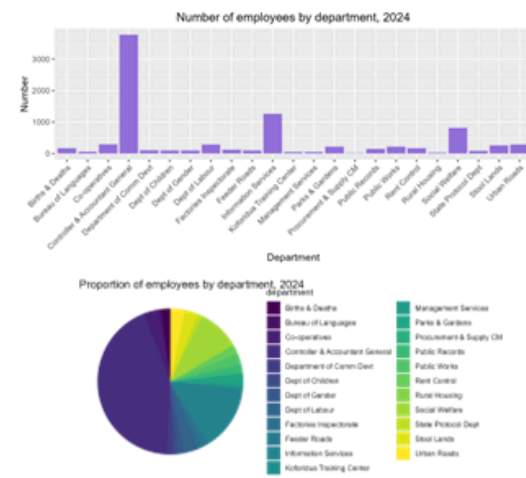
## Don't forget to save your work!

- If you haven't, add code comments with `#` to differentiate your solutions for each exercise
- Click the floppy disk to save your work
- Make sure to remember where you're saving your file



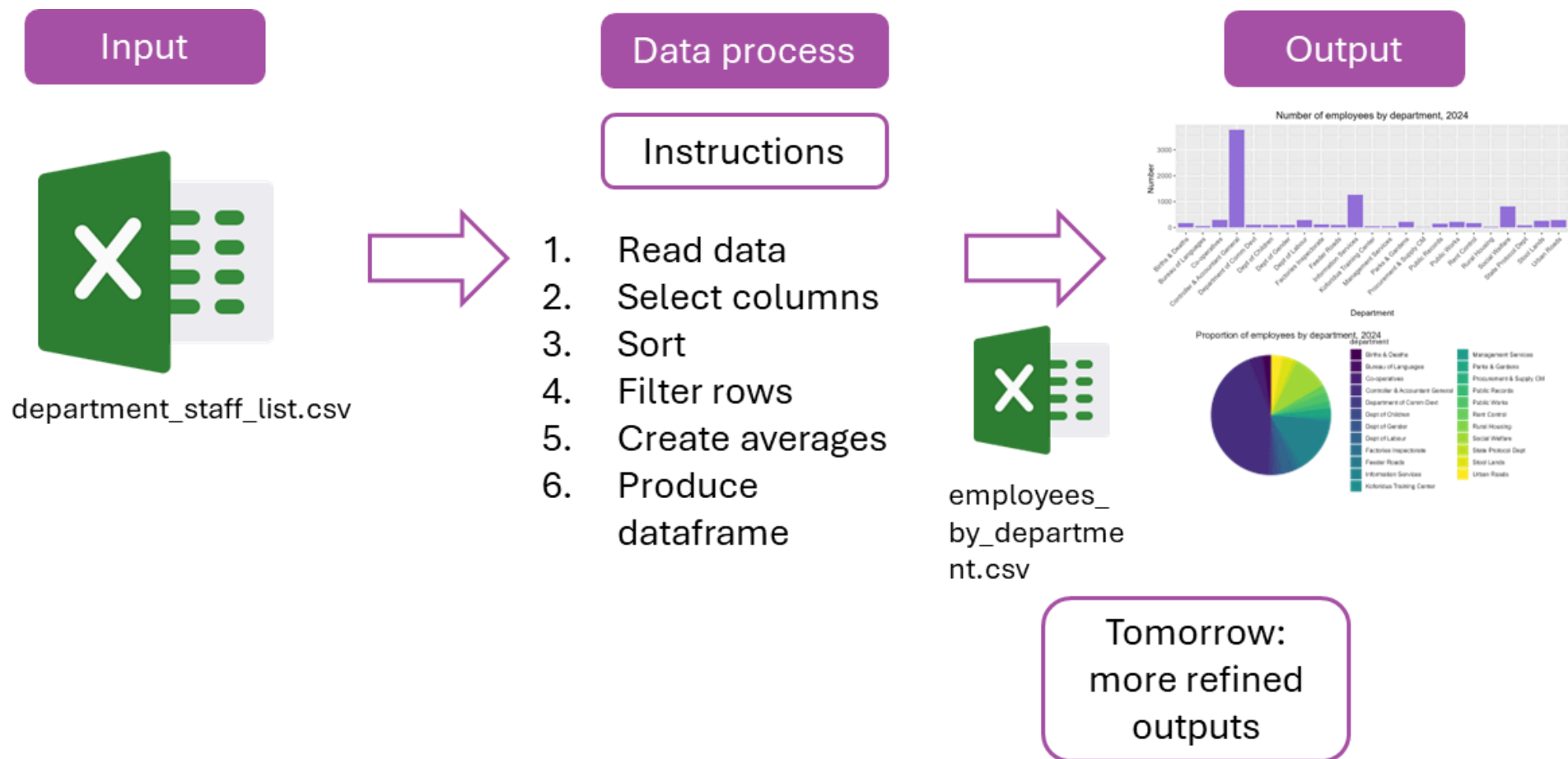
```
2024-10 - branch_review_gnana - RStudio Source Editor
exercises-session2.R*
Source on Save Run
3 # Exercise 2: Install needed packages (for this session only dplyr and tidyr)
4
5 # Exercise 3: load libraries
6 library(dplyr)
7 library(tidyr)
8
9 # Exercise 4: load data
10 # Data
11 department_staff_list <- read.csv("data/department_staff_list.csv")
12 department_staff_age <- read.csv("data/department_staff_age.csv")
13
14 # you can also to it using our point and click method.
15
16 # Exercise 5: filter and sorting
17
18 # step 1 filter
19 temp1 <- filter(department_staff_list, sex == "Female")
20
21 # step 2 sort (arrange)
22 department_female <- arrange(temp1, years_of_service)
23 # order by years of service
24
25 # Exercise 6: join datasets
```

# Government analytics pipeline



# Wrapping up

## Government analytics pipeline



Thanks! // ¡Gracias! // Obrigado!



# Appendix

---

# Appendix

## Keeping only relevant columns

Imaging I only want a list with ID and department.

Use `select()` for this:

```
temp1 <- select(department_staff_list, ID, department)
```

# Appendix

## Calculating aggregated columns: total and average income

Use `summarize()` combined with `sum()` and `mean()`:

For our current data this does not apply, but this will be really useful for budget dataframes.

This would look something like this

```
budget_2024 <- group_by(budget_data, year, product)

budget_2024 <- summarise(budget_2024,
                          total = sum(income),
                          average = mean(income))
```

# Useful links

- **R for Data Science** by Hadley Wickham and Garrett Golemund
  - Comprehensive introduction to data wrangling and visualization.
  - [Read it online for free.](#)
- **Tidyverse Cookbook**
  - Practical solutions for common data wrangling tasks.
  - [Tidyverse Cookbook GitHub.](#)
- **Tidyverse Cheat Sheets**
  - Official cheat sheets for `dplyr`, `tidyr`, and other Tidyverse packages.
  - [Tidyverse Resources.](#)